# Towards a Trustworthy and Resilient Machine Learning Classifier

## A Case Study of Ransomware Detector Creation

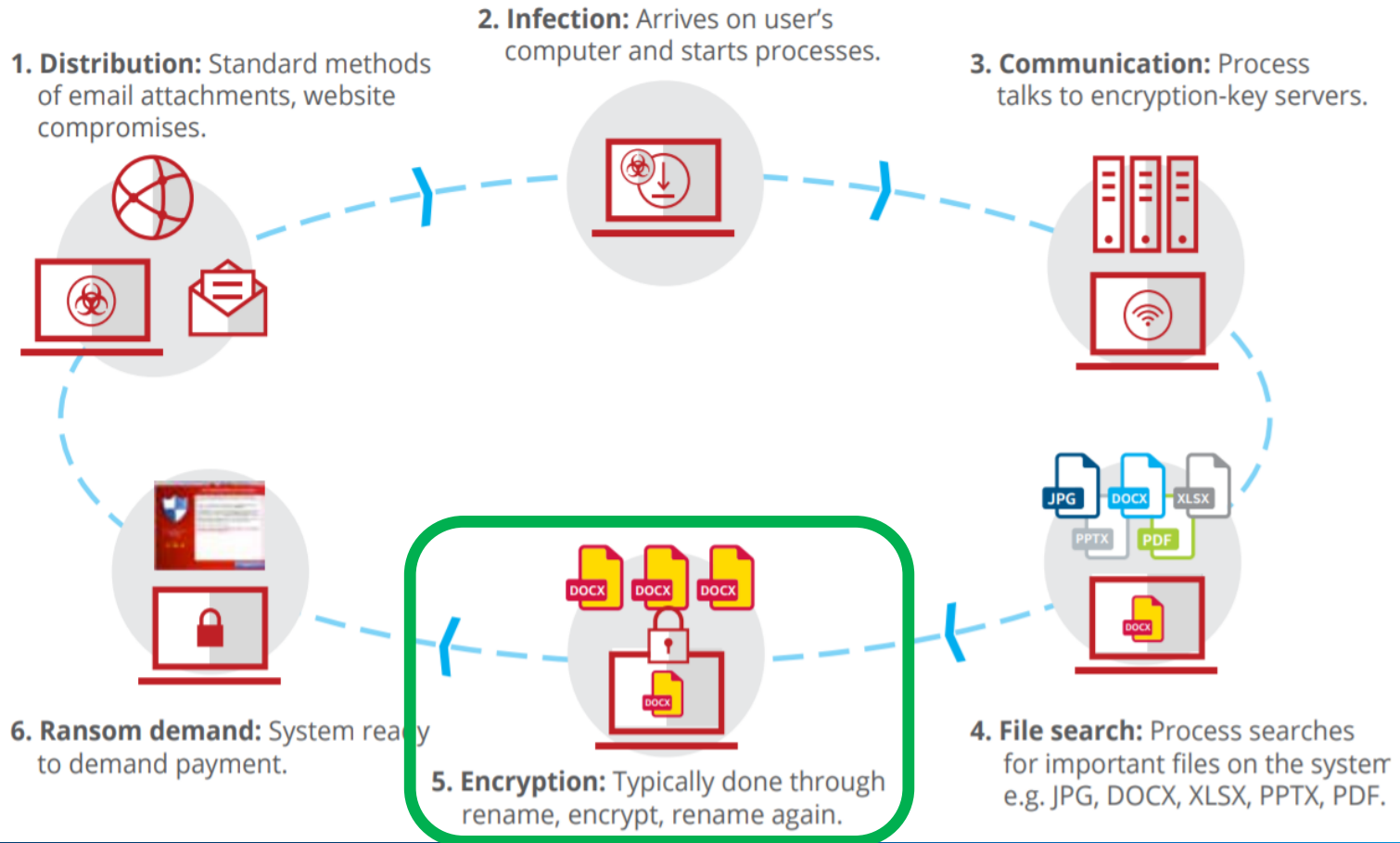**Evan Yang; Ravi Sahita**

**Intel Lab/SPR**

# Outline

- Background

- Issues of Classifier

- Model Fidelity

- Adversarial Research

- Conclusion

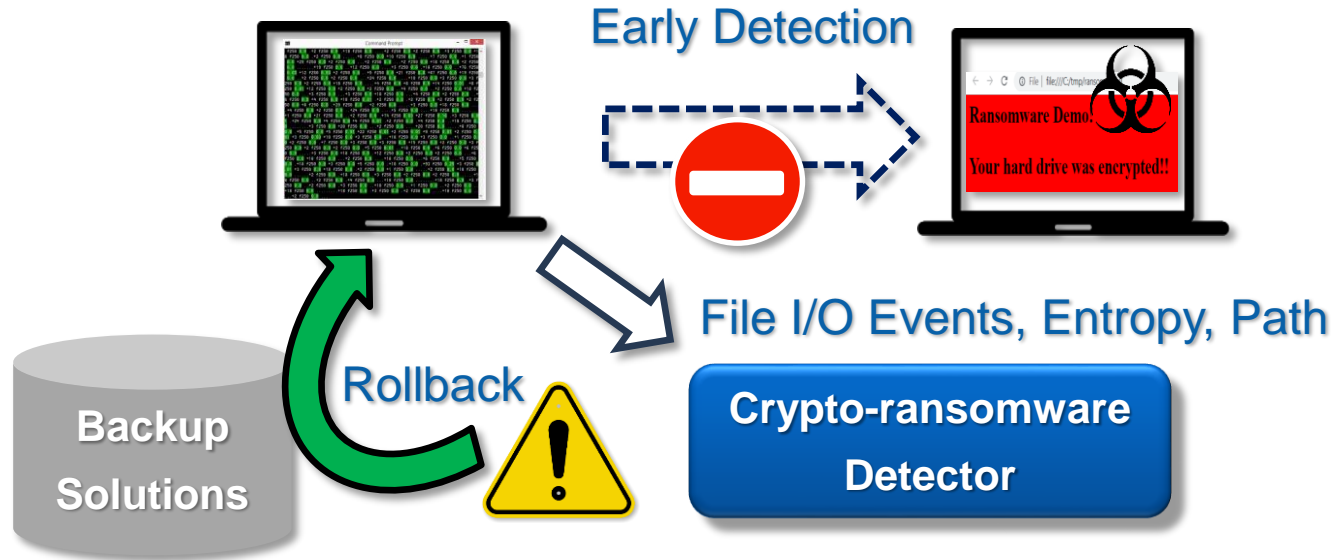# Background [(Al-rimy, B. et al. 2018)](#)

- Ransomware is a category of malware which hijacks victim's data or machine and demands monetary returns

- Taxonomy:

  – Locker-ransomware: hijack resources without encryption

  – Crypto-ransomware: encrypt files

- The damage done by crypto-ransomware is **irreversible** in most cases due to the use of cryptography
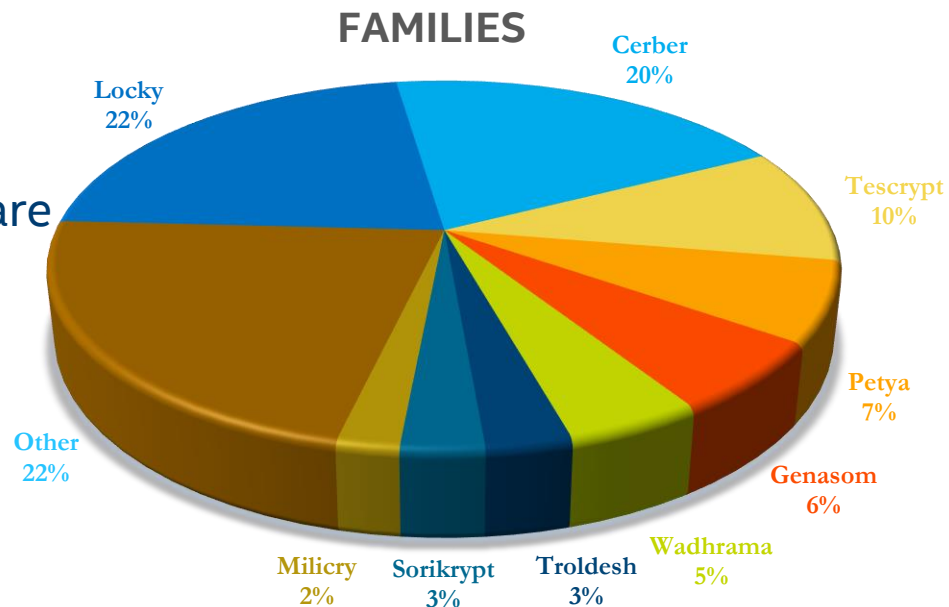
# Typical Steps of Ransomware (McAfee 2017)



**1. Distribution:** Standard methods of email attachments, website compromises.

**2. Infection:** Arrives on user's computer and starts processes.

**3. Communication:** Process talks to encryption-key servers.

**4. File search:** Process searches for important files on the system e.g. JPG, DOCX, XLSX, PPTX, PDF.

**5. Encryption:** Typically done through rename, encrypt, rename again.

**6. Ransom demand:** System ready to demand payment.

# Purpose of Detector:

- Find crypto-ransomware early by its behavior when AV missed it



Early Detection

File I/O Events, Entropy, Path

Rollback

**Backup Solutions**

**Crypto-ransomware Detector**

# Ransomware Dataset

- From VirusTotal

  - Downloaded total ~22k ransomware by Microsoft and Kaspersky's labels

  - ~5min execution for each sample

  - In bare-metal sandbox system with anti-evasion mechanism

- Decoy files to identify crypto-ransomware

- Total ~4.4k active samples:

**FAMILIES**



Pie chart — FAMILIES:
- Cerber 20%
- Tescrypt 10%
- Petya 7%
- Genasom 6%
- Wadhrama 5%
- Troldesh 3%
- Sorikrypt 3%
- Milicry 2%
- Other 22%
- Locky 22%

# Behavior Data – File Input/Output Events

- Collected by POC Windows application
  - Based on **C#.Net** framework, **FileSystemWatcher** (FSW)
  - Entropy of target files calculated by normalized Shannon entropy

- Sample data:
  - Time stamp, I/O event type, target filename, entropy etc.

"2018-04-06T12:21:28","27937 ,"Changed", c:\Windows\System32\wbem\Repository\MAPPING1.MA ",0.465655021998745, DAB00001AB000006F0200006E020000

"2018-04-06T12:21:29","28890 ,"Created", c:\temp\start_00b4d8bf603522c86b572819beac6d7c5 ded1800368071fe74ed3 280e2ca45_kasperskyransom_typepeex

"2018-04-06T12:21:29","28890 ,"Changed", c:\temp\start_00b4d8bf603522c86b572819beac6d7c5 ded1800368071fe74ed3 280e2ca45_kasperskyransom_typepeex

"2018-04-06T12:21:30","29890 ,"Changed", c:\Windows\System32\wbem\Repository\MAPPING2.MA ",0.465815580631633, DAB000056B80000700200006F020000

"2018-04-06T12:21:30","29937 ,"Changed", c:\Windows\System32\wbem\Repository\MAPPING3.MA ",0.466954218868868, DAB00006AB80000710200070020000

"2018-04-06T12:21:30","29968 ,"Changed", c:\Windows\System32\wbem\Repository\INDEX.BTR", .572493921393108,CCA 00006001000000000000000000000

"2018-04-06T12:21:30","29984 ,"Changed", c:\Windows\System32\wbem\Repository\MAPPING1.MA ",0.466994188018237, DAB00007AB80000720200071020000

# Machine Learning Analysis

- ~3.7k ransomware and similar amount of benign data (~100 applications). 80/20 split for training/testing dataset

- Featuring by event type with bucketed entropy (`-,0.2,0.4,0.6,0.8,0.9`)
  - Categorize into distinct features

- ML Algorithms for supervised learning:
  - Long-Short Term Memory (LSTM), Recurrent Neural Networks
  - Linear Support Vector Machine (SVM) with bag of N-gram, N=1 & 2
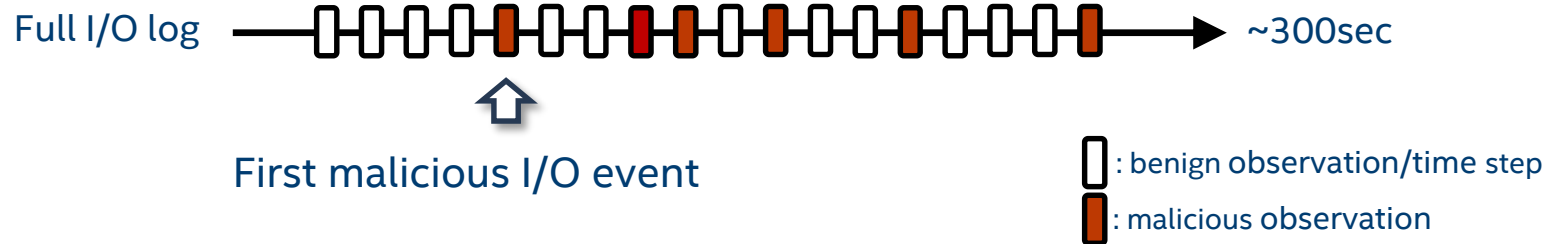
# ML Pipeline & Outcome of Supervised Learning

**ML/DL Iteration**

I/O Events → Data Processing → ML/DL Training → ML/DL Classifier → Validation

| Model | N–gram | Accuracy | FPR | Dist. Features |
|---|---|---|---|---|
| Linear SVM | 1 & 2 | 98.31% | 2.89% | 90 |
| LSTM | n/a | 98.67% | 1.38% | 9 |

# Online Detector

- A POC program utilized the ML classifier
  - Sample the I/O event stream by a sliding window
  - Real-time inference: small footprint and run fast
- Issues found after deployment:
  1. False alarms from some applications
  2. Size of sliding window affects the detection rate
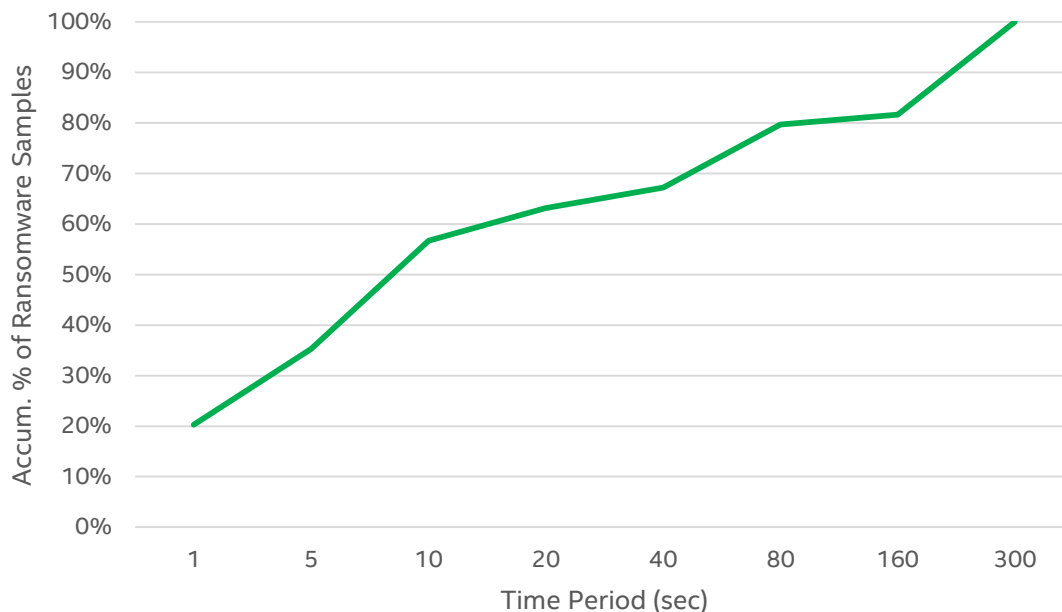  3. Cannot find ransomware early

# 1. Early Detection Issue

- Early detection is important
  - No practical value if can't detect encryption early

- When will the ransomware start doing encryption?
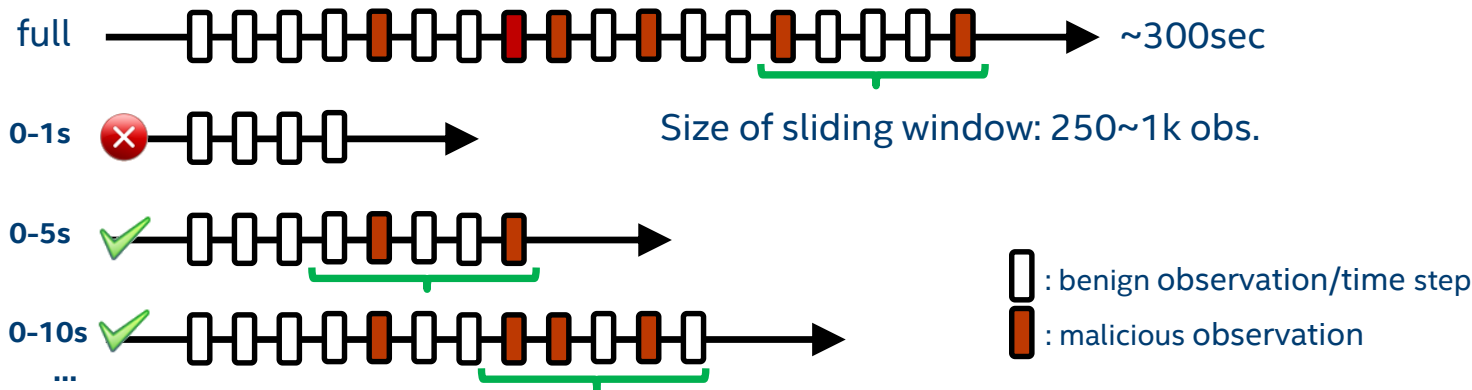  - Identify the starting time by the decoy file

Full I/O log      ~300sec

First malicious I/O event

☐ : benign observation/time step

🟧 : malicious observation

# Starting Time of Malicious Activities

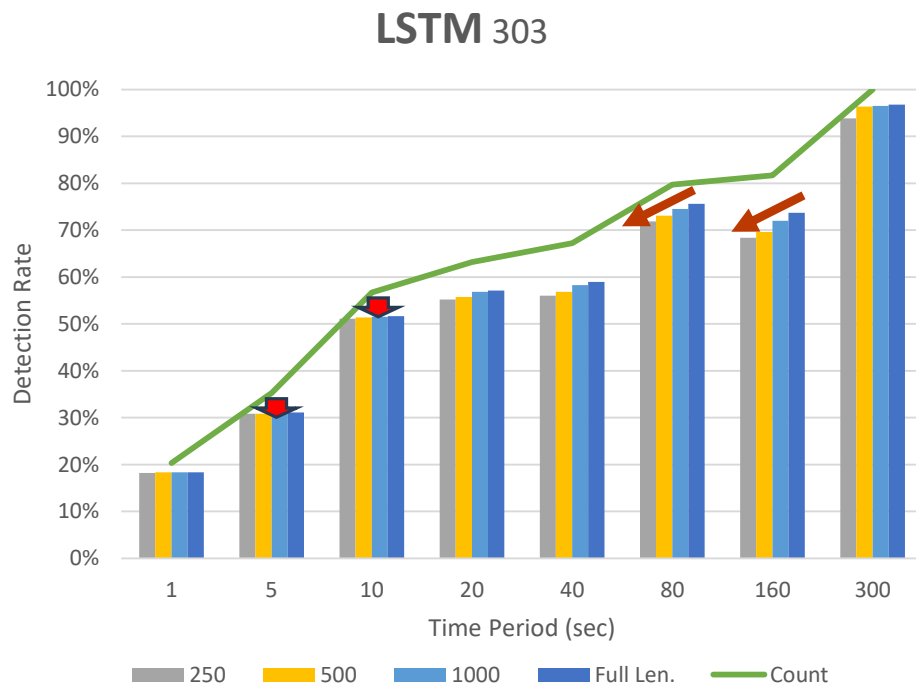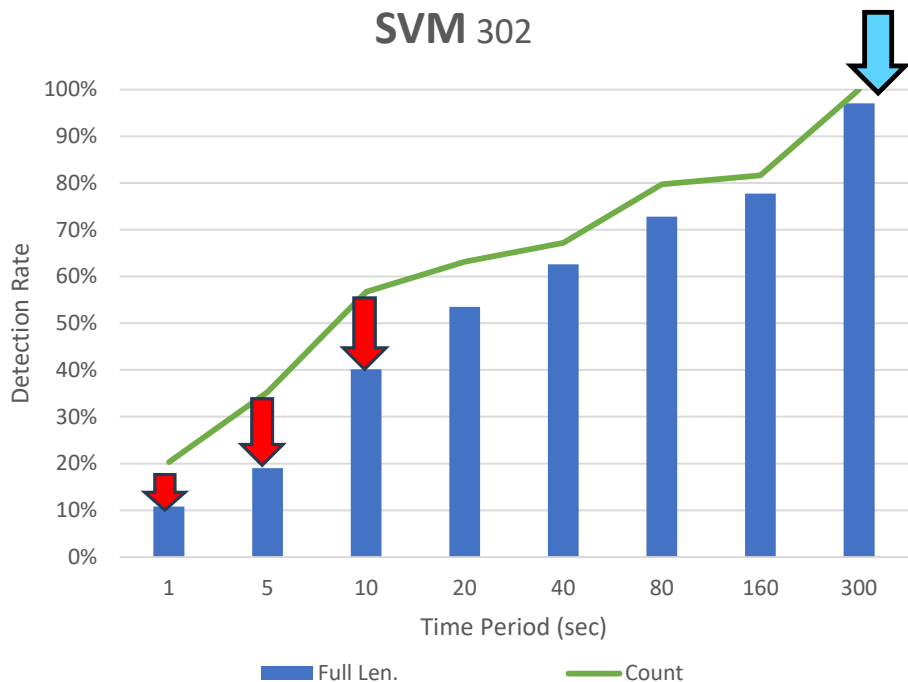- Ransomware may not show malicious activities at the beginning of execution

# Early Detection and Sliding Window Testing

- Prepare samples to measure the performance
  - From ~700 unseen *out-of-sample* ransomware logs

- Extract early-stage data from each logs by
  - different time periods
  - different sliding windows

- Example:



full    ~300sec

Size of sliding window: 250~1k obs.

0-1s

0-5s

0-10s
...

▯ : benign observation/time step

▮ : malicious observation

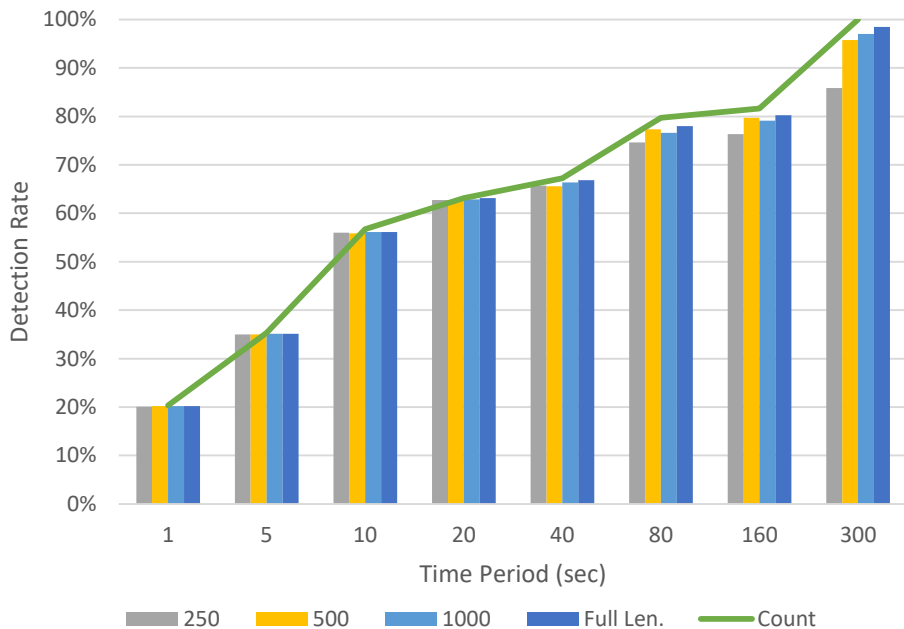# Detection Rate of Early-stage and Sliding-window

# Data Augmentation

- Synthesize samples from existing dataset for a re-train
  - Early-stage samples
  - Sliding-window samples
  - Exclude samples without malicious events
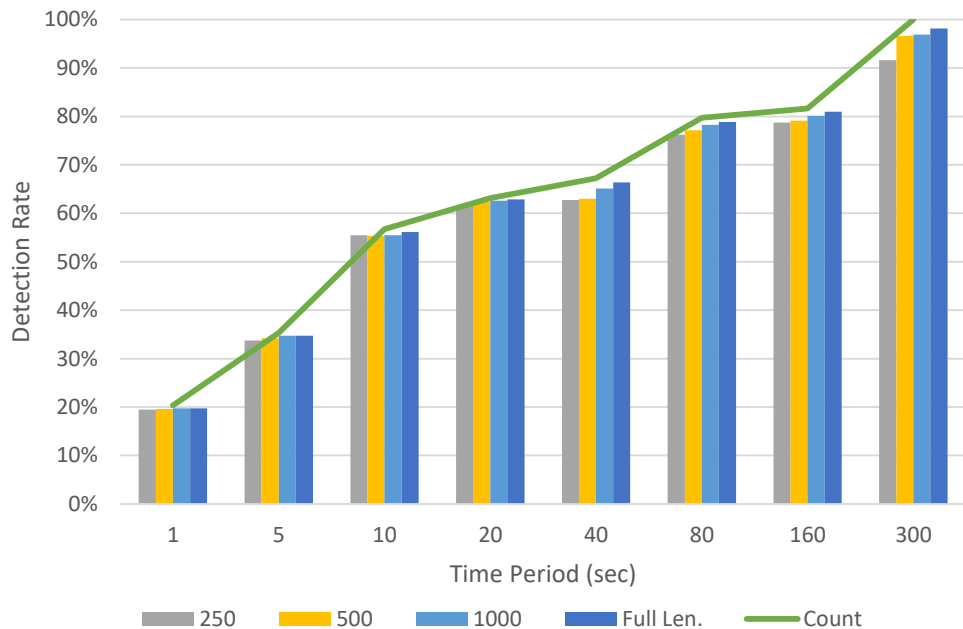- "Augmented" dataset count: 17.2k ransomware (80/20 split)

| Model | N–gram | Accuracy | FPR | Dist. Feature |
|---|---|---|---|---|
| Linear SVM | 1 & 2 | 99.13% | 1.21% | 90 |
| LSTM | n/a | 99.47% | 0.60% | 9 |

# Detection Rate by Augmented Classifier



SVM-A 319

LSTM-A 320

# 3. False Positive Issue

- Some benign-ware has similar ransomware behaviors
  - Delete or rename many files, change files with high entropy

- **Solution**: Add a new dimension to feature
  - Path: system vs. non-system folders
  - System path list: `c:\Windows, c:\ProgramData, c:\Program Files, c:\Progra~, c:\AppData, \Downloads\, \Downlo~, c:\Config.msi`

# Results with Path Flag

- Lower FPR with flag added

| Model | N–gram | Accuracy | FPR | Dist. Features |
|---|---|---|---|---|
| Linear SVM | 1 & 2 | 99.00% | 1.34% | 90 |
| Linear SVM (+ path) | 1 & 2 | 99.53% | 0.54% | 339 |
| LSTM | – | 98.26% | 3.82% | 9 |
| LSTM (+ path) | – | 98.35% | 1.80% | 18 |

- 22k out-of-sample clean execution log:
  - FPR down from 0.18% to 0.00% for SVM (40->0/22,174)
  - FPR down from 0.09% to 0.04% for LSTM (21->9/22,174)

# Model Fidelity by Integrated Gradients

- Attribution: which feature/time step contribute the most?

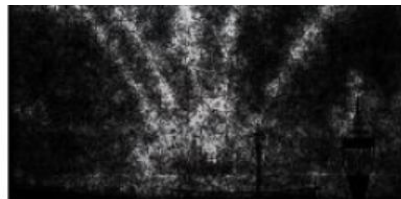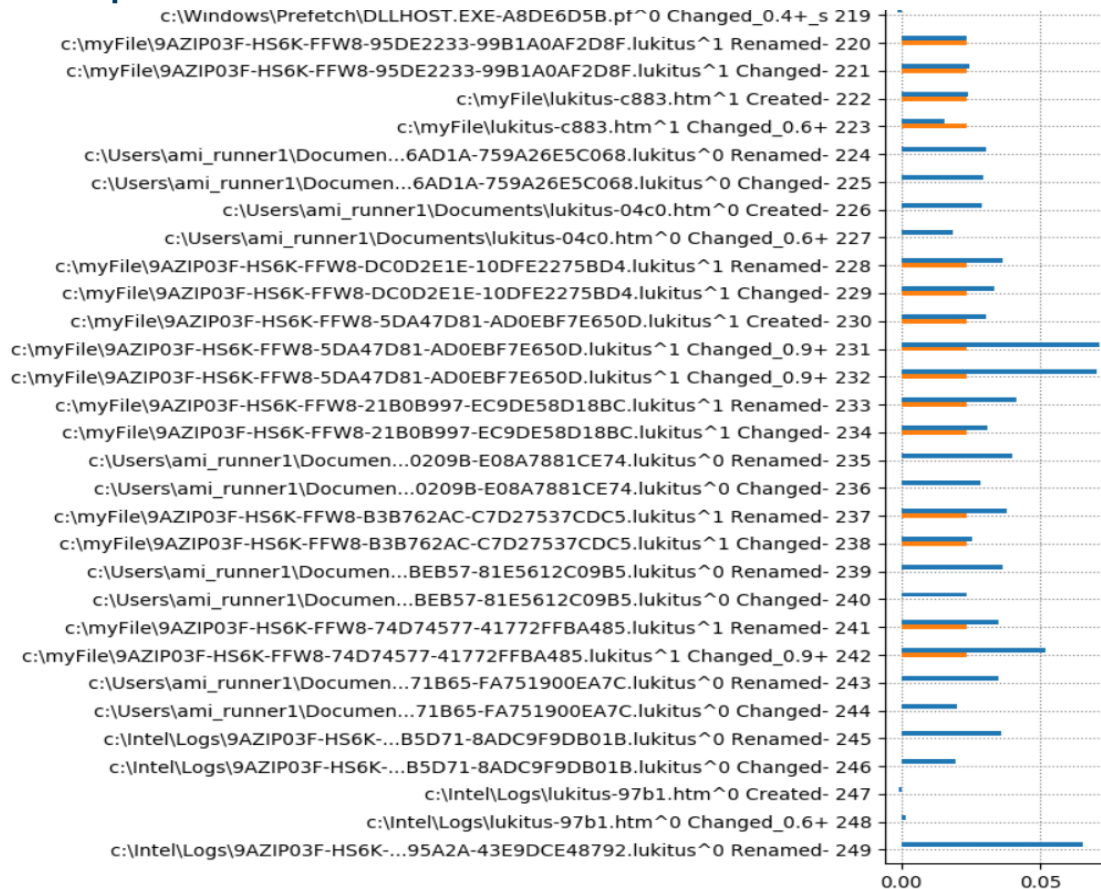| Original image | Top label and score | Integrated gradients | Gradients at image |
|---|---|---|---|
| | Top label: reflex camera | | |

how many townships have a population above 50 ? [prediction: NUMERIC]
what is the difference in population between fora and masilo [prediction: NUMERIC]
how many athletes are not ranked ? [prediction: NUMERIC]
what is the total number of points scored ? [prediction: NUMERIC]
which film was before the audacity of democracy ? [prediction: STRING]
which year did she work on the most films ? [prediction: DATETIME]
what year was the last school established ? [prediction: DATETIME]
when did ed sheeran get his first number one of the year ? [prediction: DATETIME]
did charles oakley play more minutes than robert parish ? [prediction: YESNO]

# Explanation of LSTM Models

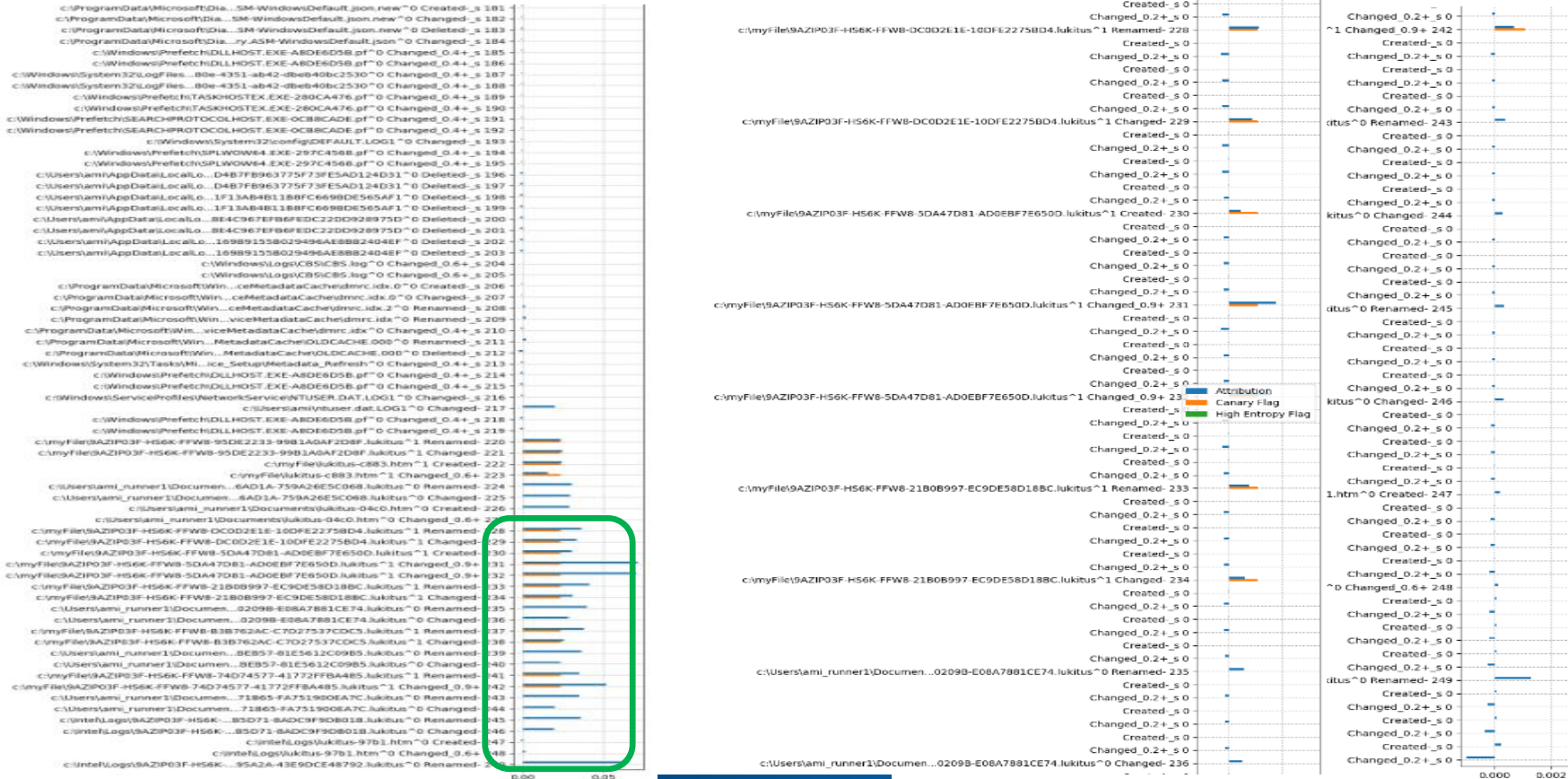- Feature attribution plot of ransomware:

# Adversarial Research

- A *simulated* ransomware, the *Red* team, was developed in C#

  – Rename, encrypt and delete files etc.

  – Evasive tricks to probe the detector (grey box attack):
    - Behavior temporal changes: e.g. slowdown the malicious activities
    - Encryption changes: e.g. insert dummy data to lower the file entropy

  – It's not difficult to evade our ML detector

- Improve model's resiliency by:

  – Discover weakness by the Red team with various conditions
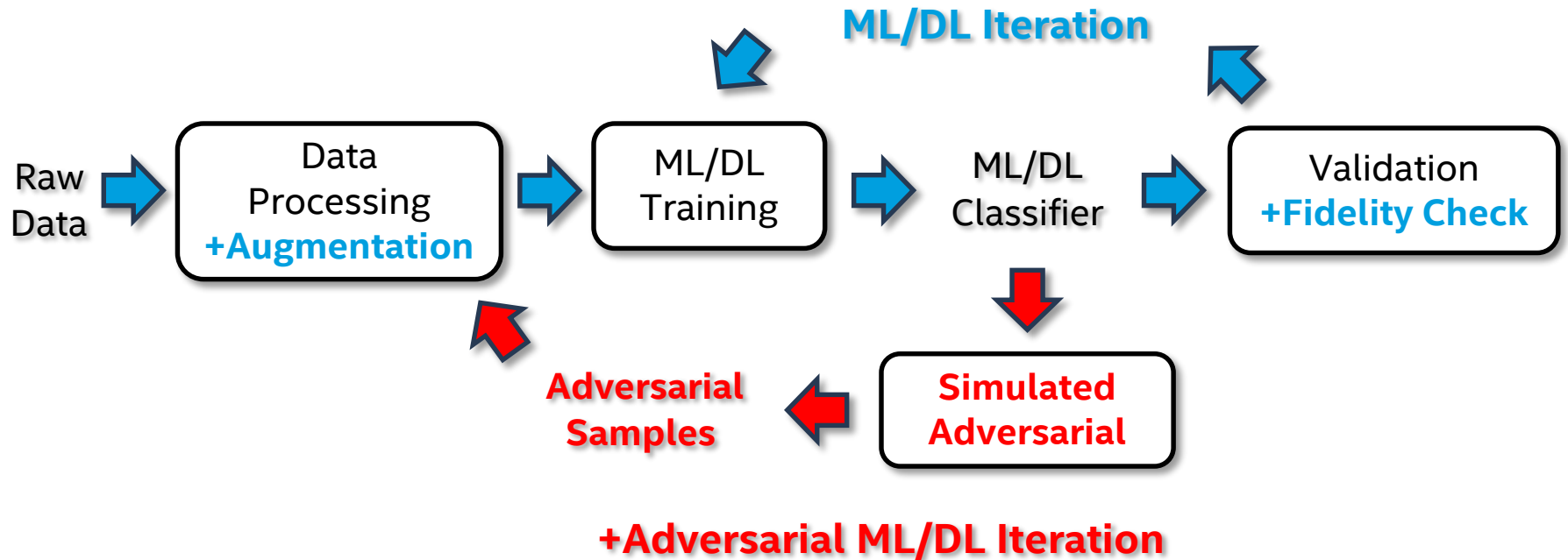
  – Re-train model by the false negatives samples

# Probing LSTM Models – by Event Insertion

Original Sample, +, 0.96

Insert 7 benign events, –, 0.01

# Conclusion: ML Pipeline +++



Raw Data → Data Processing **+Augmentation** → ML/DL Training → ML/DL Classifier → Validation **+Fidelity Check**

**ML/DL Iteration**

**Simulated Adversarial** → **Adversarial Samples**

**+Adversarial ML/DL Iteration**

# Our Team Members and Projects

- Erdem Aktas; Li Chen; Anindya Paul

- MLsploit: a platform for ML model comparison and sample sharing for adversarial research

  - github.com/mlsploit

  - github.com/intel/Resilient-ML-Research-Platform

# Thank You !

# Legal Notices and Disclaimers

- No license (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document. These materials are provided as-is, with no express or implied warranties. All products, computer systems, dates and figures specified are preliminary based on current expectations, and are subject to change without notice. The products described may contain design defects or errors known as errata which may cause the product to deviate from published specifications. Intel does not control or audit third-party data or the web sites referenced in this document. You should visit the referenced web site and confirm whether referenced data are accurate.

- © 2019 Intel Corporation.  Intel and the Intel logo are trademarks of Intel Corporation in the U.S. and/or other countries.

- *Other names and brands may be claimed as the property of others.