



Activation Analysis of a Byte-based Deep Neural Network for Malware Classification

Scott Coull

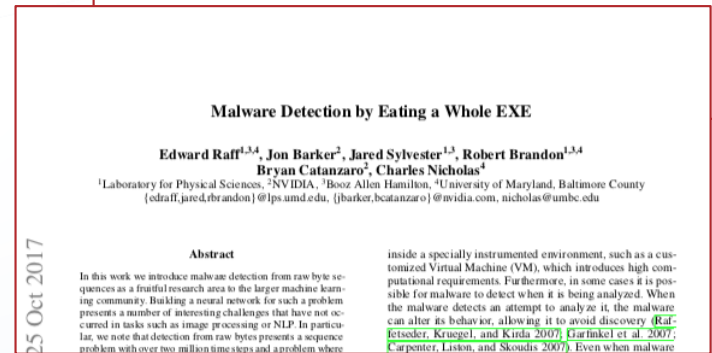
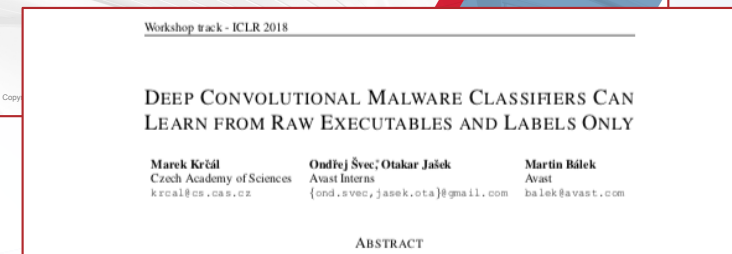
Sr. Manager, Data Science

Christopher Gardner

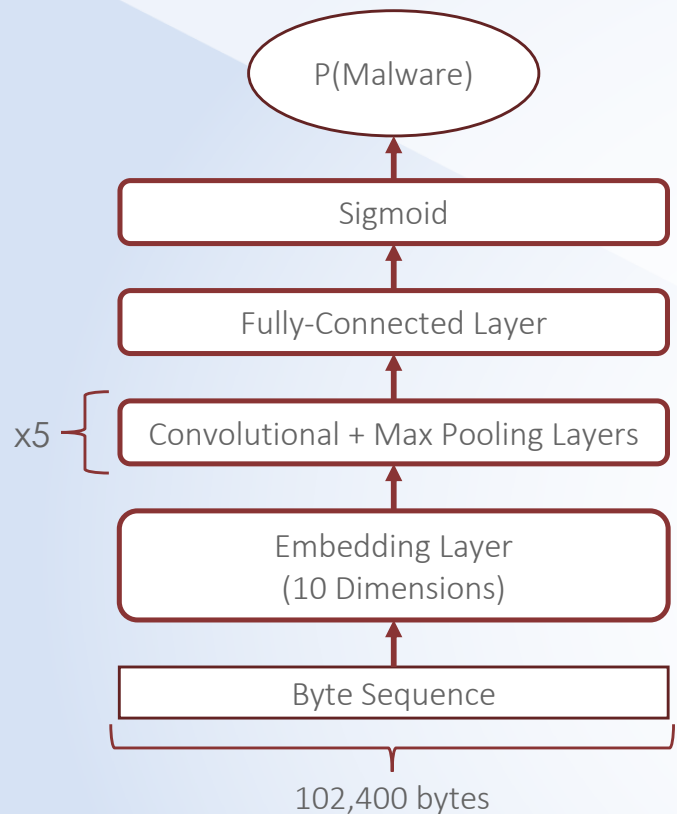
Reverse Engineer

Byte-based Malware Classifiers

- Feature engineering for malware classification tasks is **hard**. Can deep learning do it for us?
- Convolutional neural networks (CNNs) **automatically and efficiently learn feature representations** directly from data
- Recent work has shown **promising results** competitive with (though not better than) traditional machine learning
 - Accuracy: 90-96%, AUC: 0.96-0.98



CNN Models



■ Baseline

- 15.6M Windows PEs (80% goodware)
- July 2015 to July 2017
- Stratified sampling

■ Small

- 7.3M Windows PEs (50% goodware)
- July 2016 to November 2016
- No sampling

■ Baseline+Dropout

- Same data as Baseline
- Dropout layers before convolutional layers

Model Evaluation

Model	Train Data		Test Results	
	Size	Mal:Good	F1	AUC
Small	7.27M	50:50	0.943	0.98
Baseline	15.62M	20:80	0.919	0.96
Baseline+Dropout	15.62M	20:80	0.869	0.87

16.55M binaries (50:50) from June 1, 2018 to August 31, 2018

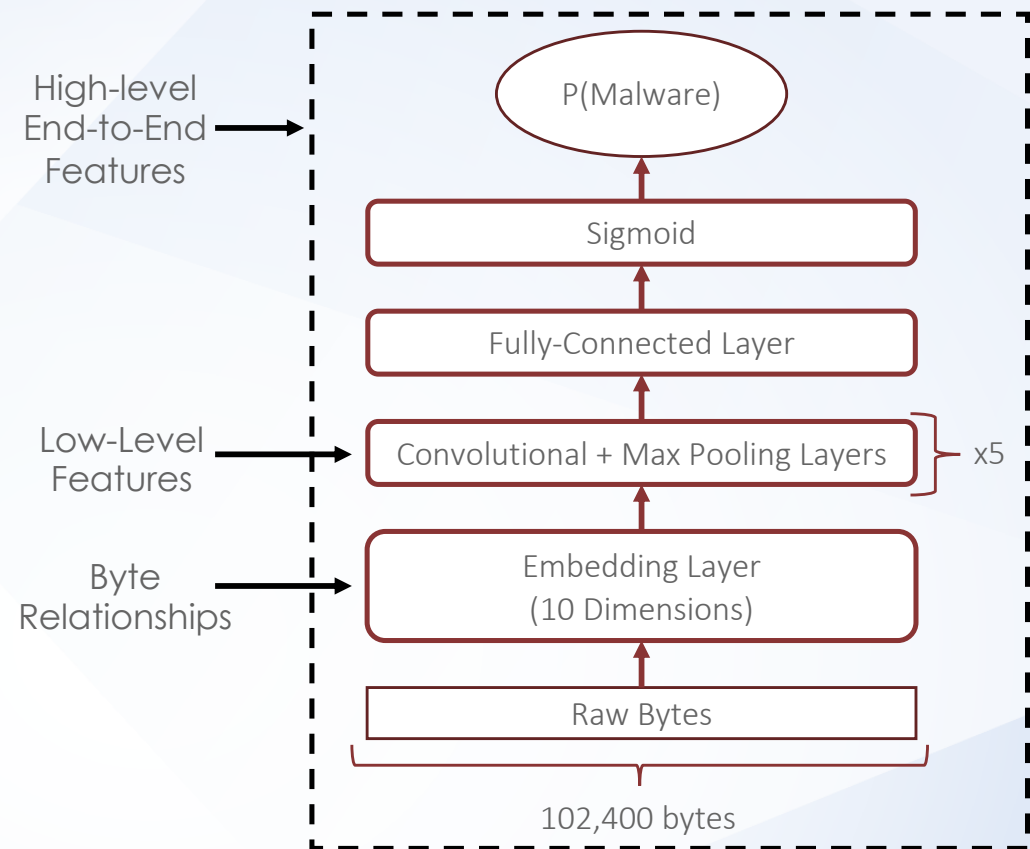
Model trained on **small dataset performs noticeably better**
despite older data and fewer samples

What are byte-based malware classifiers learning?

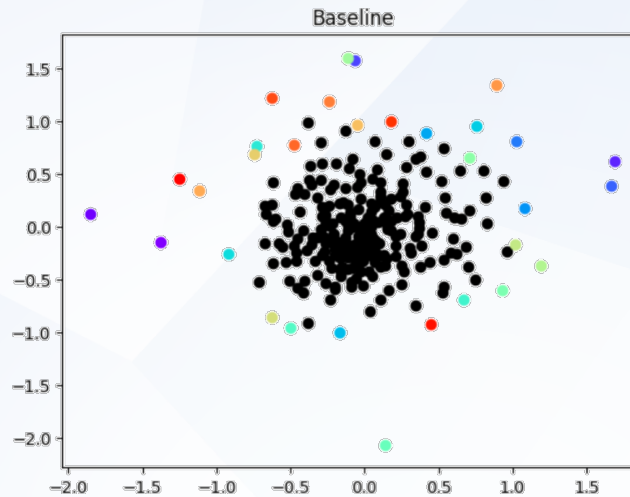
What is the impact of dataset volume and regularization on learned features?

Analysis Overview

- **Cluster** and **visualize** embedding layer with HDBSCAN¹ and MDS²
- **Disassembly of byte sequences** with large activations in first convolutional layer
- **End-to-end analysis** of byte segments with GradientSHAP³



Byte Embeddings

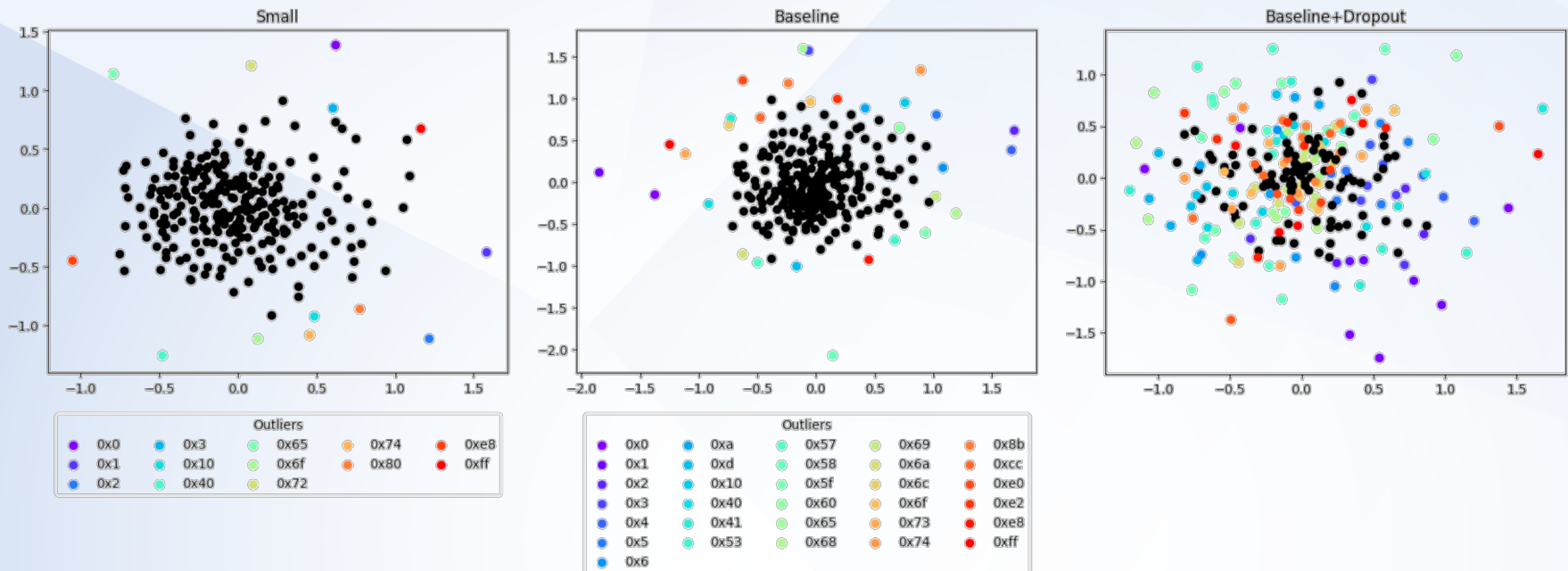


Outliers					
0x0	0xa	0x57	0xb9	0x8b	
0x1	0xd	0x58	0x6a	0xcc	
0x2	0x10	0x5f	0x6c	0xe0	
0x3	0x40	0x60	0x6f	0xe2	
0x4	0x41	0x65	0x73	0xe8	
0x5	0x53	0x68	0x74	0xff	
0x6					

eax-edx
padding
short jumps

ASCII Characters
@, \n, A, ..., j, s, t

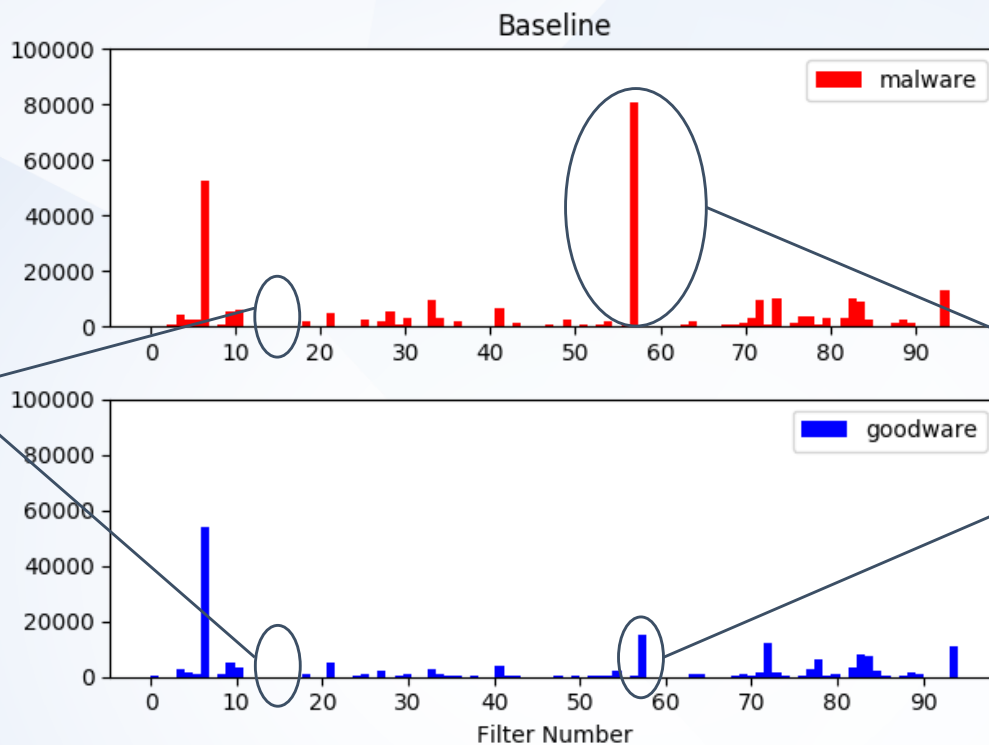
Byte Embeddings



→
Increase in number of outliers with more data/regularization
Learned features appear to be less flexible

Low-Level Feature Detectors

Distribution of Top-100 Activations Across First-Level Filters

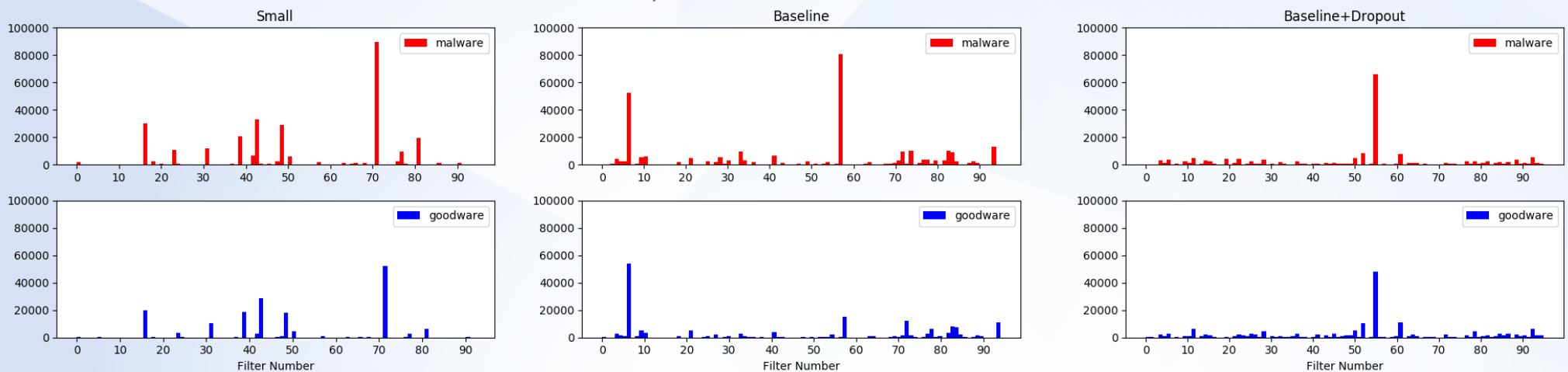


Some unused filters

Same filters,
bias toward malware

Low-Level Feature Detectors

Distribution of Top-100 Activations Across First-Level Filters



More data and regularization appears to lead to more features that are equally applicable across the two classes
Supports earlier observation about feature specificity

Low-Level Feature Detectors

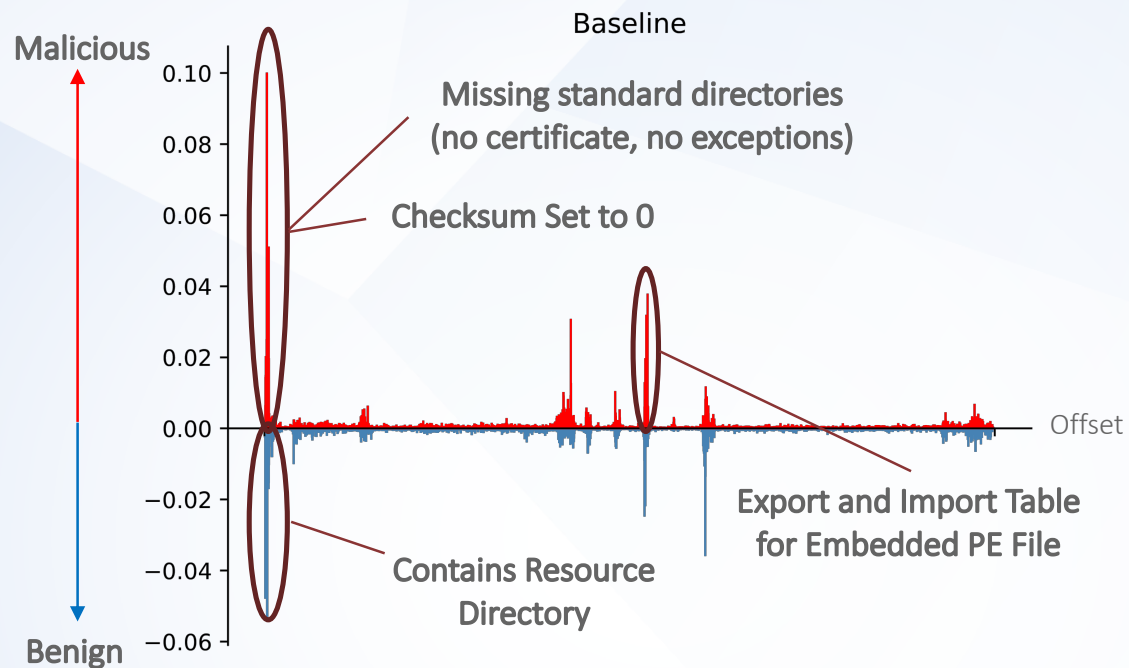
Loose filters

Specific filters

Model	Features	
	Strings	Instructions
Small	Filter 71: 'C', 'r', '@' (0x40f0c8L): tGenKey. (0x40f0d0L): CryptDec (0x40f0d8L): rypt.... (0x40f0e0L): CryptEnc (0x40f0e8L): rypt....	Filter 16: Push sequences (0x10007edbL): je,0x10007ff1 (0x10007ee1L): push,0xff (0x10007ee6L): push,edi (0x10007ee7L): push,0x10007ca5 (0x10007eecL): push,0x4
Baseline	Filter 83: 'r', 's' (0x40d850L):GetP (0x40d858L): rocAddre (0x40d860L): ss..R.Lo (0x40d868L): adLibrar (0x40d870L): yA....Gl	Filter 57: Function calls (0x4046b4L): push,0x0 (0x4046b6L): push,0x0 (0x4046b8L): push,0x1 (0x4046baL): push,0x0 (0x4046bcL): call,dword,15042
Baseline+Dropout	Filter 11: 'Directory' (0x40d9e0L): ctoryW.. (0x40d9e8L): N.Create (0x40d9f0L): ,Director (0x40d9f8L): yW....Ge (0x40da00L): tTempPat	Filter 61: mov sequences (0x408d65L): je, 0x408d6a (0x408d67L): mov, dword , edx (0x408d6aL): mov, esi, dword (0x408d6dL): mov, dword, esi (0x408d70L): mov, ecx, dword

End-to-End Features

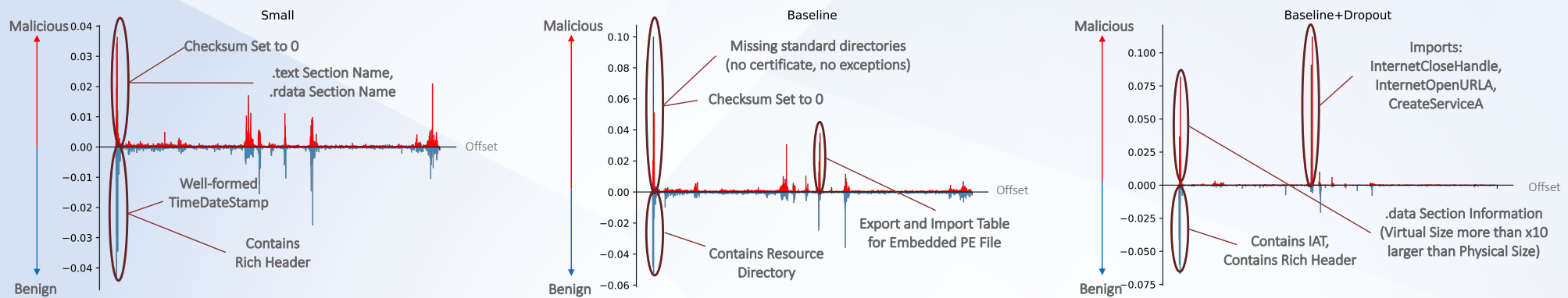
SHAP Values for WannaCry Worm



End-to-end features map closely to manual feature engineering

End-to-End Features

SHAP Values for WannaCry Worm



Data and regularization result in more focused areas of interest
Model appears to learn presence/absence of structural features

The Case of the Rich Header

- Rich header is added by Microsoft's linker and contains metadata about the binary
- Should be effectively 'random' due to XOR encryption using key derived from checksum
- Hypothesis: Hierarchical pooling can detect presence of fixed bytes around header (e.g., 'Rich')
- Proxy for whether **non-Microsoft compiler** was used, which is **common in malware**

```
HxD - [C:\Windows\System32\calc.exe]
File Edit Search View Analysis Extras Window ?
16 ANSI hex
calc.exe
Offset(h) 00 01 02 03 04 05 06 07 08 09 0A 0B 0C 0D 0E 0F
00000000 4D 5A 90 00 03 00 00 00 04 00 00 00 FF FF 00 00 MZ.....ÿÿ..
00000010 B8 00 00 00 00 00 00 00 40 00 00 00 00 00 00 00 .....@.....
00000020 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
00000030 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
00000040 0E 1F BA 0E 00 B4 09 CD 21 B8 01 4C CD 21 54 68 ..*...Ë!..LITh
00000050 69 73 20 70 72 6F 67 72 61 6D 20 63 61 6E 6E 6F is program canno
00000060 74 20 62 65 20 72 75 6E 20 69 6E 20 44 4F 53 20 t be run in DOS
00000070 6D 6F 64 65 2E 0D 0D 0A 24 00 00 00 00 00 00 00 mode....$.
00000080 08 73 A6 53 4C 12 C8 00 4C 12 C8 00 4C 12 C8 00 .*!SL.E.L.E.L.E.
00000090 45 6A 5D 00 45 12 C8 00 4C 12 C9 00 D8 13 C8 00 Ej].E.E.L.E.Ø.E.
000000A0 45 6A 5B 00 45 12 C8 00 45 6A 4B 00 57 12 C8 00 Ej].m.E.EjK.W.E.
000000B0 45 6A 4C 00 4E 12 C8 00 45 6A 5C 00 4D 12 C8 00 Ej].i.E.Ej\M.E.
000000C0 45 6A 59 00 4D 12 C8 00 52 69 63 68 4C 12 C8 00 Ej].M.E.Ej\RichL.E.
000000D0 00 00 00 00 00 00 00 00 50 45 00 00 4C 01 04 00 .....FE..L..
000000E0 9D 97 E7 4C 00 00 00 00 00 00 00 00 E0 00 02 01 -SL.....à...
000000F0 0B 01 09 00 00 2E 05 00 00 A6 06 00 00 00 00 00 00 1-.....
00000100 6C 2D 01 00 00 10 00 00 00 20 05 00 00 00 00 01 1-.....
00000110 00 10 00 00 00 02 00 00 06 00 01 00 06 00 01 00 .....
00000120 06 00 01 00 00 00 00 00 00 00 0C 00 00 04 00 00 .....
00000130 30 BD 0C 00 02 00 40 81 00 00 04 00 00 20 00 00 0%.....@.....
00000140 00 00 10 00 00 10 00 00 00 00 00 00 10 00 00 00 .....
00000150 00 00 00 00 00 00 00 00 FC 1A 05 00 54 01 00 00 .....ù...T...
00000160 00 90 05 00 98 27 06 00 00 00 00 00 00 00 00 00 .....
00000170 00 00 00 00 00 00 00 00 00 00 C0 0B 00 3C 3B 00 00 .....Å.<?..
00000180 44 3C 05 00 38 00 00 00 00 00 00 00 00 00 00 00 D<..8.....
00000190 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
000001A0 30 04 03 00 40 00 00 00 70 02 00 00 54 01 00 00 0...@...p...T...
000001B0 00 10 00 00 30 06 00 00 78 1A 05 00 40 00 00 00 ...0...x...@...
000001C0 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
Offset: 10 Overwrite
```

Summary

- CNN architectures can learn **meaningful features**
 - Imports, presence of Rich header, incorrect checksums, etc.
 - Many features mimic manually-derived features from traditional ML models
 - Partly contradicts findings by Demetrio et al. on MalConv⁴
- **Model depth, dataset, and hierarchical pooling** appear to be key
- Malware classification performance relies on **detecting malware indicators**
 - Increased data and regularization lead to more specific features that were equally applicable across the two classes but worse detection performance

FireEye Data Science is Hiring!

- Data scientist positions open at the **Senior, Staff, and Principal level**
- Perform cutting-edge ML research and apply it to cybersecurity problems
- Work on problems from across the entire cybersecurity spectrum!
 - Threat Intelligence, Email, Network, Endpoint ...

Thank you!

 scott.coull@fireeye.com

 [@DrScottCoull](https://twitter.com/DrScottCoull)

 <https://scottcoull.com>