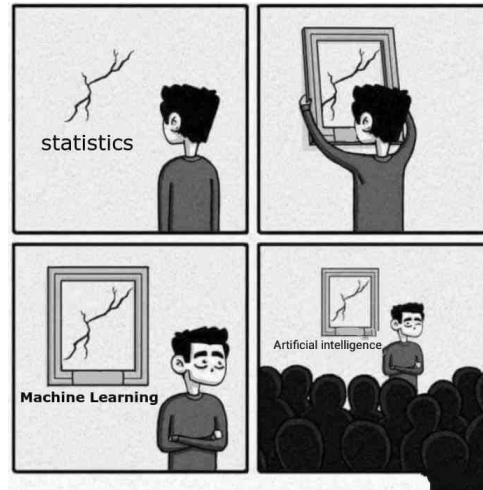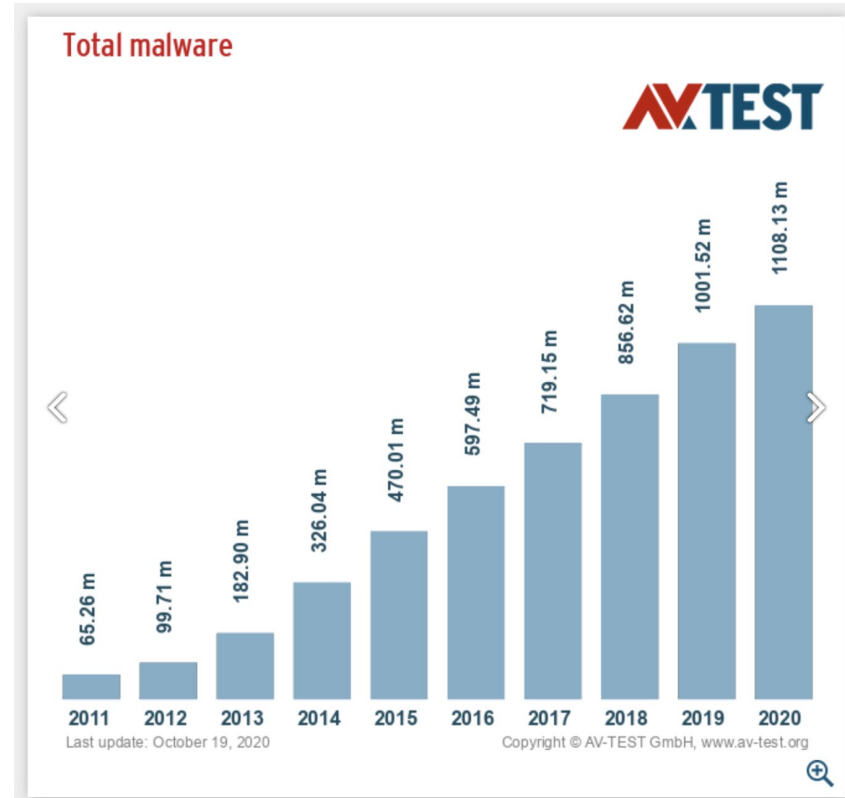# Cybersecurity and Machine Learning

# Motivation: **A lot of data!**

- Windows Executables
- Android Applications
- E-mails
- Network Traffic
- Authentication events
- Operation System data such as:
  - System calls
  - Process events
- and more...



Total malware

AV TEST

1108.13 m
1001.52 m
856.62 m
719.15 m
597.49 m
470.01 m
326.04 m
182.90 m
99.71 m
65.26 m

2011 2012 2013 2014 2015 2016 2017 2018 2019 2020

Last update: October 19, 2020

Copyright © AV-TEST GmbH, www.av-test.org

# Motivation: **A lot of data!**

- Windows Executables
- Android Applications
- E-mails
- Network Traffic
- Authentication events
- Operation System data such as:
  - System calls
  - Process events
- and more...

## 1. There are 3.9 billion active email users. (Radicati)

More than half of the global population now uses email. Radicati released updated figures early in 2019 that shows the total number of active email users has jumped to 3.9 billion. This represents accounts that have been assessed over the past three months, so there are likely many more accounts that exist but aren't frequented.

Just as a comparison, there are 3.5 billion social media users worldwide. The number of social users is impressive, but it's still fewer users than the number of email accounts.
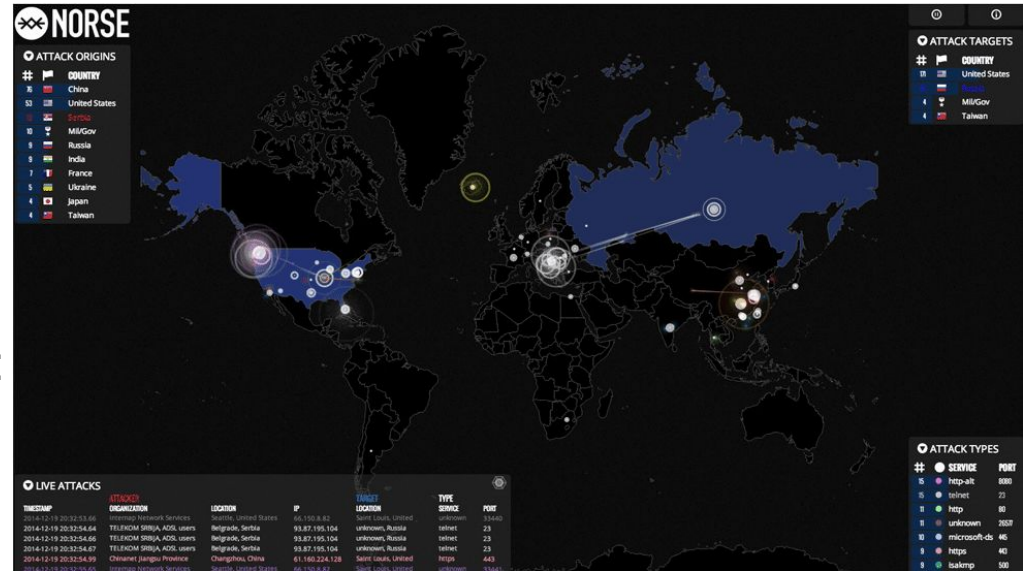
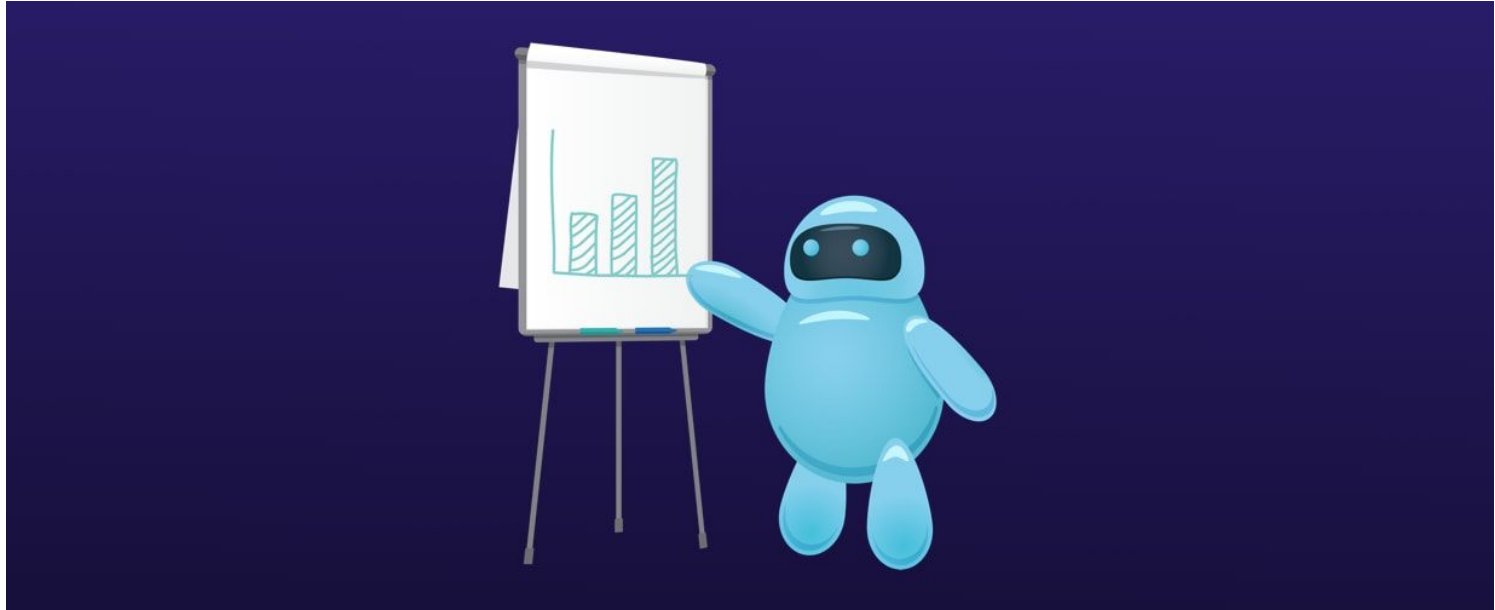If you're looking for greater penetration into your marketplace, email is a great place to start.

# Motivation: **A lot of data!**

- Windows Executables
- Android Applications
- E-mails
- Network Traffic
- Authentication events
- Operation System data such as:
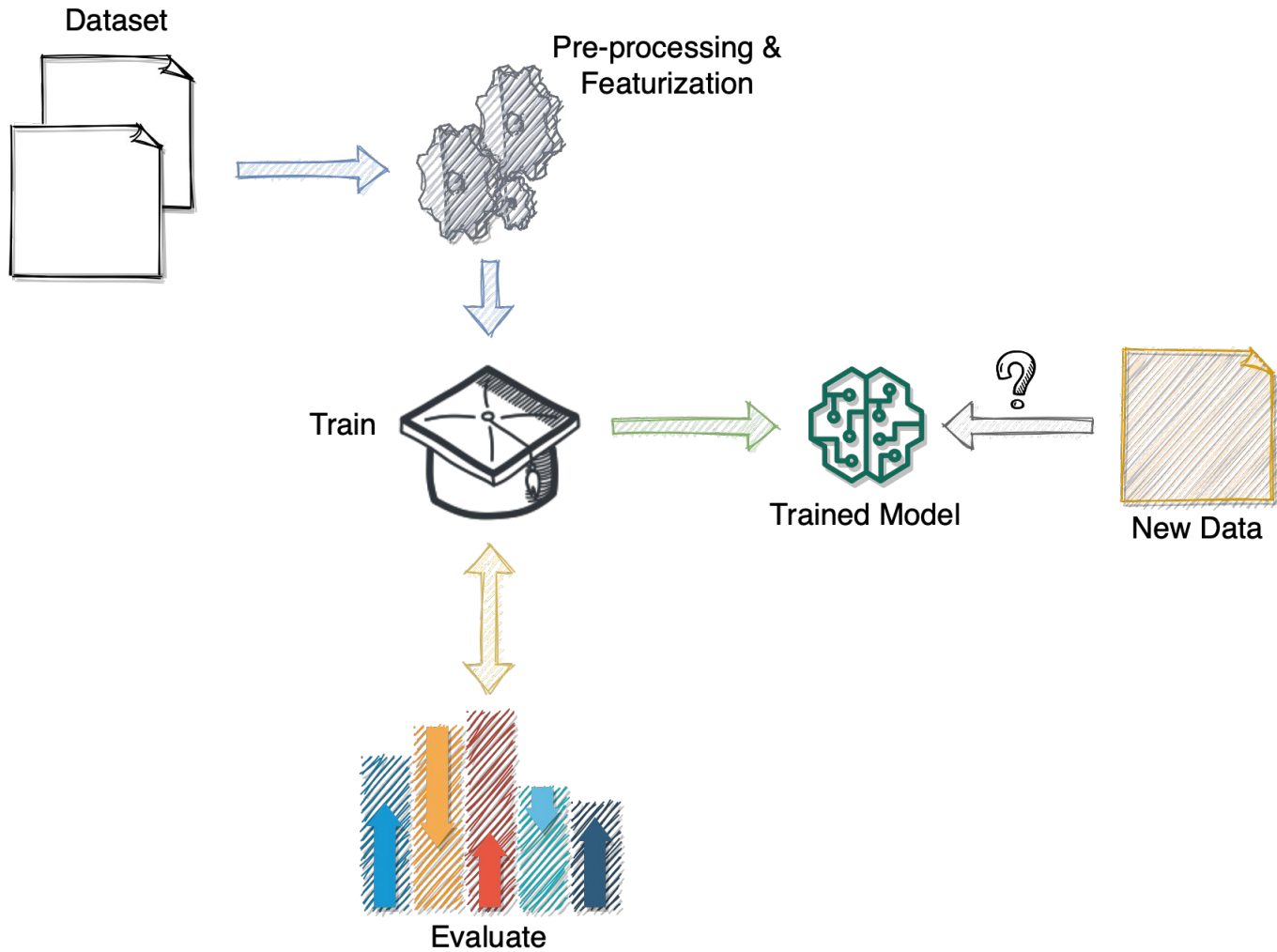  - System calls
  - Process events
- and more...

# Motivation: **A lot of data!**

- Windows Executables
- Android Applications
- E-mails
- Network Traffic
- Authentication events
- Operation System data such as:
  - System calls
  - Process events
- and more...

# How do you teach it?



[1]

Dataset

Pre-processing & Featurization

Train

Trained Model

New Data

Evaluate

# Text

In fact, the Chinese `NORP` market has the three `CARDINAL` most influential names of the retail and tech space – Alibaba `GPE` , Baidu `ORG` , and Tencent `PERSON` (collectively touted as BAT `ORG` ), and is betting big in the global AI `GPE` in retail industry space . The three `CARDINAL` giants which are claimed to have a cut-throat competition with the U.S. `GPE` (in terms of resources and capital) are positioning themselves to become the 'future AI `PERSON` platforms'. The trio is also expanding in other Asian `NORP` countries and investing heavily in the U.S. `GPE` based AI `GPE` startups to leverage the power of AI `GPE` . Backed by such powerful initiatives and presence of these conglomerates, the market in APAC AI is forecast to be the fastest-growing one `CARDINAL` , with an anticipated CAGR `PERSON` of 45% `PERCENT` over 2018 - 2024 `DATE` .

To further elaborate on the geographical trends, North America `LOC` has procured more than 50% `PERCENT` of the global share in 2017 `DATE` and has been leading the regional landscape of AI `GPE` in the retail market. The U.S. `GPE` has a significant credit in the regional trends with over 65% `PERCENT` of investments (including M&As, private equity, and venture capital) in artificial intelligence technology. Additionally, the region is a huge hub for startups in tandem with the presence of tech titans, such as Google `ORG` , IBM `ORG` , and Microsoft `ORG` .

# Assembly Instructions

# Netflow Logs

# Must speak in their language!



**Disassembled Binary**

| | | |
|---|---|---|
| 8D 8B D7 90 FE FF | lea | ecx, [ebx-16F29h] |
| 3B F9 | cmp | edi, ecx |
| **75 11** | **jnz** | **short loc_4014BB** |
| B9 39 78 00 00 | mov | ecx, 7839h |
| 2B 0D 8C 84 67 00 | sub | ecx, dword_67848C |
| 81 C3 2A 43 FF FF | add | ebx, 0FFFF432Ah |
| 3B FB | cmp | edi, ebx |
| **75 0A** | **jnz** | **short loc_4014D9** |
| C7 05 44 84 67 00 64 D3 00 00 | mov | dword_678444, 0D364h |
| 29 05 98 84 67 00 | sub | dword_678498, eax |
| . | | |
| . | | |

**Opcode Sequence**

8D, 3B, **75**, B9, 2B, 81, 3B, **75**, C7, 29, …

**N-grams with control statement shingling**

h(.)    **Hashing**    h(.)

| … | 1 | 1 | … |
|---|---|---|---|

**N-gram Frequency Vector**

Fig. 2.   Extracting opcode n-grams and hashing to reduce dimensionality

Hassen, Mehadi & Carvalho, Marco & Chan, Philip. (2017). Malware classification using static analysis based features. 1-7. 10.1109/SSCI.2017.8285426.

# What can we do?

# Malware Classification



Dataset → Pre-processing & Featurization → Train → Trained Model ← New Data
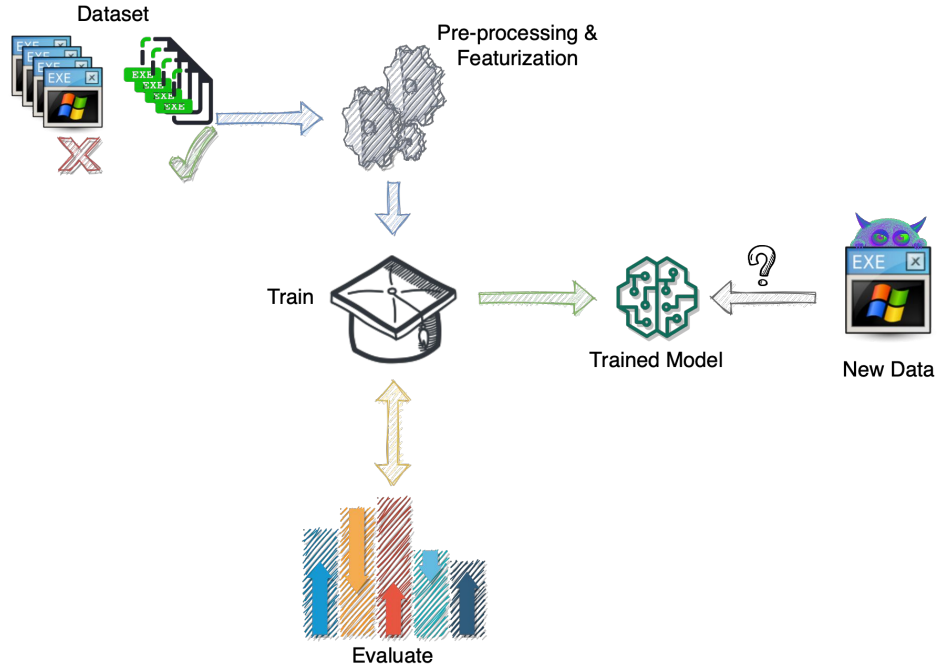
Evaluate

```
push eax #50
call DWORD PTR [ebp-0xcc] #ff9534ffffff
mov DWORD PTR [ebp-0x20], eax #8945e0
mov DWORD PTR [ebp-0xa4], 0x74726956 #c7855cffffff56697274
mov DWORD PTR [ebp-0xa0], 0x416c617f #c78560ffffff75616c41
mov DWORD PTR [ebp-0x9c], 0x636f6c6c #c78564ffffff6c6c6f63
and DWORD PTR [ebp-0x98], 0x #83a568ffffff00
lea eax, [ebp-0xa4] #8d855cffffff
push eax #50
push DWORD PTR [ebp+0xe] #ff750e
xor bh,bh #30ff
xchg ebp,eax #95
cmp bh,bh #95
.byte 0xff #ff
```

**Figure 3: Example of a disassembled 64-gram feature found in the EMBER dataset. The hex values of the raw bytes are shown in comments for each line of assembly.**

Raff, E., Fleming, W., Zak, R., Anderson, H., Finlayson, B., Nicholas, C., & McLean, M. (2019). KiloGrams: Very Large N-Grams for Malware Classification. ArXiv, abs/1908.00200.
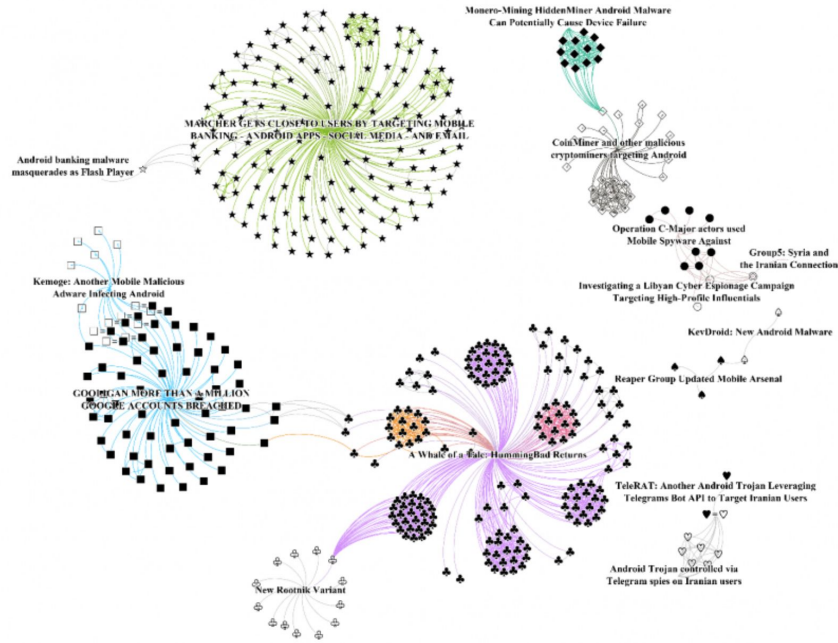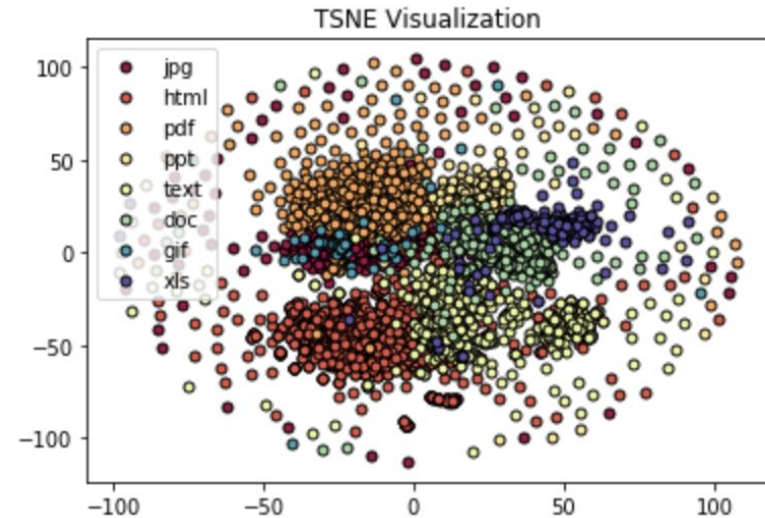
# Malware Clustering



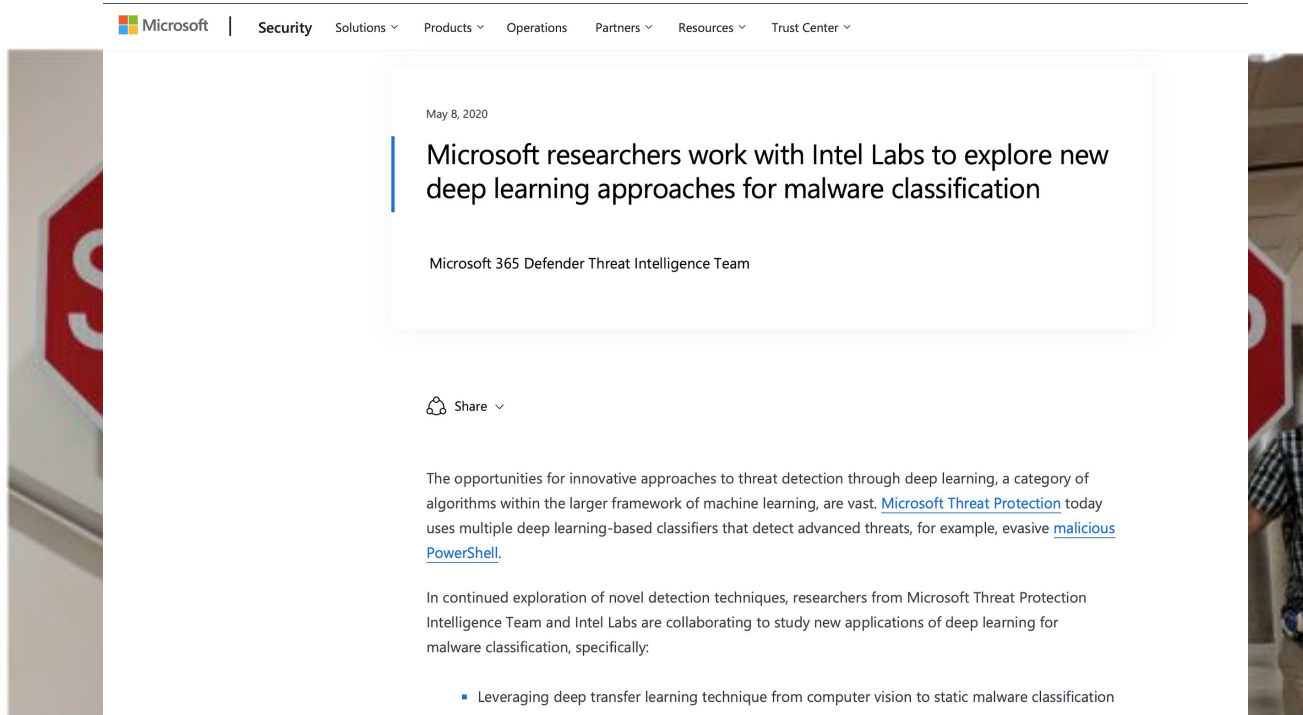Figure 17: Relationship of 15 AlienVault OTX reports.

Lee et al. Dexofuzzy: Android malware similarity clustering method using opcode sequence.



Raff, E., & Nicholas, C. (2017). An Alternative to NCD for Large Sequences, Lempel-Ziv Jaccard Distance. *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.

# Adversarial Machine Learning

Goodfellow, I. (2020, October 05). Attacking Machine Learning with Adversarial Examples. Retrieved October 21, 2020, from https://openai.com/blog/adversarial-example-research/

# Statistical User Behaviour Analysis



Figure 9: Event times for User205265. 4624l2 corresponds to *EventID* 4624 - *LogonType* 2.



Figure 8: Daily count of the fields in Figure 7.

Turcotte, M.J., Kent, A., & Hash, C. (2017). Unified Host and Network Data Set. ArXiv, abs/1708.07518.
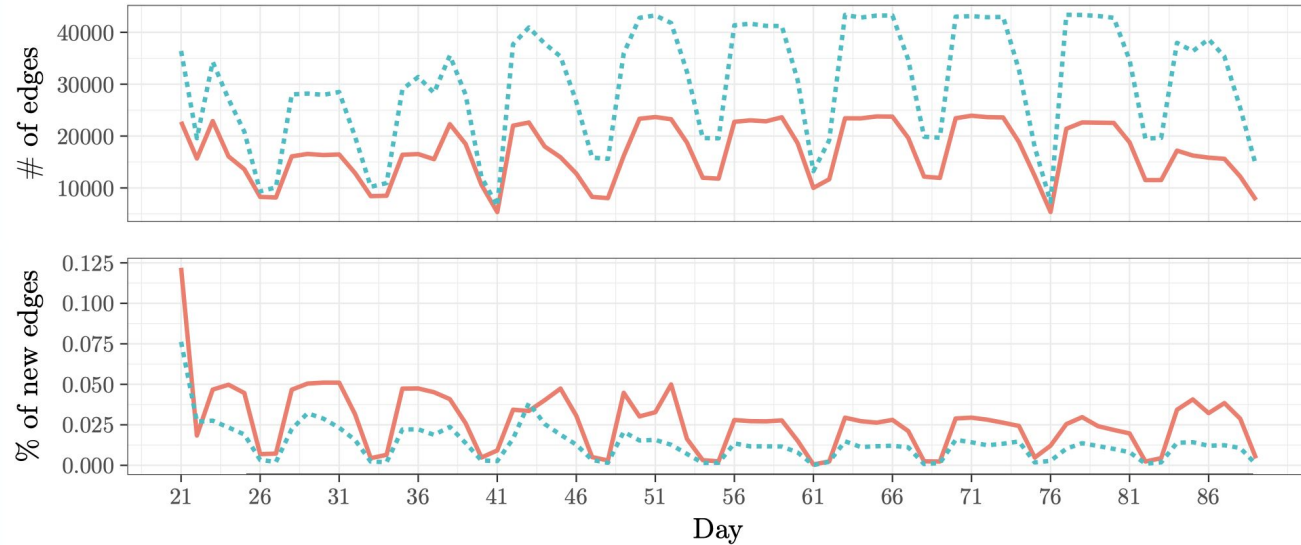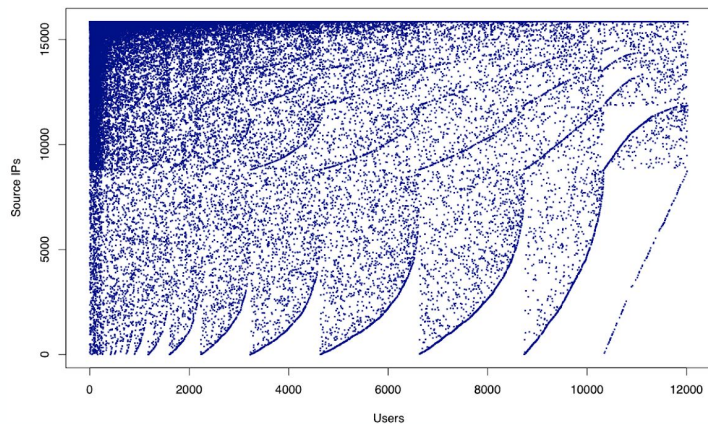
FIG 1: *Number of links per day (top), and proportion of those that are new (bottom), after 20 days of observation of the LANL computer network.* **Solid red** *curve: User − Source.* **Dashed blue** *curve: User − Destination.*

Turcotte, M.J., Kent, A., & Hash, C. (2017). Unified Host and Network Data Set. ArXiv, abs/1708.07518.

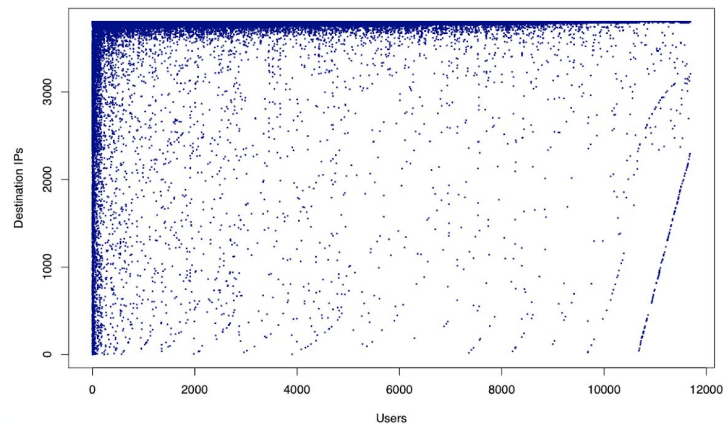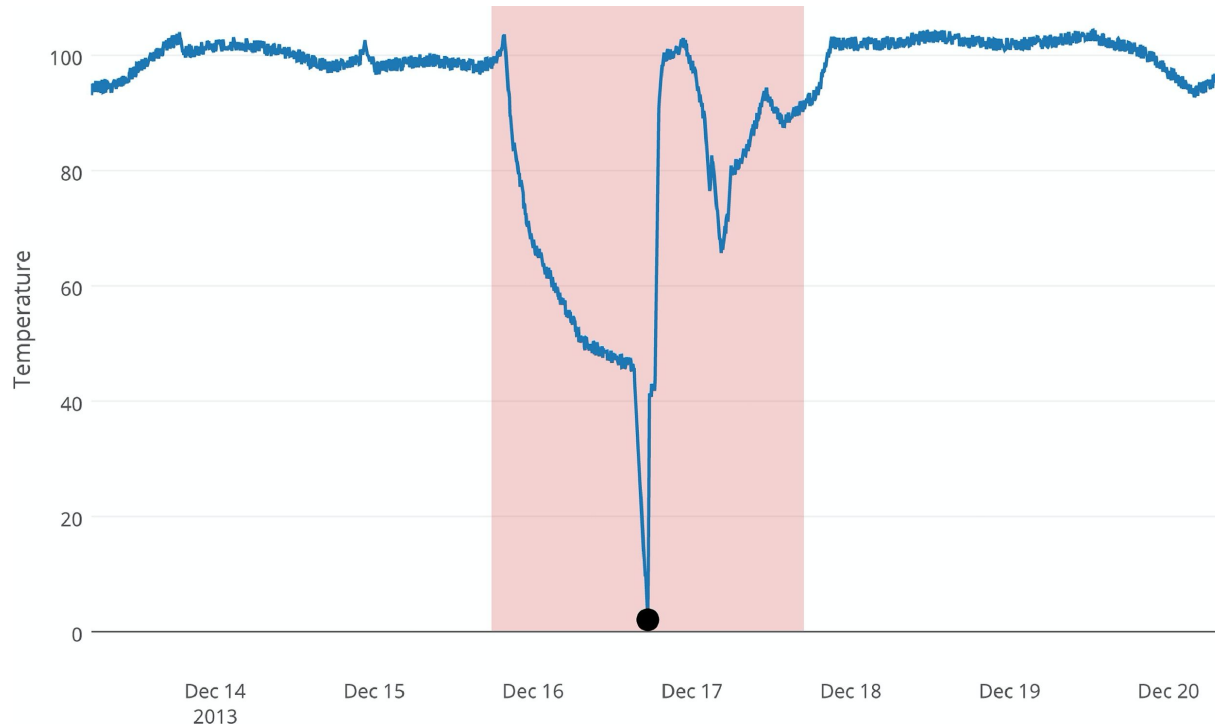# Anomaly Detection: Link Prediction



(A) *User − Source*    (B) *User − Destination*

FIG 2: *Training set adjacency matrices for the two graphs (spy-plot). Nodes are sorted by in-degree and out-degree.*

Passino, F.S., Turcotte, M.J., & Heard, N. (2020). Graph link prediction in computer networks using Poisson matrix factorisation. ArXiv, abs/2001.09456.

# Anomaly Detection: Time Series



Ahmad, Subutai & Lavin, Alexander & Purdy, Scott & Agha, Zuha. (2017). Unsupervised real-time anomaly detection for streaming data. Neurocomputing. 10.1016/j.neucom.2017.04.070.

# Woa!! Where do I sign-up?

- Books
- Malware Research Group
- Internships
- Projects on your free time



Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems 2nd Edition
by Aurélien Géron (Author)



Malware Data Science: Attack Detection and Attribution Paperback – September 25, 2018
by Joshua Saxe (Author), Hillary Sanders (Author)

# References

[1] Sarver, C. (n.d.). Introduction to Regression and Classification in Machine Learning. Retrieved from https://www.springboard.com/blog/introduction-regression-classification-machine-learning/

[2] Email Usage Statistics in 2019. (n.d.). Retrieved October 21, 2020, from https://www.campaignmonitor.com/blog/email-marketing/2019/07/email-usage-statistics-in-2019/

[3] Malware Statistics & Trends Report: AV-TEST. (n.d.). Retrieved October 21, 2020, from https://www.av-test.org/en/statistics/malware/

[4] Microsoft researchers work with Intel Labs to explore new deep learning approaches for malware classification. (2020, May 08). Retrieved October 21, 2020, from https://www.microsoft.com/security/blog/2020/05/08/microsoft-researchers-work-with-intel-labs-to-explore-new-deep-learning-approaches-for-malware-classification/

[5] Goodfellow, I. (2020, October 05). Attacking Machine Learning with Adversarial Examples. Retrieved October 21, 2020, from https://openai.com/blog/adversarial-example-research/

[6] Dexofuzzy: Android malware similarity clustering method using opcode sequence. (n.d.). Retrieved October 21, 2020, from https://www.virusbulletin.com/virusbulletin/2019/11/dexofuzzy-android-malware-similarity-clustering-method-using-opcode-sequence/

[7] Raff, E., & Nicholas, C. (2017). An Alternative to NCD for Large Sequences, Lempel-Ziv Jaccard Distance. *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.

# References

[8] Turcotte, M.J., Kent, A., & Hash, C. (2017). Unified Host and Network Data Set. ArXiv, abs/1708.07518.

[9] Passino, F.S., Turcotte, M.J., & Heard, N. (2020). Graph link prediction in computer networks using Poisson matrix factorisation. ArXiv, abs/2001.09456.

[10] Raff, E., Fleming, W., Zak, R., Anderson, H., Finlayson, B., Nicholas, C., & McLean, M. (2019). KiloGrams: Very Large N-Grams for Malware Classification. ArXiv, abs/1908.00200.

[11] Ahmad, Subutai & Lavin, Alexander & Purdy, Scott & Agha, Zuha. (2017). Unsupervised real-time anomaly detection for streaming data. Neurocomputing. 10.1016/j.neucom.2017.04.070.

[12] Hassen, Mehadi & Carvalho, Marco & Chan, Philip. (2017). Malware classification using static analysis based features. 1-7. 10.1109/SSCI.2017.8285426.