# Before We Get Started

**Why Are We Here?**

1. About TruSTAR

2. Questions are encouraged

3. We'll send you all resources (Slides, GitHub, Blog) after the presentation.

**Agenda**

1. Overview of Data Challenges
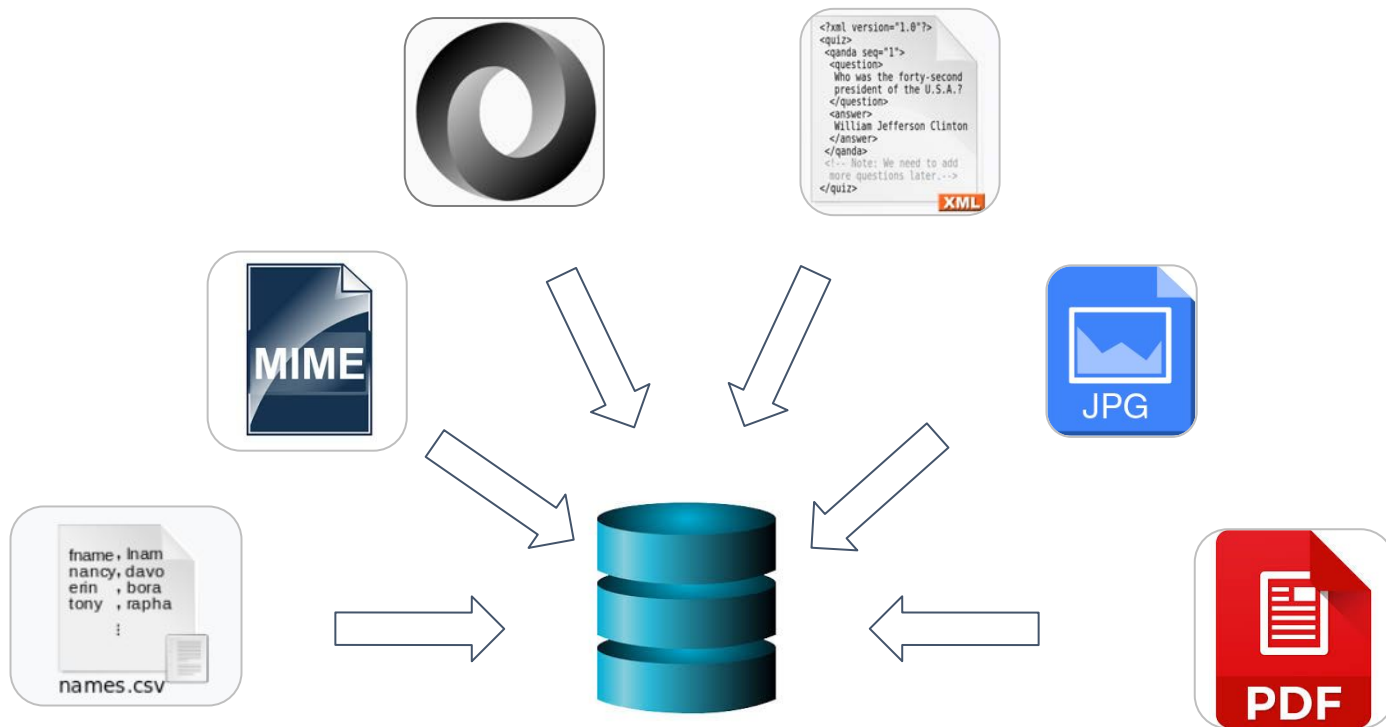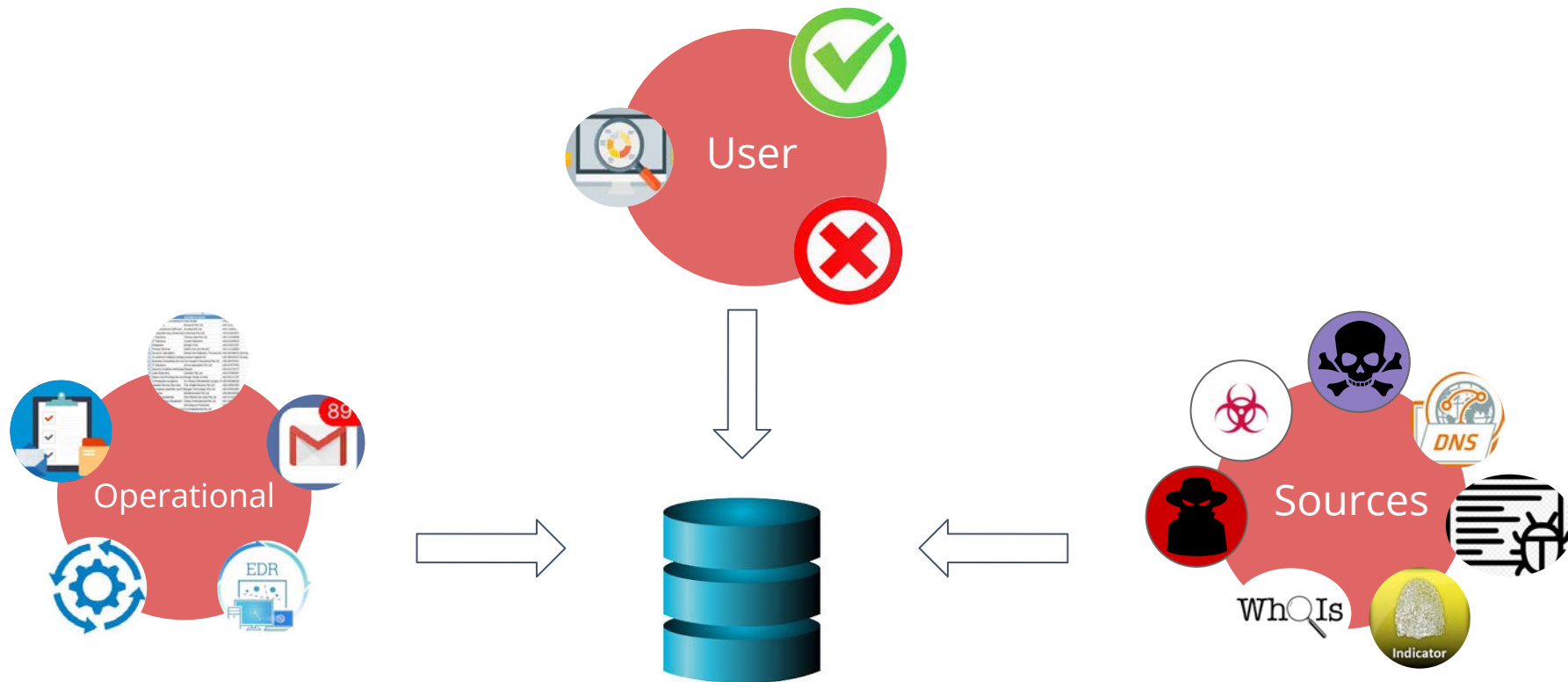2. Why do we want to solve these challenges
3. How are solving them
4. NLP

# I. Overview of Data Challenges in Intelligence

# Variety of Data Formats

# Variety of Data Sources

# Variety of Data Sources

**Structured**

| action | date | device_id |
|--------|------|-----------|
| Open | 2019-04-03 | 8 |
| Close | 2019-04-18 | 57 |
| Open | 2019-04-21 | 69 |
| Close | 2019-04-14 | 43 |
| Open | 2019-04-15 | 49 |
| Close | 2019-04-18 | 59 |
| | 2019-04-18 | 57 |

# Variety of Data Sources

**Semi-structured**

value

▼ object
    guid: 84743018-18cf-4598-995c-a2c2309f168a
    eventType: REPORT_UPDATE
  ▶ report: {"title":"alienvaulttotx-IP-221.213.54.76","body":"{\n   \"sections\": {\n     \"general\": {\n     \"area_code\": 0,\n     \"asn\": \"AS4837 CHINA UNICOM China169
    \"base_indicator\": {\n     \"access_reason\": \"\",\n     \"access_type\": \"public\",\n     \"content\": \"\",\n     \"description\": \"\",\n     \"id\": 1
    \"221.213.54.76\",\n     \"title\": \"\",\n     \"type\": \"IPv4\"\n   },\n     \"charset\": 0,\n     \"city\": \"Kunming\",\n     \"city_data\": true,\n
    \"country_code\": \"CN\",\n     \"country_code3\": \"CHN\",\n     \"country_name\": \"China\",\n     \"dma_code\": 0,\n     \"flag_title\": \"China\",\n     \"fla
    \"indicator\": \"221.213.54.76\",...
    reportEntities: null

▼ object
    guid: 84743018-18cf-4598-995c-a2c2309f168a
    eventType: REPORT_ENTITY_EXTRACT
    report: null
  ▼ reportEntities:
    ▶ entities: [{"guid":"SOFTWARE|/static/img/flags/cn.png","type":"SOFTWARE","value":"/static/img/flags/cn.png"},
      {"guid":"URL|http://whois.domaintools.com/221.213.54.76","type":"URL","value":"http://whois.domaintools.com/221.213.54.76"},{"guid":"IP|221.213.54.76","type":"IP","va
      {"guid":"URL|bl.nszones.com","type":"URL","value":"bl.nszones.com"},{"guid":"URL|asp.net","type":"URL","value":"asp.net"},{"guid":"URL|strict.dtd","type":"URL","value":
      {"guid":"SOFTWARE|welcome.png","type":"SOFTWARE","value":"welcome.png"},{"guid":"URL|https://otx20-web-
      media.s3.amazonaws.com/media/avatars/user_32073/resized/80/avatar_2e866cf4cc.png","type":"URL","value":"https://otx20-web-
      media.s3.amazonaws.com/media/avatars/user_32073/resized/80/avatar_2e866cf4cc.png"}]
    ▶ enclaveGuids: ["b077bfb9-82e4-47f2-ae92-61a939333a33"]

# Variety of Data Sources

# Extraction Challenges

## 1. Extract

135N2nfAkextd6...

**CryptoLocker**

135N2nfAkextd6...
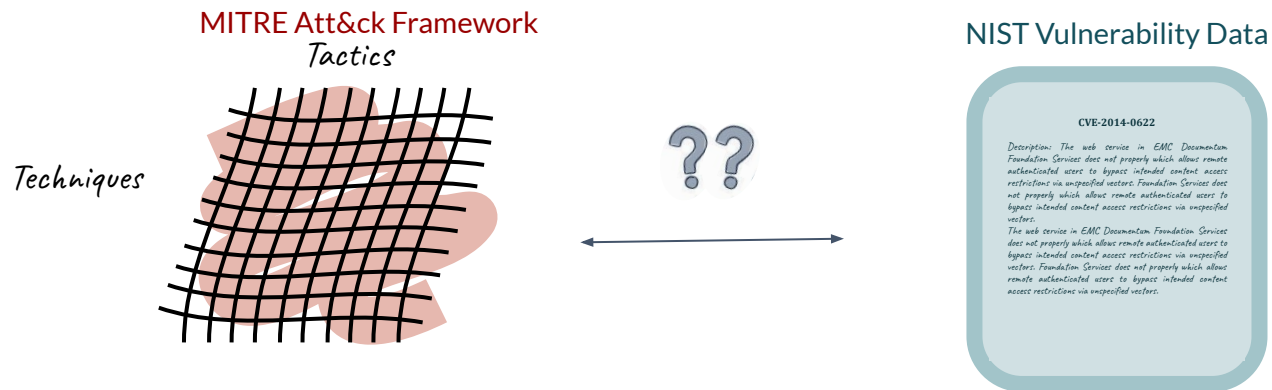135N2nfAkextd6...

**Bad Rabbit**

135N2nfAkextd6...

**NotPetya**

## 2. Disambiguate & Normalize

*Nuclear* is an Exploit Kit!

*Germany's 17 **nuclear** power plants will be shut down in 2022.*

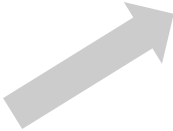| IP (4) | | |
|---|---|---|
| 6.0.5.1 | VERSION NUMBER | ✖ |
| 6.0.5.0 | VERSION NUMBER | ✖ |
| 121.2.3.5 | | |
| 121.2.3.4 | | |

# Inferring Relationships / Categorization

MITRE Att&ck Framework

*Tactics*

*Techniques*

??

NIST Vulnerability Data

**CVE-2014-0622**

*Description: The web service in EMC Documentum Foundation Services does not properly which allows remote authenticated users to bypass intended content access restrictions via unspecified vectors. Foundation Services does not properly which allows remote authenticated users to bypass intended content access restrictions via unspecified vectors.*

*The web service in EMC Documentum Foundation Services does not properly which allows remote authenticated users to bypass intended content access restrictions via unspecified vectors. Foundation Services does not properly which allows remote authenticated users to bypass intended content access restrictions via unspecified vectors.*

# Merging Sources into 1/N Standard Languages

| key | confidence | malware |
|-----|-----------|---------|
| 1.1.1.1\|ip | high | Mofksys |
| evil.com\|domain | low | Ryuk |

**Extract / Categorization**

| key | asn | country |
|-----|-----|---------|
| 1.1.1.1\|ip | 3265 | NLD |
| evil.com\|domain | 1668 | USA |

# Merging Sources into 1/N Standard Languages

Merge

| key | confidence | malware | asn | country |
|---|---|---|---|---|
| 1.1.1.1|ip | high | Mofksys | 3265 | NLD |
| evil.com|domain | low | Ryuk | 1668 | USA |

**II.**  **Why do we want to solve these challenges**

# Information Retrieval: Search & Querying

# Information Retrieval: Link Analysis

arturn.co.uk

Operational

Sources

```
object ▶ indicator
▼ object {13}
      indicator : arturn.co.uk
      type : domain
      published_date : 1553884134
      last_updated : 1553884137
   ▶ reports [0]
   ▶ actors [0]
   ▼ malware_families [1]
         0 : Mofksys
   ▶ kill_chains [0]
   ▶ ip_address_types [0]
   ▶ domain_types [0]
      malicious_confidence : high
   ▶ labels [5]
   ▶ relations [14]
```

```
object ▶ description
▼ object {3}
      short_description : Suspicious Domain
      description : Through using virustotal we discovered this
                    suspicious url arturn.co.uk
      incident_id : 43
```

# Prioritization and Triage

MITRE Att&ck Framework

*Tactics*

*Techniques*

**CVE-2017-5161**

An issue was discovered in Sielco Sistemi Winlog Lite SCADA Software. An uncontrolled search path element (DLL Hijacking) vulnerability has been identified. Exploitation of this vulnerability could give an attacker access to the system with the same level of privilege as the application that utilizes the malicious DLL.

■ Weak posture

■ Strong posture

**Sielco Sistemi Winlog Lite SCADA**

**Fuji Electric Energy Savings Estimator v1.0.0**

**BLF-Tech LLC VisualView HMI v9.9.14.0**

**SIMPlight SCADA Software version 4.3.0.27**

**Moxa SoftNVR-IA Live Viewer v3.30.3122**

**III.** **How are we solving for these challenges**

# Using regex for entity extraction

I start with a pattern to extract specific entities

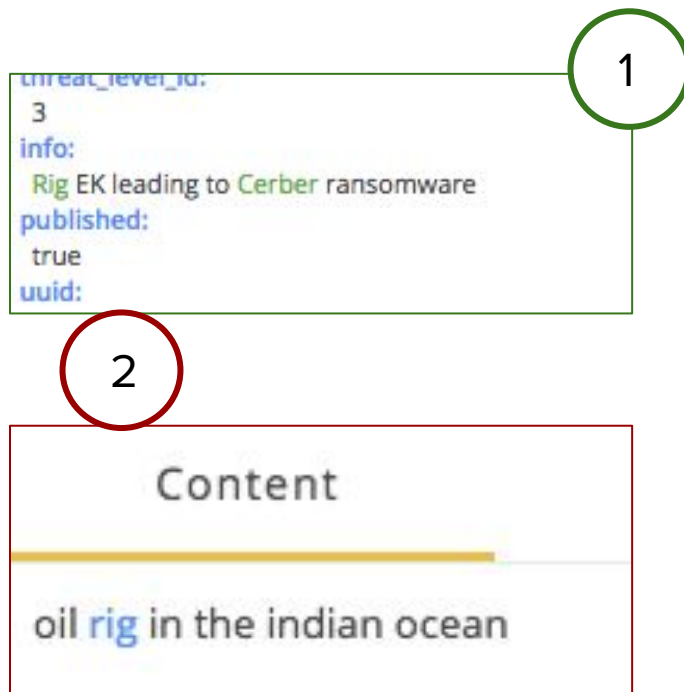$$\b\d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3}\b$$

**Upside:**

Easy to implement...

**Downside:**

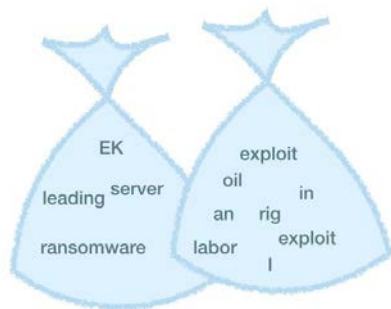Low precision due to a high number of false positives

Requires a priori knowledge

# ML to the rescue

threat_level_id:
3
info:
 Rig EK leading to Cerber ransomware
published:
 true
uuid:

**1**

**2**

Content

oil rig in the indian ocean

# ML to the rescue

**1. Transform content to Contextual bag of words**



EK
server
leading
ransomware
oil
an rig
labor
exploit
in
exploit
l

**2. Vectorize word counts and compute TF-IDF**



|  | Doc 1 | Doc 2 | Doc 3 | Doc 4 | Doc 5 | Doc 6 | Doc 7 | Doc 8 |
|---|---|---|---|---|---|---|---|---|
| Term 1 | 10 | 0 | 1 | 0 | 0 | 0 | 0 | 2 |
| Term 2 | 0 | 2 | 0 | 0 | 0 | 18 | 0 | 2 |
| Term 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| Term 4 | 6 | 0 | 0 | 4 | 6 | 0 | 0 | 0 |
| Term 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| Term 6 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| Term 7 | 0 | 1 | 8 | 0 | 0 | 0 | 0 | 0 |
| Term 8 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 |

← Word Vector (Passage Vector)

↑ Document Vector

**3. Feed into a simple model**



$$P(c \mid x) = \frac{P(x \mid c)P(c)}{P(x)}$$

Likelihood · Class Prior Probability
Posterior Probability · Predictor Prior Probability

$$P(c \mid X) = P(x_1 \mid c) \times P(x_2 \mid c) \times \cdots \times P(x_n \mid c) \times P(c)$$

**4. Predict on production**



RIG — COMMON WORD ✗

MALWARE (1)
RIG

# But it's not always enough!

New Malware names

# But it's not always enough!

MITRE Att&ck Framework

*Tactics*

*Techniques*

??

NIST Vulnerability Data

**CVE-2014-0622**

Description: The web service in EMC Documentum Foundation Services does not properly which allows remote authenticated users to bypass intended content access restrictions via unspecified vectors. Foundation Services does not properly which allows remote authenticated users to bypass intended content access restrictions via unspecified vectors.
The web service in EMC Documentum Foundation Services does not properly which allows remote authenticated users to bypass intended content access restrictions via unspecified vectors. Foundation Services does not properly which allows remote authenticated users to bypass intended content access restrictions via unspecified vectors.

# IV.   Which brings us to NLP…

Work by **Zainab Danish** ([zdanish@trustar.co](mailto:zdanish@trustar.co))

# Mitre ATT&CK Framework

- **Tactic = Why?**
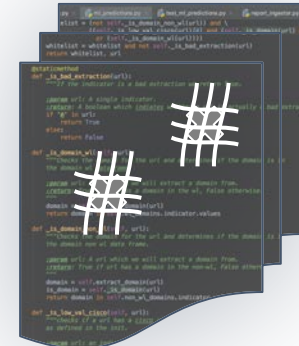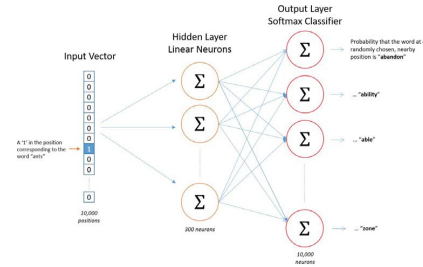- **Technique = How?**

# Mitre ATT&CK Framework



Tactics

Techniques

Initial Access

Drive-by Compromise

Exploit Public-Facing Application

Hardware Applications

Replication via Removable Media

Spearphishing Attachment

Spearphishing Link

Spearphishing via Service

Supply Chain Compromise

Trusted Relationship

Valid Accounts

Threat Actor Groups

Mitigation Strategies

Detection

References

# Consider the following...



Adversary Group — Uses → Technique

Adversary Group — Uses → Malware

Malware — Implements → Technique

Vulnerability — Implements → Technique

Technique — Accomplishes → Tactic

Malware — Targets → Vulnerability

# Word2Vec

| I | see | ants | on | the | tree.

(ants, I)
(ants, see)
(ants, on)
(ants, the)

**What?** NLP technique that seeks to teach the computer to **understand, interpret and manipulate human language**.

**Why?** **Translate words into vectors** for mathematical manipulation.

**How?** By **leveraging context** and calculating probabilities.



ants    =    $[a_1 \, a_2 \, a_3 \, a_4 \, a_5 \ldots\ldots\ldots\ldots a_d]$

# Doc2Vec

**What?** NLP technique that seeks to teach the computer to **understand, interpret and manipulate human language**.

**Why?** **Translate words and documents into vectors** for mathematical manipulation.

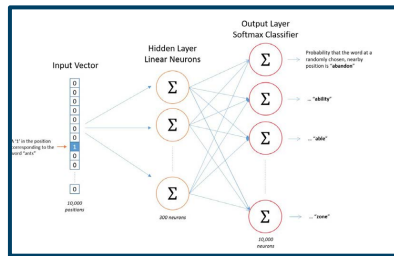**How?** By **leveraging context** and calculating probabilities.



CVE-2014-0622
Description: The web service in EMC Documentum Foundation Services does not properly implement content uploading, which allows remote authenticated users to bypass intended content access restrictions via unspecified vectors.

CVE-2014-0622
['the', 'web', 'service', 'in', 'emc', 'documentum', 'foundation', 'services', 'does', 'not', 'properly', 'implement', 'content', 'uploading', 'which', 'allows', 'remote', 'authenticated', 'users', 'to', 'bypass', 'intended', 'content', 'access', 'restrictions', 'via', 'unspecified', 'vectors']

CVE-2014-0622    =    $[a_1\ a_2\ a_3\ a_4\ a_5\ .........\ a_d]$

# Process

**CVE-2014-0622**

Description: The web service in EMC Documentum Foundation Services does not properly implement content uploading, which allows remote authenticated users to bypass intended content access restrictions via unspecified vectors.

**Data Cleaning**

**CVE-2014-0622**

['the', 'web', 'service', 'in', 'emc', 'documentum', 'foundation', 'services', 'does', 'not', 'properly', 'implement', 'content', 'uploading', 'which', 'allows', 'remote', 'authenticated', 'users', 'to', 'bypass', 'intended', 'content', 'access', 'restrictions', 'via', 'unspecified', 'vectors']

**Tokenization**



Input Vector

Hidden Layer
Linear Neurons

Output Layer
Softmax Classifier

Probability that the word at a randomly chosen, nearby position is "abandon"

... "ability"

... "able"

... "zone"

**Model Training**

CVE-2014-0622 $[a_1\ a_2\ a_3\ a_4\ a_5 \ldots\ldots\ a_d]$
CVE-2015-0765 $[a_1\ a_2\ a_3\ a_4\ a_5 \ldots\ldots\ a_d]$
...
attack-pattern12 $[a_1\ a_2\ a_3\ a_4\ a_5 \ldots\ldots\ a_d]$

**Numeric Vectors**

# CVE Clusters: Data Categorization / Tagging

**Initial number of docs:**

~100,000

**Total discovered clusters:**

~100

Browser Vulnerabilities

Overflow Vulnerabilities

Privilege Escalation Vulnerabilities

# CVE Clusters: Privilege Escalation

The compilation daemon in Scala before 2.10.7 uses weak permissions for private files in scalac-compile-server-port, which allows local users to write to arbitrary class files and consequently gain privileges.

The overlayfs implementation in the Linux kernel through 4.5.2 does not properly restrict the mount namespace, which allows local users to gain privileges by mounting an overlayfs filesystem on top of a FUSE filesystem, and then executing a crafted setuid program.

Directory traversal vulnerability in the Shared Folders feature for VMware Workstation before 5.5.4, when a folder is shared, allows users on the guest system to write to arbitrary files on the host system via the "Backdoor I/O Port" interface.

The Outlook Extension in IBM Content Collector 4.0.0.x before 4.0.0.0-ICC-OE-IF004 allows local users to bypass the intended Reviewer privilege requirement and read e-mail messages from an arbitrary mailbox by invoking the Search function.

# Connecting the two worlds...

**Dylib Hijacking**
macOS and OS X use a common method to look for required dynamic libraries (dylib) to load into a program based on search paths. Adversaries can take advantage of ambiguous paths to plant dylibs to gain privilege escalation or persistence. A common method is to see what dylibs an application uses, then plant a malicious version with the same name higher up in the search path. This typically results in the dylib being in the same folder as the application itself. If the program is configured to run at a higher privilege level than the current user, then when the dylib is loaded into the application, the dylib will also run at that elevated level. This can be used by adversaries as a privilege escalation technique.

**CVE-2017-6329**

Symantec VIP Access for Desktop prior to 2.2.4 can be susceptible to a DLL Pre-Loading vulnerability. These types of issues occur when an application looks to call a DLL for execution and an attacker provides a malicious DLL to use instead. Depending on how the application is configured, the application will generally follow a specific search path to locate the DLL. The exploitation of the vulnerability manifests as a simple file write (or potentially an over-write) which results in a foreign executable running under the context of the application.

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|\|\mathbf{B}\|} \quad > \quad \mathbf{0.5}$$

# Evaluation

100 naturally
forming clusters

1/100 chance of
getting it right at
random

1% accurate
associations

50% accurate
associations

# What next...

## Limitations:

- Language and schema differences deteriorate the quality of the models.
- Hard to scale for many to many relationships.

## Potential solutions:

- Investigate training models on recent data only.
- Investigate more sophisticated Neural Networks algorithms (RadialGANs).

# QUESTIONS?

TRU★STAR

**Access Our Codebase & Resources Here:**
**TRUSTAR.CO / NLP**

**Nicolas Kseib**

Twitter: @NKseib
Email: nkseib@trustar.co

**Zainab Danish**

Email: zdanish@trustar.co

# THANK YOU!

**TRU★STAR**

---

Sample Enterprise Customers

**rackspace**
the #1 managed cloud company

**SallieMae.**

**Prudential**

**experian.**

**STARBUCKS COFFEE**

Sample Hosted Exchange Groups

**IT ISAC**

**CSA** cloud security alliance®

**Columbus Collaboratory**

**The Academy**

**R-CISC**

**sports** -ISAO

# TRU★STAR

**Extra**

# Faster Incident Response

**CVE-2017-5161**

An issue was discovered in Sielco Sistemi Winlog Lite SCADA Software. An uncontrolled search path element (**DLL Hijacking**) vulnerability has been identified. Exploitation of this vulnerability could give an attacker access to the system with the same level of privilege as the application that utilizes the malicious **DLL**.

## Dylib Hijacking

macOS and OS X use a common method to look for required dynamic libraries (dylib) to load into a program based on search paths. Adversaries can take advantage of ambiguous paths to plant dylibs to gain privilege escalation or persistence. A common method is to see what dylibs an application uses, then plant a malicious version with the same name higher up in the search path. This typically results in the dylib being in the same folder as the application itself. If the program is configured to run at a higher privilege level than the current user, then when the dylib is loaded into the application, the dylib will also run at that elevated level. This can be used by adversaries as a privilege escalation technique.

## DETECTION

Objective-See's Dylib Hijacking Scanner can be used to detect potential cases of dylib hijacking. Monitor file systems for moving, renaming, replacing, or modifying dylibs. Changes in the set of dylibs that are loaded by a process (compared to past behavior) that do not correlate with known software, patches, etc., are suspicious. Check the system for multiple dylibs with the same name and monitor which versions have historically been loaded into a process.

## MITIGATION

Prevent users from being able to write files to the search paths for applications, both in the folders where applications are run from and the standard dylib folders. If users can't write to these directories, then they can't intercept the search path.

## REFERENCES

**Malware Persistence on OS X :**
https://www.rsaconference.com/writable/presentations/file_upload/ht-r03-malware-persistence-on-os-x-yosemite_final.pdf
**Writing Bad Malware for OS X:**
https://www.blackhat.com/docs/us-15/materials/us-15-Wardle-Writing-Bad-A-Malware-For-OS-X.pdf
 **Mitre-attack:**
 https://attack.mitre.org/techniques/T1157

# MITRE Att&ck Framework: Possible Improvements

**Exact match**

| attack_id | attack_name | attack_description |
|---|---|---|
| attack-pattern--03f4a766-7a21-4b5e-9ccf-e0cf422ab983 | Acquire or compromise 3rd party signing certificates | Code signing is the process of digitally signing executables and scripts to confirm the software author and guarantee that the code has not been altered or corrupted. Users may trust a signed piece of code more than an signed piece of code even if they don't know who issued the certificate or who the author is. (Citation: DiginotarCompromise) |
| attack-pattern--e5164428-03ca-4336-a9a7-4d9ea1417e59 | Acquire or compromise 3rd party signing certificates | Code signing is the process of digitally signing executables or scripts to confirm the software author and guarantee that the code has not been altered or corrupted. Users may trust a signed piece of code more than an signed piece of code even if they don't know who issued the certificate or who the author is. (Citation: Adobe Code Signing Cert) |

**Similar meaning**

| attack_id | attack_name | attack_description |
|---|---|---|
| attack-pattern--784ff1bc-1483-41fe-a172-4cd9ae25c06b | Acquire OSINT data sets and information | Open source intelligence (OSINT) is intelligence gathered from publicly available sources. This can include both information gathered on-line, such as from search engines, as well as in the physical world. (Citation: RSA-APTRecon) |
| attack-pattern--2b9a666e-bd59-4f67-9031-ed41b428e04a | Acquire OSINT data sets and information | Open source intelligence (OSINT) provides free, readily available information about a target while providing the target no indication they are of interest.  Such information can assist an adversary in crafting a successful approach for compromise. (Citation: RSA-APTRecon) |

# Organizational Prioritization

**Dylib Hijacking**

macOS and OS X use a common method to look for required dynamic libraries (dylib) to load into a program based on search paths. Adversaries can take advantage of ambiguous paths to plant dylibs to gain privilege escalation or persistence. A common method is to see what dylibs an application uses, then plant a malicious version with the same name higher up in the search path. This typically results in the dylib being in the same folder as the application itself. If the program is configured to run at a higher privilege level than the current user, then when the dylib is loaded into the application, the dylib will also run at that elevated level. This can be used by adversaries as a privilege escalation technique.

**CVE-2017-5161**

An issue was discovered in Sielco Sistemi Winlog Lite SCADA Software. An uncontrolled search path element (**DLL Hijacking**) vulnerability has been identified. Exploitation of this vulnerability could give an attacker access to the system with the same level of privilege as the application that utilizes the malicious **DLL**.

❌ **Sielco Sistemi Winlog Lite SCADA**

❌ **Fuji Electric Energy Savings Estimator v1.0.0**

❌ **BLF-Tech LLC VisualView HMI v9.9.14.0**

❌ **SIMPlight SCADA Software version 4.3.0.27**

❌ **Moxa SoftNVR-IA Live Viewer v3.30.3122**

# Improvement Strategies

- Heuristics used

# Unexpected (but cool) outcomes

# Evaluation

Human evaluation

- On CVE similarity
- On technique relationship

# Vulnerability Data: NIST NVD

"A weakness in the computational logic (e.g., code) found in software and hardware components that, when exploited, results in a negative impact to confidentiality, integrity, or availability. Mitigation of the vulnerabilities in this context typically involves coding changes, but could also include specification changes or even specification deprecations (e.g., removal of affected protocols or functionality in their entirety)."
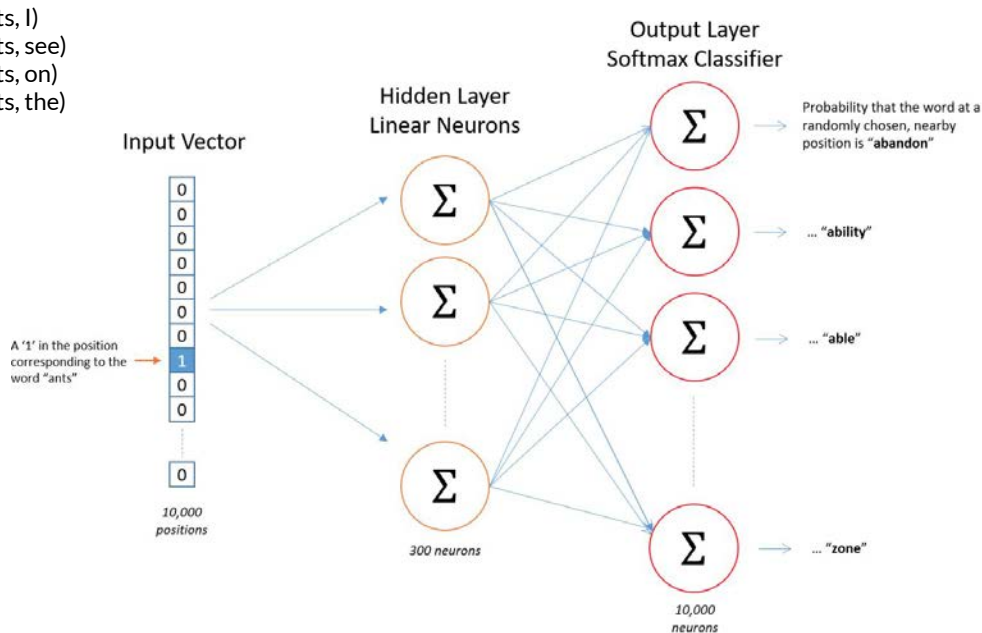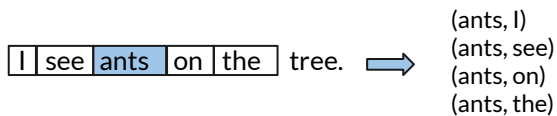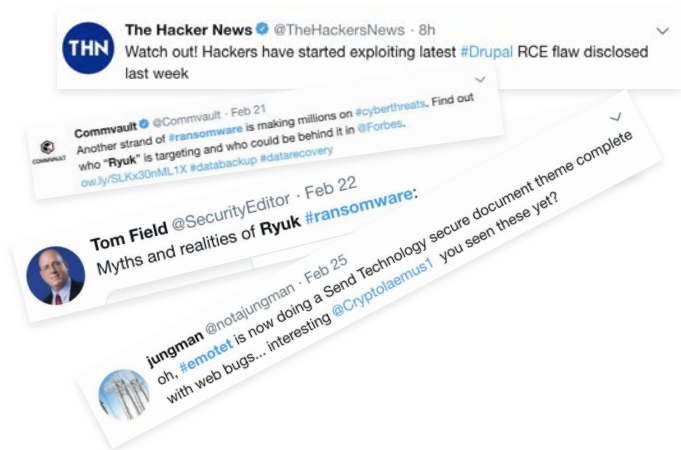
# Speaking of NLP



IoC Extraction

Recommendation

Data Summarization

Enrichment

# Method to the madness?


Unstructured Threat Intel

STIX - Structured Threat Intelligence Expression is **language and serialization format** used to exchange cyber threat intelligence (**CTI**).

# Word2Vec

# Word2Vec: Malware Example



Entity Extraction



Related Trends

# MITRE Att&ck Framework: Possible Improvements

Clustering shows:

- Overlapping ATT&CK techniques
  - Exact matches
  - Conveying the same meaning
- Overarching categories for some techniques



| attack_id | technique_name | description |
| --- | --- | --- |
| attack-pattern--2b9a666e-bd59-4f67-9031-ed41b428e04a | Acquire OSINT data sets and information | Open source intelligence (OSINT) provides free, readily available information about a target while providing the target no indication they are of interest. Such information can assist an adversary in crafting a successful approach for compromise. (Citation: RSA-APTRecon) |
| attack-pattern--028ad431-84c5-4eb7-a364-2b797c234f88 | Acquire OSINT data sets and information | Data sets can be anything from Security Exchange Commission (SEC) filings to public phone numbers. Many datasets are now either publicly available for free or can be purchased from a variety of data vendors. Open source intelligence (OSINT) is intelligence gathered from publicly available sources. This can include both information gathered on-line as well as in the physical world. (Citation: SANSThreatProfile) (Citation: Infosec-osint) (Citation: isight-osint) |

# Method to the madness?