



Towards a Data-Driven Recommender System for Handling Ransomware and Similar Incidents

IEEE International Conference on Intelligence and Security Informatics

Martin Husák

husakm@ics.muni.cz

Czech CyberCrime Centre of Excellence (C4e)

Institute of Computer Science, Masaryk University, Brno, Czech Republic

November 2, 2021

Outline

Introduction

Requirements and Design

Calculations and Recommendations

Example of Using the System

Conclusion

Section 1

Introduction

Motivation

Ransomware and similar threat

- The rising **complexity** and **variety** of cyberattacks complicate **incident handling**.
- IDS and secure perimeter are bypassed by **social engineering attacks**, e.g., phishing.
- The malware further **spreads in the network**, exploiting surrounding computers.
- There is little chance of mitigating the spread of infection.

Incident handling

- Rapid incident response prevents spread of infection and reduces attack impact.
- Effective **triage and prioritization** of threats and incidents are of utmost importance.
- The behavior of malware can be **anticipated** to some extent.
- Social engineering is **difficult to detect** – we depend on **user reports**.

Approach

Anticipating the behavior of the malware

- A typical malware uses a few attack vectors and spreads in close proximity first.
- The **lateral movement** of an attacker can be observed, traced, and even projected.
- However, that requires detailed **knowledge of the local environment** and collaboration with users and administrators (complicated in large networks).

Recommender system for incident handling

- The incident handlers would appreciate any piece of information that would guide them through the network and pinpoint nodes that are immediately threatened.
- The key question of an incident handler is:
if this device is infected, which other devices can be infected or threatened?
- Recommendation of the list of devices at risk **can be automated**.

Section 2

Requirements and Design

Requirements

Data Collection

The system should collect or be able to access the data on the network and hosts in it. The required data items include a list of network segments, network topology, hosts' services, software, and vulnerabilities, and contacts on administrators.

Rich Information

The data shall hold as much information as possible. It is often unfeasible to have all the data available; the system shall work even with incomplete data.

Interconnection of Heterogeneous Data

The data shall be stored in a way that allows for the interconnection of heterogeneous data. The data should be accessible at any time and updated at least once a day.

Architecture

The proposed recommender system consists of three parts:

■ Data Collection

- Set of tools to monitor the network and the hosts¹
- Uses [NetFlow](#) network traffic monitoring and [Nmap](#) network scanner
- Timed execution to update the data regularly

■ Data Representation and Storage

- Data model to represent links in heterogeneous data
- Database to store the data

■ Recommendation Subsystem

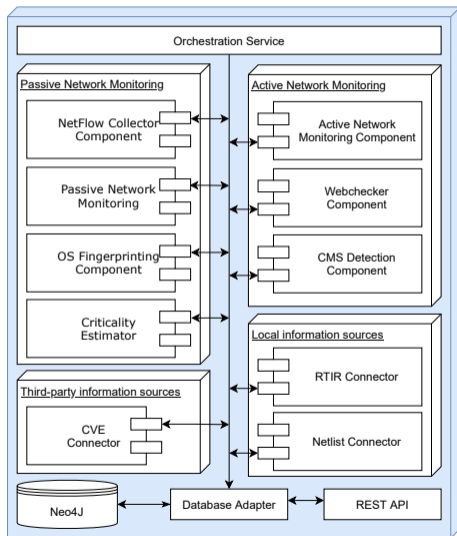
- Calculates similarity and distance of network nodes
- Prioritizes the hosts by their risk score

¹Husák, M., Laštovička, M., & Tovarňák, D. (2021, August). System for Continuous Collection of Contextual Information for Network Security Management and Incident Handling. In The 16th International Conference on Availability, Reliability and Security.

Data Collection

For the whole **network**, the system collects:

- List of the active **hosts** in the network,
- **Network topology** (via Nmap from several observation points),
- List of **network segments** with:
 - location (e.g., department, building, server room),
 - purpose (e.g., workstations, servers, IP pool of VPN, ...),
 - contact on responsible person (e.g., local IT administrator),
- History of **security incidents**:
 - including a list of involved devices,
 - available via RTIR or similar system.



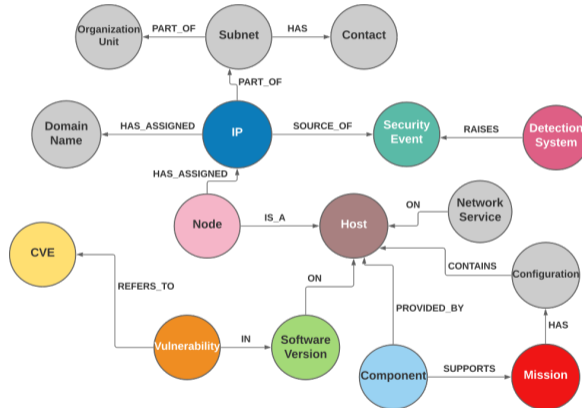
Data Collection

For each **host** in the network, the system collects the following:

- Fingerprint of the **operating system** (via NetFlow or Nmap),
- List of open **ports and services**, including the name and version of the underlying software (via NetFlow and NBAR2 signatures or Nmap),
- For web server: name and version of **Content Management System** (via WhatWeb),
- Name and version of a **web browser** used on the system (via NetFlow),
- Name of the **antivirus** software on the system and its latest update (via NetFlow),
- List of **vulnerabilities** (via vulnerability scanner or estimated from fingerprints),
- Location, purpose, contacts on administrators or main users (if available).

Data Representation and Storage

Data are stored in [Neo4j](#) graph database and structured using the [CRUSOE](#) data model²



² Komárková, J., Husák, M., Laštovička, M., & Tovarňák, D. (2018, August). CRUSOE: Data model for cyber situational awareness. In Proceedings of the 13th International Conference on Availability, Reliability and Security.

Recommendation Subsystem

The recommendation subsystem is a service that:

1. receives an identifier of a **host in the network** (e.g., IP address) on the input,
2. looks up the host in the database,
3. looks up devices in the **proximity** of the host,
4. calculates their **similarity** to the host on the input,
5. **prioritizes** the found hosts by their risk score,
6. returns a sorted list of **similar devices in close proximity** as the output.

More details in the following section.

Section 3

Calculations and Recommendations

General Idea

The recommendations are based on the **proximity** and **similarity** of the hosts in the network to the host on the input; similar hosts in close proximity are prioritized.

Proximity

Two hosts can be close to each other in physical and logical network topology, e.g., in the same room or in the same IP range. Alternatively, the two machines can be close to each other if they are controlled by the same users or administrators.

Similarity

The similarity is based on the similarity in software equipment, role, profile, or shared history of the two hosts. Similarity in software equipment is a prevalent feature due to the fact that the attackers typically exploit certain services or software.

Risk Score

- Formally, the hosts are sorted by their **risk score** (R) calculated as a quotient of the similarity (S) and distance (D) of the two hosts:

$$R = \frac{S}{D} = \frac{s_1 * s_2 * \dots * s_n}{\min\{d_1, d_2, \dots, d_n\}}$$

- In practice, it would be advantageous to assign **weights** to S and D or their elements.
- The weights could be extracted from real-world scenarios and tuned in operations.
- The weights are left for future work.

Distance Calculation

When the ransomware is reported, we do not know yet how it spreads:

- Malware spreading over the network will typically spread in the same subnet.
- Malware infecting files and drives will spread to machines used by the same user.
- Malware in email attachments will spread in the same department.

Distance

The distance between the two hosts is the minimal value of various distance metrics.

- **Breadth-first graph traversal** is used to find hosts with minimal distance in any of the distance metrics (in the implementation using graph database).
- The distance in logical network topology is the **length of the path** in the graph.
- Arbitrary distance metrics can be added as needed:
physical distance, location in the same room, similarity of IP addresses, ...

Similarity Calculation

- The malware often uses exploits of specific software or services.
 - If malware uses SSH brute-forcing, then Linux machines with SSH servers are at risk.
- We do not know the exact software equipment and may only assume similarities.
 - If the malware exploits Outlook email client, we shall look up all Windows machines.

Similarity

The similarity is calculated as a product of partial similarities $s_1 * s_2 * \dots * s_n$.
Each partial similarity is a value in the range $\langle 0, 1 \rangle$.

- The similarity of software equipment and network services are the main features.
- CPE strings represent pieces of software running on a host.

Similarity Calculation

Examples of similarity metrics

- CPE string similarity
 - CPE is an array of strings (vendor, product, version, ...) weighted 0.5, 0.25, 0.125, ...
 - The metric is the **sum of weights** of equal strings from the left to the first difference.
- CPE categories
 - If there is always **1 main CPE for each category**, then simple CPE similarity is used.
 - Categories can be OS, browser, antivirus, ...
- Service similarity
 - If a service is provided by one host but not the other, a **default value of 0.8** is used.
 - CPE strings are compared if both hosts provide the service.
- Similarities in vulnerabilities or past incidents
 - The number of **common CVEs** divided by the number of **unique CVEs** in the network.
 - The number of **common past incidents** divided by the **total number of past incidents**.

Section 4

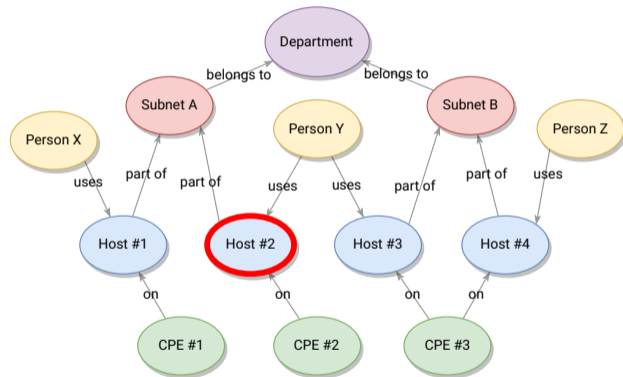
Example of Using the System

Distance calculation

Host #2 is reported to be infected, it's distance to other hosts is:

- 2 to Host #1 (same subnet)
- 2 to Host #3 (same user)
- 4 to Host #4 (subnets belonging to the same department)

Host #4 is too far – the calculation of similarity between Host #2 and Hosts #1 and #3 follows.



Similarity calculation

The OS fingerprint of Host #2 is compared to fingerprints of Hosts #1 and #3

- Host #1 and Host #2 share the same vendor, product, and version
the similarity is $0.5 + 0.25 + 0.125 = 0.875$
- Host #2 and Host #3 share only the vendor – their similarity is 0.5

CPE format	cpe:part:vendor:product:version:update:edition:language
Weights	0.5, 0.25, 0.125, 0.0625, 0.03125, 0.03125
CPE #1	cpe:2.3:o:microsoft:windows_7:-:sp2:*:*
CPE #2	cpe:2.3:o:microsoft:windows_7:-:sp1:*:*
CPE #3	cpe:2.3:o:microsoft:windows_10:-:*:*

The list of similar devices in close proximity for Host #2 goes as follows:

- Host #1, risk score is $0.875/2 = 0.4375$
- Host #3, risk score $0.5/2 = 0.25$

Section 5

Conclusion

Conclusion

Summary

- We proposed a design of a **recommended system for incident handling**.
- If a compromise of a host in the network is reported, the system instantly provides a **prioritized list** of other **hosts at risk** for rapid response.
- The system considers various ways of attack propagation or lateral movement.

Future work

- This paper is merely the first step in the future research.
- The system is **under development** and will be **evaluated in operations**.
- The weights of the metrics will be inferred from **past incidents**.
- Integration with other incident handling tools will follow.

MUNI
C4E



EUROPEAN UNION
European Structural and Investment Funds
Operational Programme Research,
Development and Education

MŠMT
MINISTRY OF EDUCATION,
YOUTH AND SPORTS

C4E.CZ