

Article

Not peer-reviewed version

---

# Using a Graph Engine to Visualize the Reconnaissance Tactic of the MITRE ATT&CK Framework from UWF-ZeekData22

---

[Sikha S Bagui](#)\*, [Dustin Mink](#), [Subhash C Bagui](#), Michael Plain, [Jadarius Hill](#), Marshall Elam

Posted Date: 9 June 2023

doi: 10.20944/preprints202306.0696.v1

Keywords: Graph databases; Data Visualization; MITRE ATT&CK Tactics; Star Motif; Clique Motif; Reconnaissance Tactic



Preprints.org is a free multidiscipline platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This is an open access article distributed under the Creative Commons Attribution License which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

*Article*

# Using a Graph Engine to Visualize the Reconnaissance Tactic of the MITRE ATT&CK Framework from UWF-ZeekData22

Sikha S. Bagui <sup>1,\*</sup>, Dustin Mink <sup>1</sup>, Subhash C. Bagui <sup>2</sup>, Michael Plain <sup>1</sup>, Jadarius Hill <sup>1</sup> and Marshall Elam <sup>1</sup>

<sup>1</sup> Department of Computer Science, University of West Florida, Pensacola, FL 32514, US

<sup>2</sup> Department of Mathematics and Statistics, University of West Florida, Pensacola, FL 32514, US

\* Correspondence: bagui@uwf.edu

**Abstract:** There has been a great deal of research in the area of using graph engines and graph databases to model network traffic and network attacks, but the novelty of this research lies in visually or graphically representing the Reconnaissance Tactic (TA0043) of the MITRE ATT&CK framework. Using the newly created dataset, UWF-Zeekdata22, based on the MITRE ATT&CK framework, patterns involving network connectivity, connection duration, and data volume were found and loaded into a graph environment. Patterns were also found in the graphed data that match the Reconnaissance as well as other tactics captured by UWF-Zeekdata22. The Star motif was particularly useful in mapping the Reconnaissance tactic. The results of this paper show that graph databases/graph engines can be essential tools for understanding network traffic and trying to detect network intrusions before they happen. Finally, an analysis of the run-time performance of the reduced dataset used to create the graph databases showed that the reduced datasets performed better than the full dataset.

**Keywords:** Graph databases; Data Visualization; MITRE ATT&CK Tactics; Star Motif; Clique Motif; Reconnaissance Tactic

## 1. Introduction

In the past decade, the number of IoT (Internet of Things) devices connected to the internet has significantly increased. It is expected that 43 billion IoT devices will be connected by the end of 2023 [1]. As the number of connected devices grows, so will network traffic and the amount of data transmitted. Because IoT devices are used in industries that use sensitive data, for example, health care and the financial sector, not only it is imperative that the data maintains its integrity and is uncompromised during transit and at rest, but it is also important that we try to prevent network attacks before they happen. To do this properly, not only do we need to possess the ability to distinguish between regular network traffic and attack traffic, but we also need to possess the ability to detect attacks before they happen.

Many studies have been performed on identifying attack network traffic after the attacks have happened [2–5], but in this work we are trying to study the step before that, that is, who is trying to gather information about our system so that they can perform an attack. Hence, our aim in this work is to analyze the Reconnaissance Tactic (TA0043) of the MITRE ATT&CK framework. The Reconnaissance tactic of the MITRE ATT&CK framework is used to gather information about vulnerabilities in a system [6], mostly by active scanning. Understanding the nature of reconnaissance being performed in a system is very important to be able to prevent future attacks before they happen. In this work we use a graph engine or graph database to present visual representations of the Reconnaissance tactic. Though the focus is on the Reconnaissance tactic, we also present visual representations of regular network traffic and other attack traffic labeled as per the MITRE ATT&CK framework.

Graph databases by definition are no-SQL databases based on a network structure and are based on mathematical graph theory. Graphs are composed of 3 different types of objects: vertices, edges,

and properties. Vertices, or points, are used to represent entities of data that correspond to some object. Edges, or lines, represent relationships between various vertices; these connections may be unidirectional or bidirectional [7]. Properties are attributes of the objects. In this work, vertices correspond to different machine IPs that are communicating, edges represent the connections between different machines, and properties are different attributes that correspond to the edges such as connection duration.

Graphs and graph databases can be utilized to generate graph models to represent relationships. In addition to visualizations representing attack/non-attack data, graph data models can be extremely useful, especially in cybersecurity, because these models can be utilized for pattern recognition, machine learning, and other analysis. Graph databases can be used to generate predictions to distinguish between regular network traffic patterns and attack patterns [8].

Though there has been a great deal of research in the area of using graph engines and graph databases to model network traffic and network attacks, the novelty of this research lies in visually or graphically representing the Reconnaissance Tactic (TA0043) of the MITRE ATT&CK framework. Using the newly created dataset, UWF-ZeekData22 [9,10], labeled based on the MITRE ATT&CK framework, patterns involving network connectivity, connection duration, and data volume were found from the Conn Log files of UWF-ZeekData22 [9,10], and loaded into a graph environment. Hence, to elaborate on the novelty of this research, it can be stated that:

- To date, tactics from the MITRE ATT&CK framework have not been visualized graphically. This work focuses on presenting graphic visualizations of the MITRE ATT&CK Reconnaissance Tactic (TA0043) using graph representation.
- Essential feature selection is performed so that this work generates a graph data model using only a very limited set of network connection features. Feature generation was also performed using the limited set of network connection features.

Though this is beyond the scope of this work, the benefits of this graphical representation can be realized as follows in the future:

- The graph models could be effectively used to train machine learning models, especially in the Big Data environment, in order to accurately predict when network traffic is nefarious.
- The reduction of the network data to only a few features (feature selection) that could be used to identify a Reconnaissance tactic would be computationally beneficial in machine learning analysis, especially in the Big Data environment.
- And above all, these graph models can be used to develop a more robust Threat Intelligence Platform (TIP) that would be able to visually detect the attacks before they happen, by recognizing the attack patterns in the data. A TIP is a technology solution that collects, aggregates and organizes threat intelligence.

Finally, in this work, an analysis is done of the runtime performance of creating the graph representations with the reduced set of data.

The rest of this paper is organized as follows. Section 2 presents previous works related to graph databases; section 3 presents the dataset and the software used to process the data; section 4 presents the pre-processing that was used on this dataset; section 5 presents the algorithmic approach to creating the graphs; section 6 presents data visualizations using graph databases; section 7 presents the runtime performance for creating the graph databases; section 8 presents the conclusions and section 9 presents the future works.

## 2. Related Works

Utilizing graphs to represent network connectivity, for the purpose of identifying anomalies has been the topic of many research articles [7,11,12,13,14]. Interpretation of the graph data to detect anomalies has been a challenging task in relation to summarizing normal data while retaining enough information to detect anomalies [12]. Identifying motifs and comparing multiple graphs for similarity using various motifs becomes challenging as graph sizes increase [11]. A named entity recognizer (NER) was proposed by one group of authors, allowing for the training of an extractor to obtain useful information from the MITRE ATT&CK framework. A multi-step approach to building a

knowledge base included collection and analysis, construction of an ontology from the information gathered, and finally, generation of a cybersecurity knowledge deduction engine [7]. Another group of researchers approached the problem by an abstracted graph approach, where flexible attack profiles were created and used to detect simulated attacks. Utilizing a graph database, the team proposed the possibility of not only identifying the attacker but also the possibility of detecting other impacted system components [13]. Finally, an approach was proposed to compare similarities between graphs using a novel neural network approach. Important vertices would be identified by a specific similar metric and a pairwise vertex comparison would be utilized to identify similarity. The group concluded that the first steps were made at bridging the gap between graph deep learning and the graph search problem [14].

In this paper, the idea is to get away from solely using edges. This paper presents the network hops between source and destination which resulted in an attack in the MITRE ATT&CK framework. The paper also demonstrates the successful utilization of motifs to visually identify behavior patterns representing an attack tactic. And finally, an analysis is performed of the runtime performance of creating the graph representations and databases with the reduced set of data.

### 3. The Dataset: UWF-ZeekData22

Since graph data models depend on the connections between data points, the Conn log files of the UWF-ZeekData22 [9,10] dataset were used for generating the graphs. UWF-ZeekData22 [9,10] was generated by the Cyberrange group associated with the University of West Florida and the full data set is available at [10]. This dataset has 9,280,869 attack records and 9,281,599 benign records with a total of 18,562,468 records.

The data schema of the Conn log files is presented in Table 1. To generate the graphs, only four fields from the Conn Log files were used in addition to count: id.orig\_h (the source IP, referred to as srcIP in this paper), id.resp\_h (the destination IP, referred to as dstIP in this paper), duration, and orig\_bytes (referred to as bytes).

**Table 1.** UWF-ZeekData22: Schema of the Conn Log files [9, 10].

Attribute Name	Description of Attribute	Used to Create Graph DB
ts	Time of first packet	
uid	Unique identifier of connection	
id.orig_h	IP address of packet sender	Yes
id.orig_p	Outgoing port number	
id.resp_h	IP address of packet receiver	Yes
id.resp_p	Incoming port number	
proto	Transport layer protocol of connection	
service	Application protocol sent over connection	
duration	How long connection lasted	Yes
orig_bytes	Payload bytes originator sent	Yes
resp_bytes	Payload bytes responder sent	
conn_state	Possible connection states	
local_orig	If connection is originated locally	
local_resp	If connection is responded to locally	
missed_bytes	Representative of packet loss	
history	History of connections	
orig_pkts	Number of packets originator sent	
orig_ip_bytes	Number of IP level bytes originator sent	
resp_pkts	Number of packets responder sent	
resp_ip_bytes	Number of IP level bytes responder sent	
community_id		

id	Connection's 4-tuple of endpoint addresses/ports	
tunnel_parents	*uid* values for encapsulating parent(s) connections used	

### 3.1. Distribution of UWF-ZeekData22 by Tactics

Table 2 presents tactics available in UWF-ZeekData22. For this analysis, initially, the data was divided into four categories by attack tactic: Reconnaissance, Discovery, No Attack, and all attack tactics. Reconnaissance and Discovery were selected since they had more data. No Attack was selected to visualize how a normal network traffic would appear without abnormal traffic included. The All Attack Tactics dataset was selected to visualize how normal and abnormal network traffic would appear. Since the volume of data for Discovery was eventually not considered enough for a robust analysis, this category was also not further analyzed in this work. Hence finally a full analysis is presented of only the Reconnaissance tactic, non-attack data and all data (which also includes the Reconnaissance and Discovery). The other categories were also not analyzed individually due to the minimal amounts of occurrences of the other tactics.

**Table 2.** UWF-ZeekData22 Tactics [10]

Attack Tactic	Count
None (Not an attack)	9,281,599
Reconnaissance	9,278,722
Discovery	2,086
Credential Access	31
Privilege Escalation	13
Exfiltration	7
Lateral Movement	4
Resource Development	3
Defense Evasion	1
Initial Access	1
Persistence	1

### 3.2. Software Utilized to Process Data

Python and pySpark were utilized as GraphFrames is readily available in this environment. In order to visualize the graph data, GraphStreams [15] was used since it has a feature-rich library. GraphStreams [15] was implemented in the Java environment.

## 4. Preprocessing

Using the Conn dataset from UWF-Zeekdata22[9,10], a unique list of source and destination IP addresses were generated using a simple hashmap. A graph was created using the unique list as the graph vertices, naming the vertices based on whether they were a source IP or destination IP. Once the graph vertices were created, edges were established and weighted based on the following dominant attributes:

- Destination ip (id.resp\_h) and originating bytes (orig\_bytes), used as per [16].
- Total number of connections between the unique source and destination
- Total duration of the connection(s) between the vertices
- Total number of bytes of the connections between vertices
- The attack tactic

First, this information was used to generate a pySpark vertex and edge list. Then, this information was used to create a graphFrame in order to determine vertex and edge relationships and graph shapes. The objective was to look for two primary structures in the graphs, star motifs and clique motifs. Star motifs are where a single vertex connects to multiple vertices and clique motifs are



where the largest set of interconnected vertices is identified. Stars in a graph are defined as having  $n-1$  vertices with a degree of 1 and a single vertex having a degree of  $n-1$  [17]. The Bron-Kerbosch algorithm [18] was utilized to find maximal cliques. This algorithm finds the largest connected vertices that produce the unique clique.

Additional effort was taken to scan the vertices and edges to find and eliminate intermediate vertices, revealing true endpoints in the graph. In order to do this, cycles had to be identified and eliminated. The approach taken initially was to use Depth-First-Search (DFS), but due to the number of vertices in the graph, a dynamic algorithmic approach was taken to minimize recursive code. The dataset was reduced to tables of unique source and destination addresses and accumulated connections, durations, and bytes transmitted. These vertices were then used to construct a graph, eliminating any edges that result in a cycle. Eliminating cycles provided for a minimally connected graph which was easier and faster to traverse when connecting the source of an attack to its destination. Elimination of the cycles did not impact the underlying graph as all vertices were still reachable by other adjacent vertices [19]. Elimination of the cycles reduced the edges needed to create the graph and thus produced a more concise graph. This allowed for identifying motifs of interest as they stood out from the background of random interconnections that were not of interest [20].

#### 4.1 Binning Methodology

Binning allowed for continuous data to be represented in various discrete categories or bins. In order to best characterize the data, the following attributes of the edge connections were binned: number of connections, average duration, and average bytes. In order to bin the data, the methodology outlined by the authors of [16] was utilized, however, a stationary mean was implemented instead of a moving mean. The standard deviation was first calculated by using the formula:

$$stddev = \sqrt{\sum_i^n \frac{(x_i - \bar{x})^2}{n}} \quad (1)$$

where  $x$  is the attribute that is being binned,  $\bar{x}$  is the average of the attribute, and  $n$  is the number of data points. Six bins were then constructed using the calculated standard deviation as follows:

$$bin_1 = (-\infty, \bar{x} - stddev * 2) \quad (2)$$

$$bin_2 = [\bar{x} - stddev * 2, \bar{x} - stddev) \quad (3)$$

$$bin_3 = [\bar{x} - stddev, \bar{x}) \quad (4)$$

$$bin_4 = [\bar{x}, \bar{x} + stddev) \quad (5)$$

$$bin_5 = [\bar{x} + stddev, \bar{x} + stddev * 2) \quad (6)$$

$$bin_6 = [\bar{x} + stddev * 2, \infty) \quad (7)$$

Each of the three edge attributes was assigned a bin determined by which bin the attribute's value landed in. Because the data had a large variance and thus a large deviation, the first two bins were negative for some of the attributes.

After using equation (1) to calculate the standard deviation for the count attribute for the full Reconnaissance dataset, equations (2) – (7) were used to calculate the bins for the count attribute as follows:

$$stddev = 265048.551, \quad \bar{x} = 16963.973$$

$$bin_1 = (-\infty, \bar{x} - stddev * 2) = (-\infty, 16963.973 - 2 * (265048.551))$$

$$bin_2 = [\bar{x} - stddev * 2, \bar{x} - stddev) = [16963.973 - 2 * 265048.551, 265048.551) = [-513133.129, -248084.578)$$

$$bin_3 = [\bar{x} - stddev, \bar{x}) = [16963.973 - 265048.551, 16963.973) = [-248084.578, 16963.973)$$

$$bin_4 = [\bar{x}, \bar{x} + stddev) = [16963.973, 16963.973 + 265048.551) = [16963.973, 282012.524)$$

$$bin_5 = [\bar{x} + stddev, \bar{x} + stddev * 2) = [16963.973 + 265048.551, 16963.973 + 2 * 265048.551) = [282012.524, 547061.074)$$

$$bin_6 = [\bar{x} + stddev * 2, \infty) = [16963.973 + 2 * 265048.551, \infty) = [547061.074, \infty)$$

To find which bin a value is in, the bin that overlaps the value is found. As an example, the value 1280 is between the values -248084.578 and 16963.973; therefore, the value resides in  $bin_3$ .

## 5. Algorithmic Approach to Creating the Graphs

### 5.1 Overview of Approach

UWF-ZeekData22 [9,10] was reduced to the source and destination IPs only, by removing intermediary vertices and cycles in an effort to remove network noise. To remove the intermediary vertices, a Depth First Search (DFS) algorithm approach was taken, adding only edges that did not result in a cyclic graph. Due to the number of vertices in the graph, a dynamic algorithmic approach was taken to minimize recursive code. The dataset was reduced to tables of unique source and destination addresses and accumulated connections, durations, and bytes transmitted. These vertices were then used to construct graphs, eliminating any edges resulting in cycles. Graphical representations are presented of the Reconnaissance Tactic, as well as all attack and non-attack traffic.

### 5.2 Workflow

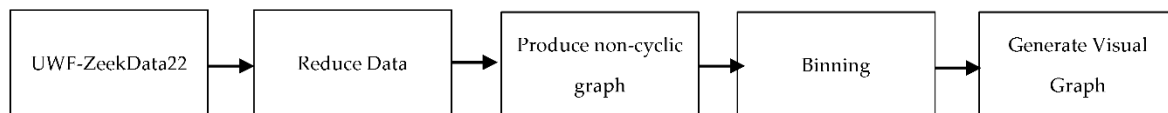


Figure 1. Workflow

#### 5.2.1 Reduce the Data

Since UWF-ZeekData22 [9,10] is a large dataset, one of the first objectives was to see if any kind of feature reduction could be applied. Hence, only the connection counts, bytes transferred, and connection data were aggregated to reduce the number of data points that would feed into the next graphing phase. Specifically, the duration and orig\_bytes features from the Conn Log files of UWF-ZeekData22 [9,10] were aggregated by the unique source to destination key. These features were totaled and, additionally, new features were generated using duration and orig\_bytes. The additional new features were average duration and average bytes.

#### 5.2.2 Produce a non-cyclic graph

Graphs were created using the IP addresses obtained in the previous phase, populating the edges with the aggregated counts, bytes, and duration values. As each edge was added to the graph, a check was performed to determine if the new edge produced a cycle. If a cycle was created, the edge was removed from the graph. The final graph data was then written out as a CSV file for the next phase.

### 5.2.3 Binning

The CSV file from the previous phase was analyzed and binned as explained in the pre-processing section. The resulting bins replaced the original graph data and a new CSV file was produced for the next phase.

### 5.2.4 Generate Visual Graph

The resulting graph data, now binned on count, bytes, and duration, was loaded into the GraphStream application and visualization of the graphs was produced and used in this work.

## 5.3 Algorithmic Approach to Creating the Graphs

Each unique source to destination edge was identified and mapped. With each unique edge between the source and destination, a summation of attributes that were to be tracked was stored. A graph  $G$ , of unique vertices, was created. Iterating through all source vertices, an edge was added to the graph, from source to destination, and tested for the creation of a cycle in the graph. If a cycle was detected, then the last edge was removed. The final resulting graph produced the longest path between a given source vertex and its furthest destination vertex, which did not result in a cycle. This allowed for the elimination of intermediate vertices and the detection of the final destination of an attack from a source.

If calling isCyclic method (Algorithm 1) for the Graph results in true, then a cycle has been encountered and the last vertex must be removed to remove the cycle. Analysis was done to determine if any meaningful correlation could be attributed to the attack tactics port numbers used by the source or destination. It was found that this information did not add any value to the graph and therefore port was eliminated as a possible attribute of interest.

#### Algorithm 1: isCyclic

**Input:** Graph  $G$ , vertex  $V$  to add

**Output:** Boolean true if after adding  $V$ , the graph is cyclic,  
updated  $G$ , with vertex  $V$  added

Add  $V$  to  $G$

Create and initialize visited array, recursionStack array

Mark all vertices as unvisited in both visited and recursionStack

forall vertex  $v$  in  $G$

Return isCyclicUtil( $v$ , visited, recursionStack)

isCyclicUtil(vertex, visited array, recursionStack)

if vertex visited before return false

if vertex is in recursionStack return true

Mark vertex as visited for vertex

Mark recursionStack as visited for vertex

forall children of vertex

if isCyclic(childVertex, visited array, recursionStack)

Return true

Set recursionStack for vertex to false

Return false

## 6. Resulting Graph Visualizations



GraphStream [15] was utilized to generate graphical visuals for each of the subsets of the edges. GraphStream is a Java library used for modeling, visualizing, and analyzing dynamic networks of various sizes [15].

The data was fitted to different motif models to determine if various attacks could be characterized by specific shapes. In the motifs (Figures 2-8) that follow, the color of each edge represents the intensity/bin of the corresponding attribute that the graph represents. The colors, orange for bin 1, yellow for bin 2, green for bin 3, blue for bin 4, purple for bin 5, and red for bin 6, were used in order of least to highest intensity to represent the bin value ranges.

### 6.1 Star Motif

The Reconnaissance tactic resembles the star motif, in which there is a central vertex which the connections originate from. As seen in Figure 2, all connections originate from the central vertex of 143.88.2.10. This indicates active scanning [21], typical of a Reconnaissance tactic. In active scanning, an adversary probes a victim infrastructure's network traffic by mechanisms such as port scanning. Port scanning classifies each port into the state of open, closed, filtered, unfiltered, open/ filtered, or closed/ filtered [22]. This helps an attacker determine which ports on a network are open and can be utilized to receive and send data. Figures 2, 3 and 4 represent the Reconnaissance motif by connection count, average duration, and average bytes respectively.

#### 6.1.1 Visualizing the Reconnaissance Tactic by Connection Count

Figure 2 depicts the Reconnaissance tactic radiating from a single vertex, 143.88.2.10, to multiple other vertices in the graph. The number of connections from point to point is generally in the average range of connections with the exception of a few which were in the extreme range of binning. Looking deeper into the data, it can be seen that each connection generally involves a different port, therefore this graph is representative of a port scan, typical of a Reconnaissance tactic. This graph had some areas of interest, represented by the red connections (bin = 6), where considerably more connections occurred than the normal connection count (bin = 3) which was 1,024 connections. Each of these bin6 connections was in excess of 1 million. One outlier in the data was a connection between 143.88.5.12 and 143.88.5.1 (bin = 5) with ½ million connections. Example data points can be seen in Table 3. For the Reconnaissance tactic, the maximum connection count was 3,112,192, while the minimum connection count was zero, and the average connection count was 33,927.946.

It can also be noted from Figure 2 that 143.88.2.10 is mostly pointing to 143.88.7.\* addresses. The graph is actually pointing to the entire range of the subnet which is from 143.88.7.0-255. The red lines indicate where most of the bytes are being transmitted back and forth. This is highly likely because the 4 IP addresses belonged to running virtual machines on the victim's network, and a reply from the victim's network is indicative of an open port of a victim's host.

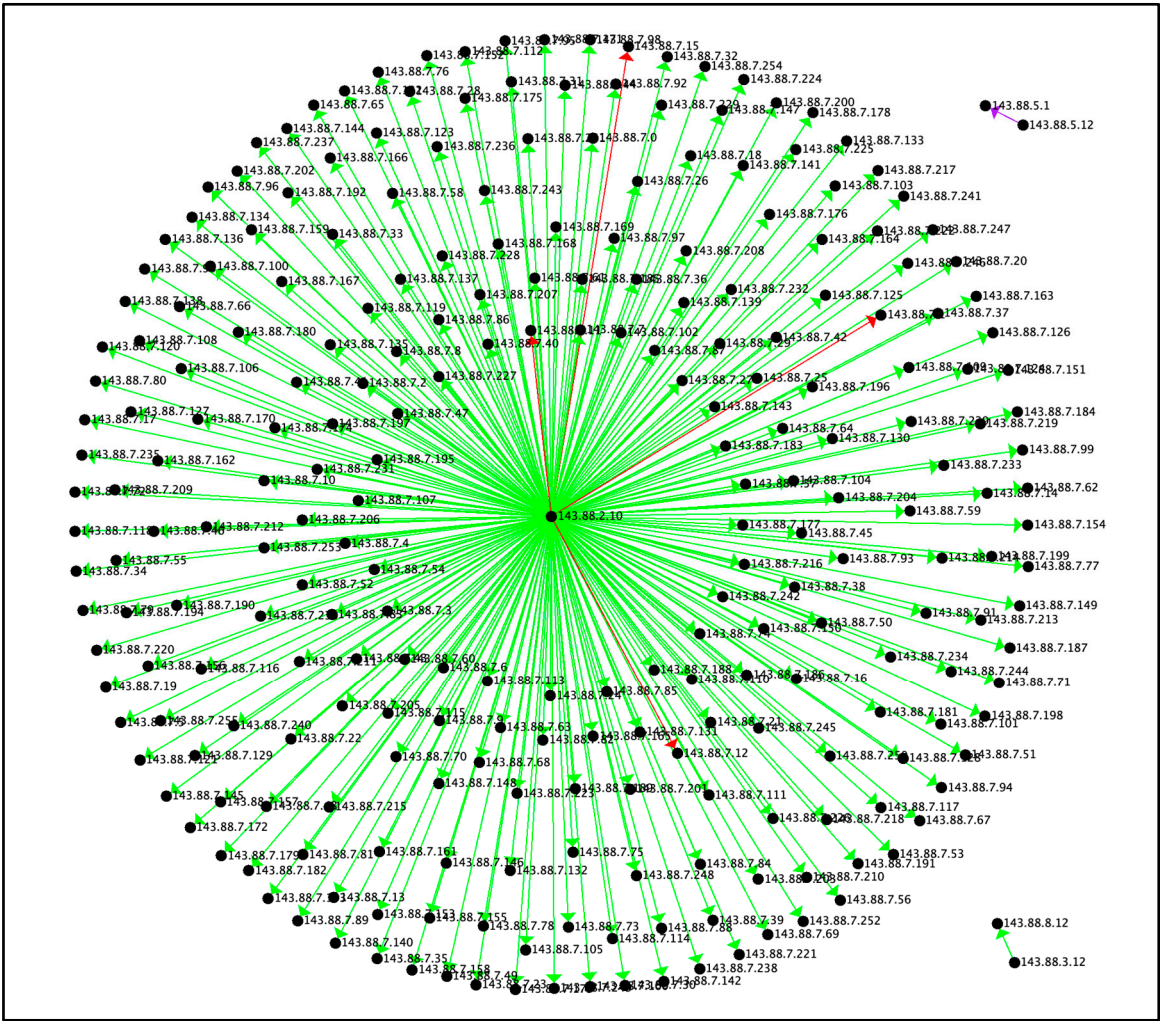


Figure 2. Reconnaissance Tactic by Connection Count

Table 3. Reconnaissance Points of Interest (Count)

Id	From	To	Total_Dur	Avg_Du r	Total_Bytes	Avg_Bytes	Count	CountBi n
edge_0	143.88.2.1 0	143.88.7.1 5	353248.5154	0.2126	265458232832 0	1597759.972 2	166144 0	6
edge_1	143.88.2.1 0	143.88.7.1 1	972063.5371	0.3123	5579520	1.7928	311219 2	6
edge_2	143.88.2.1 0	143.88.7.1	279987.9888	0.1338	8567808	4.0934	209305 6	6
edge_3	143.88.2.1 0	143.88.7.1 2	778386.2988	0.6914	925758636800	822247.5387	112588 8	6
edge_25 7	143.88.5.1 2	143.88.5.1	943576.7243	1.8777	36458752	72.5507	502528	5

6.1.2 Visualizing the Reconnaissance Tactic By Average Duration

Figure 3 presents the Reconnaissance Tactic by average duration. The average duration of the connections in the star motif did not identify areas of interest as green (bin=3) and blue (bin=4) are

average behaviors in this graph. The blue connections in Figure 3 correspond to the high connections found in Figure 2, although the duration per connection is considerably higher, ranging from 300 to 1700 times longer than the other connections in green. The connections in green transferred 0 bytes whereas the connections in blue transferred data from between 2 bytes to 1.5 MB of data per connection. Sample data points for Reconnaissance points of interest based on average duration are presented in Table 4. The maximum duration was 972,063.54, minimum duration was 0.04, and average duration was 12,947.3263.

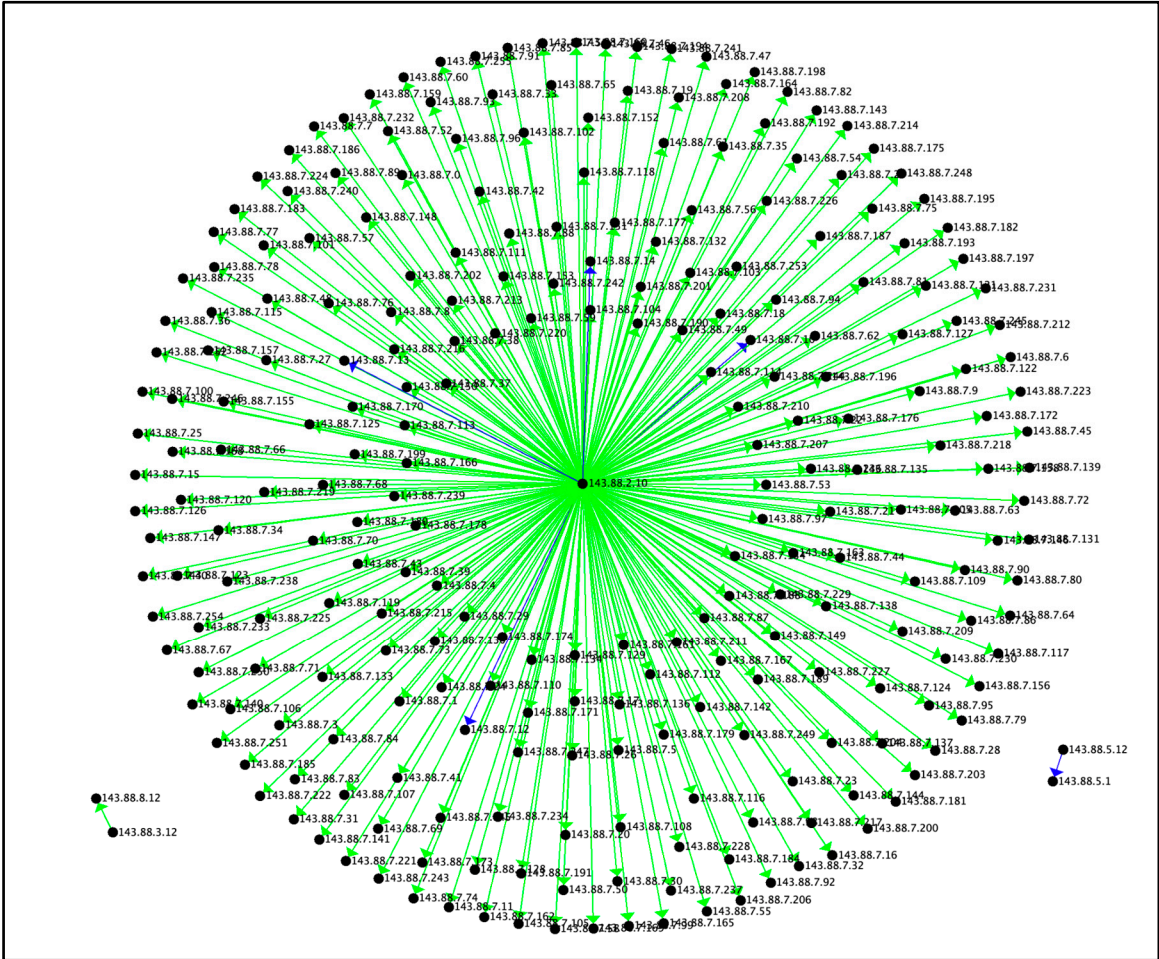


Figure 3. Reconnaissance Tactic by Average Duration

Table 3. Reconnaissance Points of Interest (Count)

Id	From	To	Total_Du r	Avg_D ur	Total_Bytes	Avg_Bytes	Count	CountBi n
edge_0	143.88.2. 10	143.88.7. 15	353248.51 54	0.2126	26545823283 20	1597759.97 22	166144 0	6
edge_1	143.88.2. 10	143.88.7. 11	972063.53 71	0.3123	5579520	1.7928	311219 2	6
edge_2	143.88.2. 10	143.88.7. 1	279987.98 88	0.1338	8567808	4.0934	209305 6	6
edge_3	143.88.2. 10	143.88.7. 12	778386.29 88	0.6914	92575863680 0	822247.538 7	112588 8	6



edge_2 57	143.88.5. 12	143.88.5. 1	943576.72 43	1.8777	36458752	72.5507	502528	5
--------------	-----------------	----------------	-----------------	--------	----------	---------	--------	---

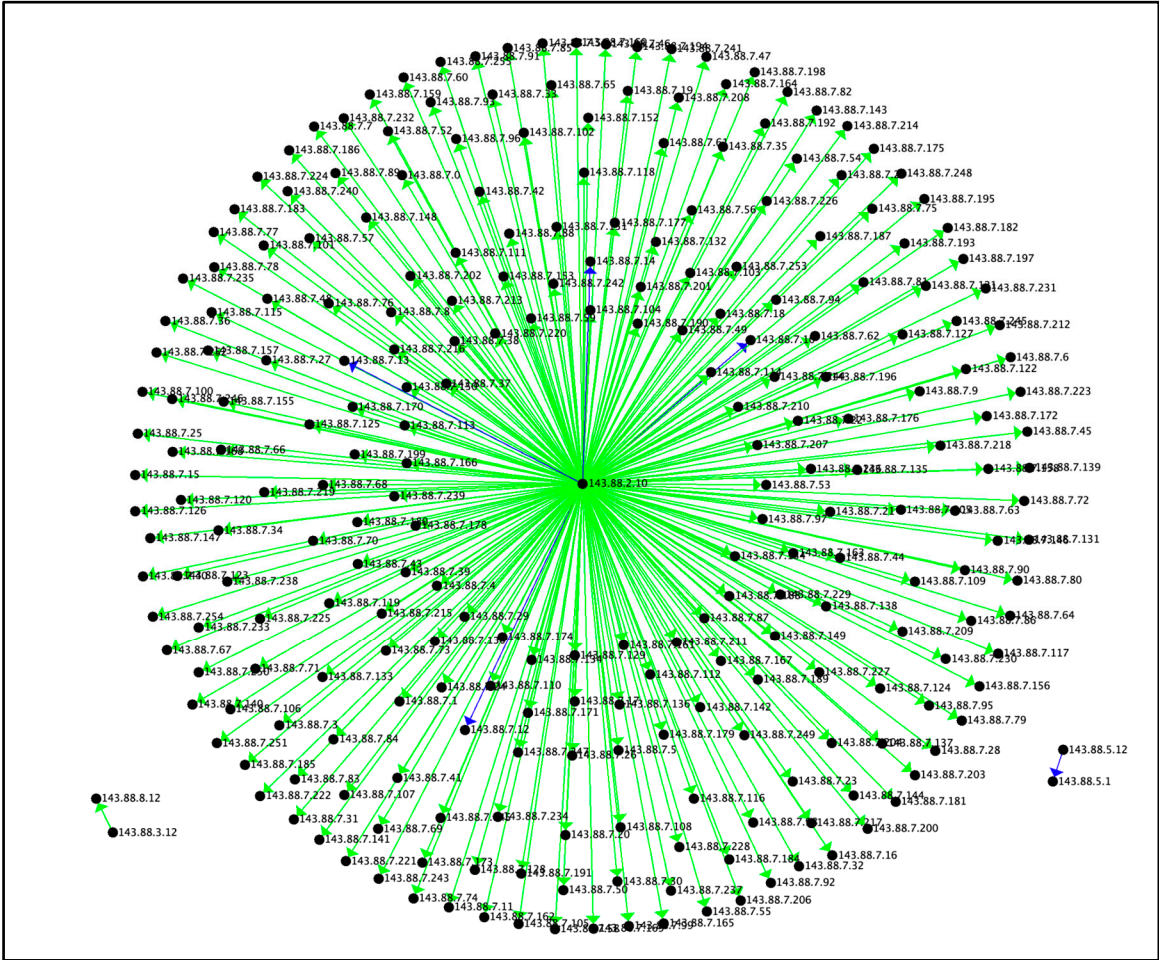


Figure 3. Reconnaissance Tactic by Average Duration

Table 4. Reconnaissance Points of Interest (Average Duration)

Id	From	To	Total_Du r	Avg_D ur	Total_Bytes	Avg_Byte s	Count	CountBi n
edge_3	143.88.2.10	143.88.7.1 2	778386.29 88	0.6913	92575863680 0	822247.538 7	112588 8	4
edge_4	143.88.2.10	143.88.7.1 0	1792.9392 7	1.4007	798720	624	1280	4
edge_4 2	143.88.2.10	143.88.7.1 4	3080.24	3.0080	0	0	1024	4
edge_4 3	143.88.2.10	143.88.7.1 3	3080.264	3.0080	0	0	1024	4
edge_2 57	143.88.5.12	143.88.5.1	943576.7	1.8776	36458752	72.5506877 2	502528	4

6.1.3 Visualizing the Reconnaissance Tactic by Average Bytes

Figure 4 presents the Reconnaissance tactic by average bytes. As depicted in Figure 4, only two areas of interest were identified. In both cases, the number of bytes transferred, per connection was 0.8 MB to 1.5 MB. It is possible that the attacker found that these IP addresses had exposed ports and thus was available to use them to send and/or receive data to/from the network. Example data points for the Reconnaissance points of interest based on average bytes are presented in Table 5. The maximum bytes transferred were 2,654,582,328,320, minimum bytes transferred were zero and the average bytes transferred were 13,877,478,833.

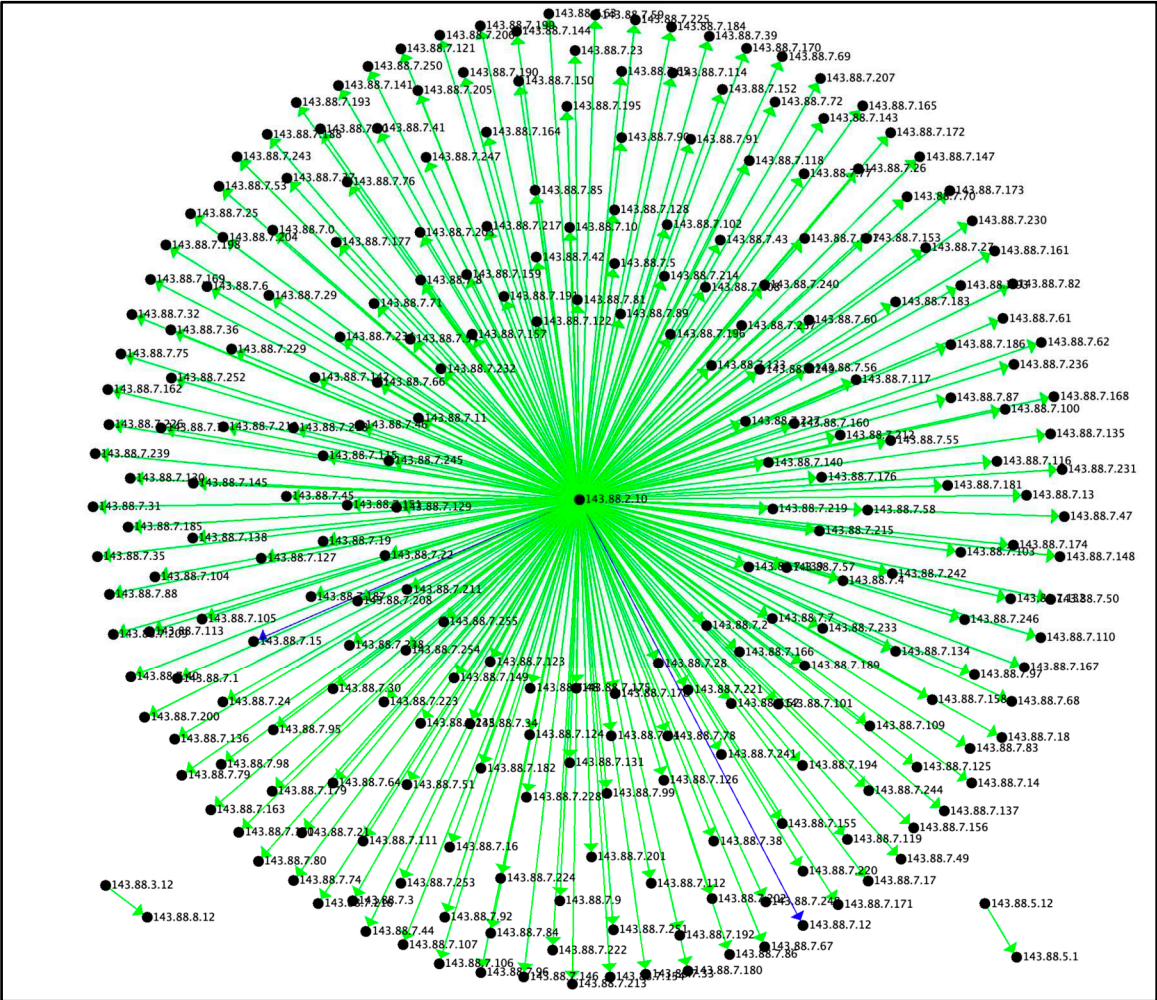


Figure 4. Reconnaissance Tactic by Average Bytes

Table 5. Reconnaissance Points of Interest (Average Bytes)

Id	From	To	Total_Dur	Avg_Dur	Total_Bytes	Avg_Bytes	Count	CountB in
edge_0	143.88.2.10	143.88.7.15	353248.5	0.212616	26545823283	1597760	166144	3
edge_3	143.88.2.10	143.88.7.12	778386.3	0.691353	92575863680	822247.5	112588	4
edge_257	143.88.5.12	143.88.5.1	943576.7	1.87766	36458752	72.55069	502528	4

6.2 Clique Motif

Figure 5 depicts the cliques found in UWF-ZeekData22. The bottom left set of IP addresses are reverse shells coming back to the 143.88.2.10 address, which are attackers on the kali linux machine used to scan and attack the victim’s network. The connections in the red box are interesting because they are able to gain a connection to the University of West Florida’s (UWF’s) IP address which is the 143.88.0.\* subnet. The group of connections in the top right are IPv6 addresses. The IPv6 address is the successor of the regular IPv4 address [23]. With the limited number of IPv4 addresses, in order to accommodate for the increasing number of devices on the internet, the Internet Engineering Task Force (IETF) developed Internet Protocol version 6 (IPv6) address. IPv6 uses a 128-bit compared to IPv4, which uses a 32-bit address.

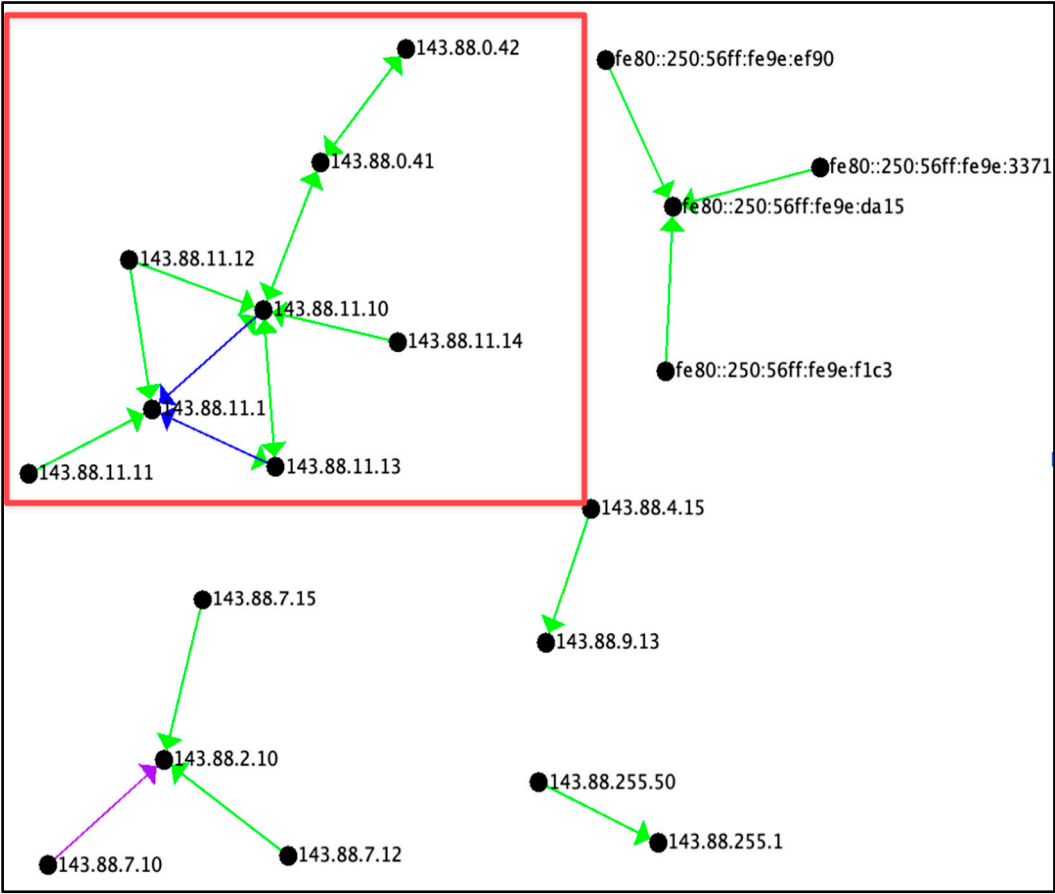


Figure 5. Maximal Cliques Found by Connection Count for UWF-ZeekData22

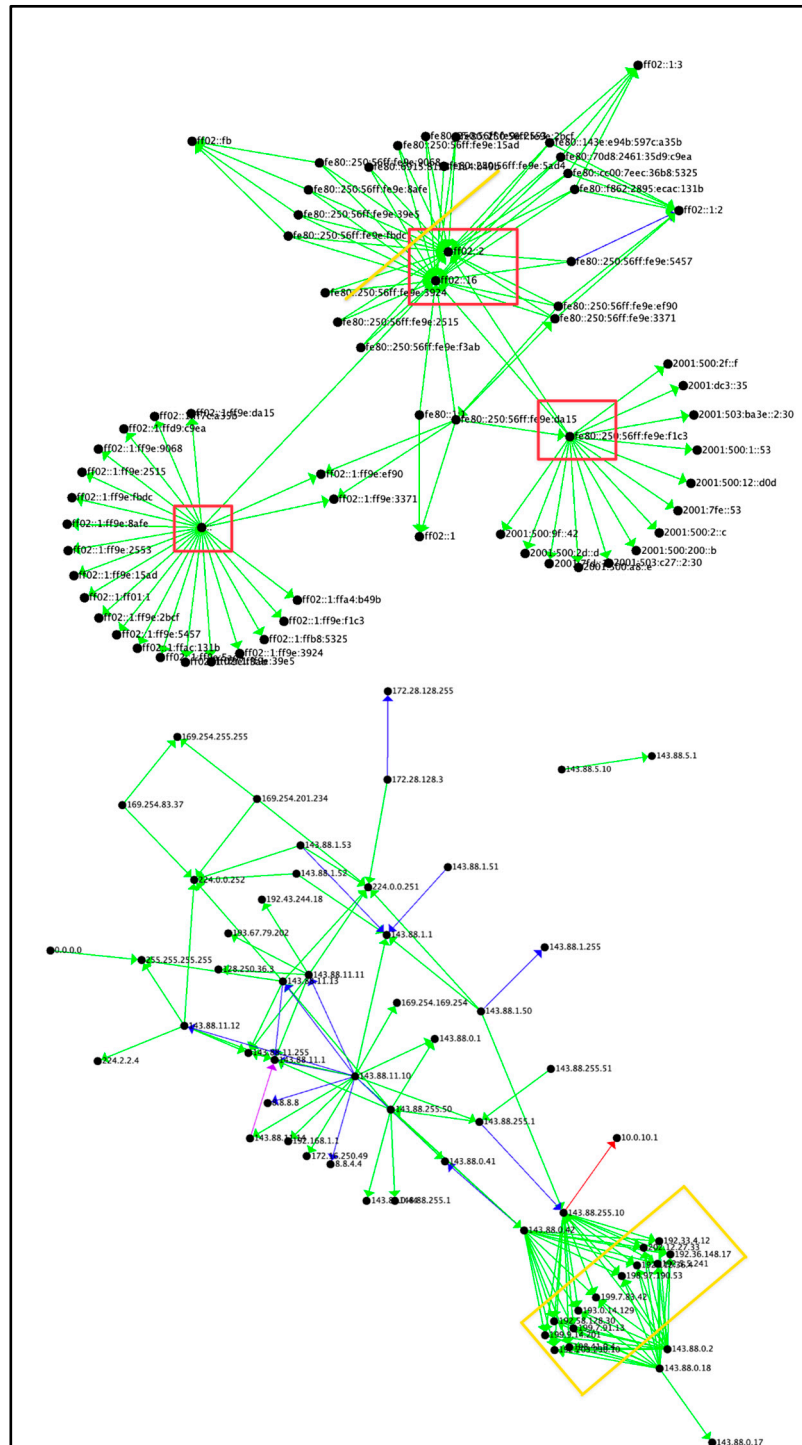
6.3 Visualizations of Non-Attacks By Count

Figure 6 depicts the count of connections that were categorized as non-attacks, and shows a large cluster of different connections of IPv6 addresses. There are several areas of interest identified by the colored boxes. The IP addresses within the red boxes are routers or switches that are redirecting traffic to different subnets, ff02::fb and ff02::1:3. And these subnets are possibly redirecting it to servers or load balancers.

As cycles were removed from the data, they appeared unidirectional. The yellow boxed area (bottom right) represents servers that were behind a load balancer. The load balancer evenly distributes traffic to the various servers.

Two data points for the non-attack by connection count are presented in Table 6. The maximum count was 6,724,017, minimum count was 1 and average count was 4,273,817.





**Figure 6. Non-Attack by Connection Count**

## 6.4 Visualizing Attacks By Count

Figure 7 depicts the full picture of the attack data binned with respect to the number of occurrences (Count). The star motif in the red box is the Reconnaissance port scan example shown in Figure 2. The top right of Figure 7 has more IPv6 addresses, compared to Figure 6.

Example data points for all attack tactics by count are presented in Table 7. The maximum count was 6,724,017, minimum count was 1, and average count was 3,864,567.

**Table 6.** Non-Attack Points of Interest (Count)

Id	From	To	Total_Dur	Avg_Dur	Total_Bytes	Avg_Bytes	Count	CountBin
edge_21	143.88.11.14	143.88.1.1.1	1267576.92	2.60	40376997	82.73	488029	5
edge_35	143.88.255.10	10.0.10.1	114.42	0.00	605569716	90.06	6724017	6

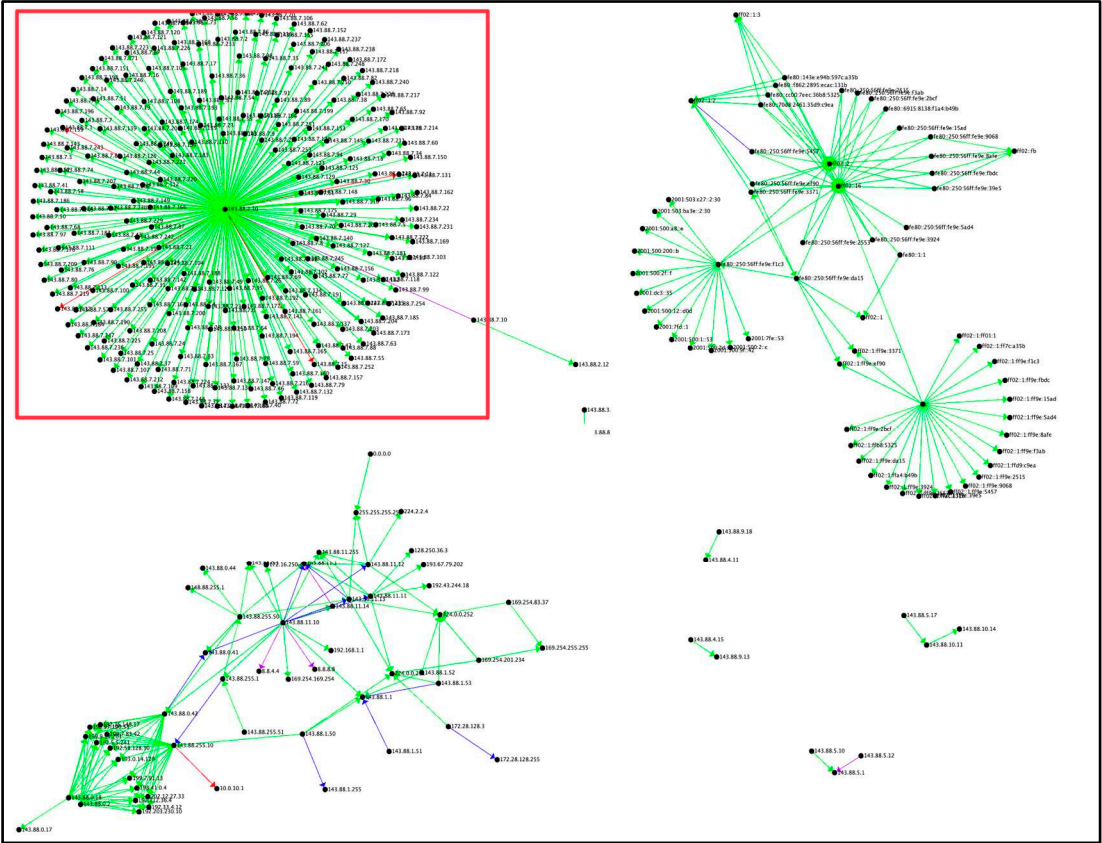


Figure 7. All Attack Tactics by Connection Count

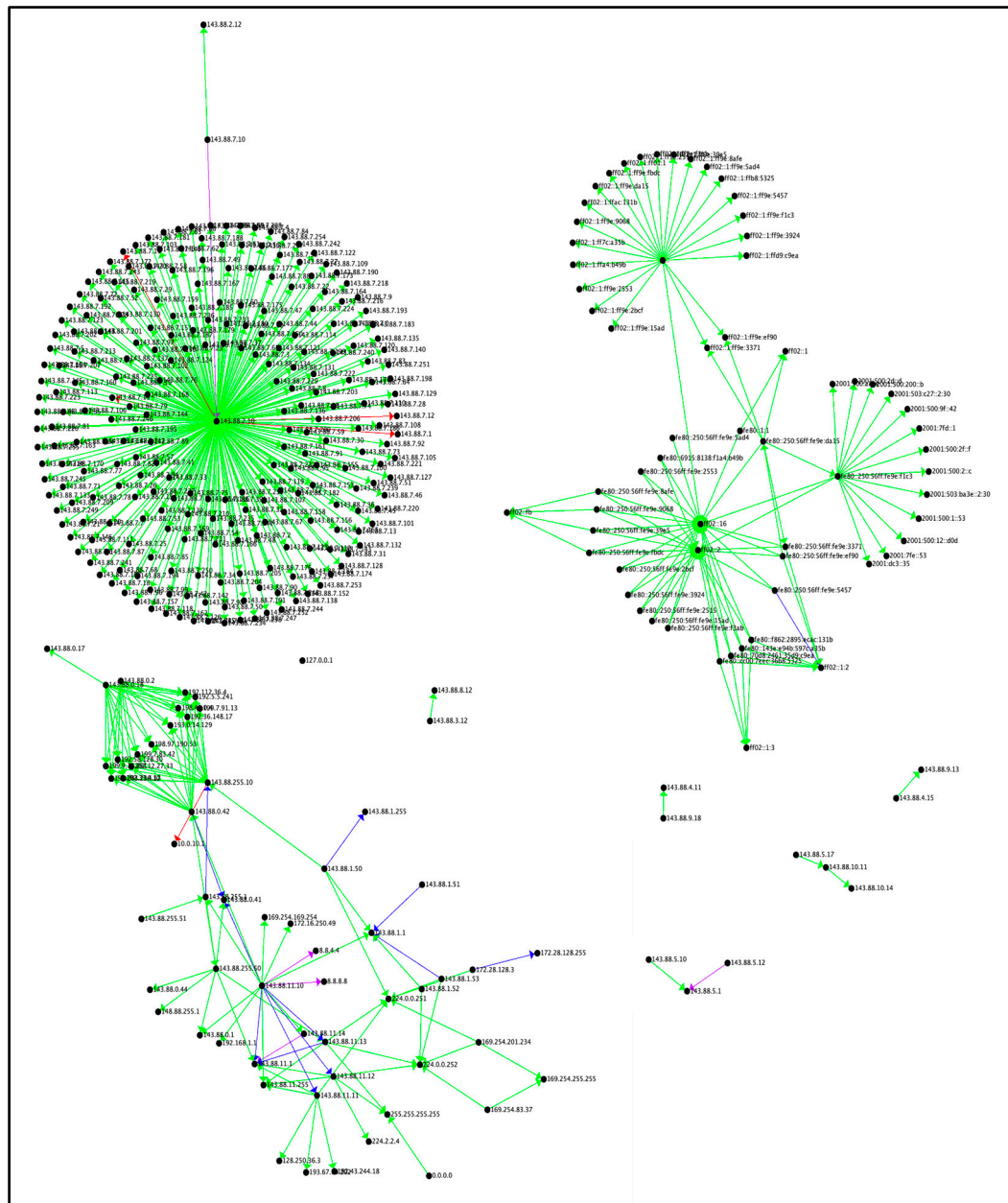
Table 7. All Attack Tactics Points of Interest (Count)

Id	From	To	Total_Dur	Avg_Dur	Total_Bytes	Avg_Bytes	Count	CountBin
edge_3	143.88.7.10	143.88.2.10	1216.984	0.002334	24576	0.047128	521472	5
edge_6	143.88.2.10	143.88.7.15	353248.5	0.212616	2654582328320	1597760	1661440	6
edge_7	143.88.2.10	143.88.7.11	972063.5	0.31234	5579520	1.792794	3112192	6
edge_8	143.88.2.10	143.88.7.1	279988	0.13377	8567808	4.093444	2093056	6
edge_9	143.88.2.10	143.88.7.12	778386.3	0.691353	925758636800	822247.5	1125888	6

edge_262	143.88.5.12	143.88.5.1	943576.7	1.87766	36458752	72.55069	502528	5
edge_267	143.88.11.10	8.8.8.8	588871.3	1.293066	43664530	95.88023	455407	5
edge_268	143.88.11.10	8.8.4.4	590266.6	1.300276	43591546	96.02614	453955	5
edge_284	143.88.11.14	143.88.11.1	1267577	2.597	40376997	82.73483	488029	5
edge_298	143.88.255.10	10.0.10.1	114.4165	0.000	605569716	90.06071	6724017	6

6.5 Visualizations of the Noncyclic Counts

Figure 8 represents the final count of connections for all identified attacks, with all cycles removed. All edges were added in this graph except for any edges that returned to a previously visited vertex. This allowed for the visualization of one-way traffic from the source to the destination. Adding the return cycles would have produced additional noise and could obscure the true target of the attack.



**Figure 8.** Noncyclic All Tactics by Connection Count

## 6.6 Summarizing the Graphic Visualizations

In this dataset, UWF-ZeekData22, the Star motif represents the Reconnaissance tactic well. The Reconnaissance tactic basically radiates from a single vertex, 143.88.2.10, to multiple other vertices in the graph. Figures 2, 3, and 4 are star motifs that depict the Reconnaissance tactic, but from different angles - connection count, duration, and byte count, respectively. The Clique motif was not useful in graphing the Reconnaissance tactic.

## 7. Runtime Performance

This section presents the runtime performance of the process of creating the graph databases, starting from file processing to the visualization of the graphs. In every case it can be noted that the truncated data, which is our reduced dataset used to create the graphs, performed better than the full data.

Table 8 presents the execution time for processing, including writing resulting output files, running on a quad core i5 intel processor at 2.4 GHz with 16 GB of DDR4 3200 ram. For both Phase 1

(File Processing) and Phase 2 (Graph Processing), it can be noted that the reduced data (with fewer attributes, used to create the graphs) performed better than the full data, which had all the attributes.

**Table 8.** Execution Time For Processing

	Phase 1 - File processing		Phase 2 - Graph processing	
	Duration (milliseconds)		Duration (milliseconds)	
	Reduced Data (84.3k rows)	Full Data (18.56 M rows)	Reduced Data (84.3k rows)	Full Data (18.56 M rows)
All rows	702 milliseconds	64,955	60	65
Reconnaissance	546	64,535	55	54
IP: 143.88.2.10	543	62,402	51	47

After file processing and graph processing, the resulting datafile is reduced to vertices and summed by connection count, connection duration, and bytes transmitted. These summed amounts are then binned across the vertices and graphed. Table 9 presents the execution time for binning and generating the resulting csv files after data processing, executing on 10-Core Intel Core i9 at 3.6 GHz with 32 GB of 2667 MHz DDR4 ram. It can once again be noted that the reduced data performed better than the full data.

**Table 9.** Execution Time for Binning and Generating Resulting csv Files

	Duration for Graph Stream part (milliseconds)		Row Count	
	Reduced Data	Full Data	Reduced Data	Full Data
All rows	39	41	374	480
Reconnaissance	39	40	255	258
IP: 143.88.2.10	38	38	254	256

Table 10 presents the execution time for generating GraphStream visuals after data binning, running on Quad-Core Intel Core i7 at 2.8 GHz with 16 GB of 2133 MHz LPDDR3 ram. Here we can see that the reduced data performed better for the Reconnaissance and the IP address 143.88.2.10.

**Table 10.** Execution Time for Generating Visuals

	Duration (milliseconds)		Row Count	
	Reduced Data	Full Data	Reduced Data	Full Data
All rows	7,904	6,967	374	480
Reconnaissance	7,510	7,644	255	258
IP: 143.88.2.10	6,834	7,241	254	256

## 8. Conclusion

The objective of this research was to determine if UWF-Zeekdata22 [9,10] could be mapped into a graph that could then be analyzed to yield consistent and identifiable patterns. Patterns involving network connectivity, connection duration, and data volume were found when the UWF-Zeekdata22 dataset was extracted and loaded into a graph environment. Patterns were also found in the graphed data that matched the attack tactics captured by UWF-Zeekdata22. The Reconnaissance tactic was represented well by the Star Motif.



There were some interesting discoveries when reviewing the resulting graphs. In the non-attack data, it was possible to identify normally occurring interactions between vertices in the graph. This could potentially be used to teach a ML what behaviors to ignore. This could potentially help identify zero-day attacks as they would not “look” like a learned normal behavior of the network. These graphs also provide insight into what the structure/topology of the network resembles.

Finally, an analysis of the run-time performance of the reduced dataset, using only four features from UWF-ZeekData22's Conn Log files and two additionally generated features plus count, showed that the reduced dataset performed better than the full dataset. Hence, a set of four connection features and two additionally generated features plus the count was enough for the graph engine to generate the graphs.

## 9. Future Works

The results in this paper show that graph databases/graph engines can be essential tools for understanding network traffic and detecting various network intrusions. The amount of data available for use in the analysis of this paper was fairly limited, so one area for future research is to apply the principles of this paper to multiple datasets and compare the results. Another area for further research is to use the models generated from this analysis to train machine learners. The learners would then be run against various simulated attack/non-attack data to determine the accuracy of the models.

**Supplementary Materials:** The dataset can be downloaded at [datasets.uwf.edu](https://datasets.uwf.edu)

**Author Contributions:** This work was conceptualized by S.S.B., D.M., S.C.B., M.P. and J.H.; methodology was mainly done by S.S.B., D.M., S.C.B., M.P. and J.H.; software was done by M. P. and J.H.; validation was done by D.M. and M.E.; formal analysis was done by M.P., J.H. and D.M.; investigation was done by S.S.B. D.M., M.P. and J.H.; data curation was done by M.P.; writing—original draft preparation was done by M.P. and J.H.; writing—review and editing was done by S.S.B., D.M., S.S.B., M.P., J.H. and M.E.; visualization was done by M.P. and J.H.; supervision was done by S.S.B., D.M. and S.S.B.; project administration was done by S.S.B. and D.M.; funding acquisition was done by S.S.B., D.M. and S.S.B. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research was funded by the National Centers of Academic Excellence in Cybersecurity, 2021 NCAE-C-002: Cyber Research Innovation Grant Program, Grant Number: H98230-21-1-0170

**Institutional Review Board Statement:** *Not applicable.*

**Informed Consent Statement:** *Not applicable.*

**Data Availability Statement:** The data is available at [datasets.uwf.edu](https://datasets.uwf.edu)

**Conflicts of Interest:** The authors declare no conflict of interest.

## References

1. Huong, T.T.; Bac, T.P.; Long, D.M.; Thang, B.D.; Binh, N.T.; Luong, T.D.; Phuc, T.K. LockEdge: Low-Complexity Cyberattack Detection in IoT Edge Computing. *IEEE Access* **2021**, *9*, 29696–29710, doi:10.1109/access.2021.3058528.
2. Leevy, J.L.; Hancock, J.; Zuech, R.; Khoshgoftaar, T.M. Detecting Cybersecurity Attacks across Different Network Features and Learners. *Journal of Big Data* **2021**, *8*, doi:10.1186/s40537-021-00426-w.
3. Bagui, S.; Simonds, J.; Plenkers, R.; Bennett, T.A.; Bagui, S. Classifying UNSW-NB15 Network Traffic in the Big Data Framework Using Random Forest in Spark. *International Journal of Big Data Intelligence and Applications* **2022**, *2*, 39–61, doi:10.4018/ijbdia.287617.
4. Zhang, J.; Sun, J.; He, H. Clustering Detection Method of Network Intrusion Feature Based on Support Vector Machine and LCA Block Algorithm. *Wireless Personal Communications* **2021**, *127*, 599–613, doi:10.1007/s11277-021-08353-y.
5. Kevric, J.; Jukic, S.; Subasi, A. An Effective Combining Classifier Approach Using Tree Algorithms for Network Intrusion Detection. *Neural Computing and Applications* **2016**, *28*, 1051–1058, doi:10.1007/s00521-016-2418-1.
6. MITRE ATT&CK Reconnaissance, Tactic TA0043 - Enterprise Available online: <https://attack.mitre.org/tactics/TA0043>.



7. Jia, Y.; Qi, Y.; Shang, H.; Jiang, R.; Li, A. A Practical Approach to Constructing a Knowledge Graph for Cybersecurity. *Engineering* **2018**, *4*, 53–60, doi:10.1016/j.eng.2018.01.004.
8. Oracle Corporation 17 Use Cases for Graph Databases and Graph Analytics. **2021**. Available online: <https://www.oracle.com/a/ocom/docs/graph-database-use-cases-ebook.pdf> (accessed on 19 August 2022).
9. Bagui, S.S.; Mink, D.; Bagui, S.C.; Ghosh, T.; Plenkers, R.; McElroy, T.; Dulaney, S.; Shabanali, S. Introducing UWF-ZeekData22: A Comprehensive Network Traffic Dataset Based on the MITRE ATT&CK Framework. *Data* **2023**, *8*, 18, doi:10.3390/data8010018.
10. University of West Florida UWF-ZeekData22 Available online: <https://datasets.uwf.edu> (accessed on 20 August 2020).
11. Coupette, C.; Vreeken, J. Graph Similarity Description. In Proceedings of the Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining; ACM: New York, NY, USA, August 14 2021.
12. Lee, M.-C.; Nguyen, H.T.; Berberidis, D.; Tseng, V.S.; Akoglu, L. GAWD. In Proceedings of the Proceedings of the 2021 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining; ACM: New York, NY, USA, November 8 2021.
13. Schindler, T. Anomaly Detection in Log Data Using Graph Databases and Machine Learning to Defend Advanced Persistent Threats Available online: <https://dl.gi.de/handle/20.500.12116/4016>.
14. Bai, Y.; Ding, H.; Bian, S.; Chen, T.; Sun, Y.; Wang, W. SimGNN. In Proceedings of the Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining; ACM: New York, NY, USA, January 30 2019.
15. GraphStream - A Dynamic Graph Library Available online: <https://graphstream-project.org/>.
16. Bagui, S.; Mink, D.; Bagui, S.; Ghosh, T.; McElroy, T.; Paredes, E.; Khasnavis, N.; Plenkers, R. Detecting Reconnaissance and Discovery Tactics from the MITRE ATT&CK Framework in Zeek Conn Logs Using Spark's Machine Learning in the Big Data Framework. *Sensors* **2022**, *22*, 7999, doi:10.3390/s22207999.
17. Sur, S.; Srimani, P.K. Topological Properties of Star Graphs. *Computers & Mathematics with Applications* **1993**, *25*, 87–98, doi:10.1016/0898-1221(93)90188-2.
18. Bron, C.; Kerbosch, J. Algorithm 457: Finding All Cliques of an Undirected Graph. *Communications of the ACM* **1973**, *16*, 575–577, doi:10.1145/362342.362367.
19. Mackaness, W.A.; Beard, K.M. Use of Graph Theory to Support Map Generalization. *Cartography and Geographic Information Systems* **1993**, *20*, 210–221, doi:10.1559/152304093782637479.
20. von Landesberger, T.; Görner, M.; Rehner, R.; Schreck, T. A System for Interactive Visual Analysis of Large Graphs Using Motifs in Graph Editing and Aggregation. *Proceedings of Vision Modeling Visualization Workshop* **2009**, *9*, 331–340.
21. MITRE ATT&CK Active Scanning, Technique T1595 - Enterprise Available online: <https://attack.mitre.org/techniques/T1595/>.
22. Chapter 4. Port Scanning Overview Available online: <https://nmap.org/book/port-scanning.html#port-scanning-intro>.
23. Frankel, S.; Green, D. Internet Protocol Version 6. *IEEE Security & Privacy Magazine* **2008**, *6*, 83–86, doi:10.1109/msp.2008.65.

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.