

Detecting the Behavioral Relationships of Malware Connections

Sebastian Garcia
Czech Technical University, FEE
Karlovo namesti 13, 121 35 Praha 2
Prague, Czech Republic
sebastian.garcia@agents.fel.cvut.cz

Michal Pechoucek
Czech Technical University, FEE
Karlovo namesti 13, 121 35 Praha 2
Prague, Czech Republic
michal.pechoucek@agents.fel.cvut.cz

ABSTRACT

A normal computer infected with malware is difficult to detect. There have been several approaches in the last years which analyze the behavior of malware and obtain good results. The malware traffic may be detected, but it is very common to miss-detect normal traffic as malicious and generate false positives. This is specially the case when the methods are tested in real and large networks. The detection errors are generated due to the malware changing and rapidly adapting its domains and patterns to mimic normal connections. To better detect malware infections and separate them from normal traffic we propose to detect the behavior of the *group* of connections generated by the malware. It is known that malware usually generates various related connections simultaneously and therefore it shows a group pattern. Based on previous experiments, this paper suggests that the behavior of a group of connections can be modelled as a directed cyclic graph with special properties, such as its internal patterns, relationships, frequencies and sequences of connections. By training the group models on known traffic it may be possible to better distinguish between a malware connection and a normal connection.

CCS Concepts

•**Security and privacy** → *Malware and its mitigation; Intrusion detection systems; Security protocols;*

Keywords

Malware Detection, Behavioral Models, Groups of Connections, Machine Learning

1. INTRODUCTION

The detection of infected computers is usually problematic not because the malware patterns can not be detected, but because too much normal patterns are mis-detected [6]. This problem arises because the patterns in the network generated by malware are very similar to the patterns of normal

traffic [7]. Moreover, it is known that normal connections usually use malware-like techniques, on purpose, to connect to their normal servers [5]. The problem is, then, how to better separate normal computers from infected computers given that the traffic of some connections looks very much the same.

Experiments with malware and normal captures have shown that most malware generates more than one connection [3], and that the connections seem to be coordinated in some way [1]. This suggested that a possible relationship between connections may be significant for a new detection method. The experiments done by the Stratosphere Project [2] have shown that it was already possible to model the behavior of individual connections for detection purposes. These detections were not exempt of errors, and those errors were the motivations for looking a detection alternative. The Stratosphere behavioral models are the base of our current analysis. We propose in this work to study the behavior of all the connections in a **group** simultaneously, in order to discover the properties of the group and to better distinguish the malware from the normal computers. We hypothesize that even though some malware connections can have the same behavior as normal connections, the group of connections itself must show intrinsic properties of the malware as a coordinated and logic unit.

There are previous analysis of groups in relation to malware. Most of these works grouped hosts to find all the infected bots, but they do not group connections from the same host. Moreover, they do not try to analyze the behavior of the *group*. The most related analysis has been the study of DNS requests. An example was the detection of malicious domain groups by using the temporal correlation of queries, given some well-known seed malicious domains [7]. The main difference with our work are (1) that they need a known malicious domain, and (2) that we analyze the behavior of the group. Another example analyzed the relationships between flows to form a graph structure of each host [4]. The authors summarized the behavior of a host in a vector in order to later cluster the malware families.

2. BEHAVIOR OF A CONNECTION

The base models for this work are the models created by Stratosphere. These models study which are the patterns of an individual connection in the network when all its flows are group together [1]. The stratosphere method only uses the flows in the network to create its models, and does not use the content of the packets. The base unit of work in what is called a *connection*. It is defined as all the flows that share

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

PrAISe '16, August 29 - 30, 2016, The Hague, Netherlands

© 2016 Copyright held by the owner/author(s). Publication rights licensed to ACM. ISBN 978-1-4503-4304-6/16/08...\$15.00

DOI: <http://dx.doi.org/10.1145/2970030.2970038>


```

10.0.2.105-197.149.90.166-12165-tcp
10.0.2.105-91.240.236.122-443-tcp
10.0.2.105-93.115.172.232-443-tcp
10.0.2.105-67.222.201.105-443-tcp
10.0.2.105-90.24.13.21-443-tcp
10.0.2.105-202.79.131.125-443-tcp
10.0.2.105-202.70.89.57-443-tcp
10.0.2.105-43.248.24.50-443-tcp
10.0.2.105-197.254.56.126-443-tcp
10.0.2.105-190.95.138.66-443-tcp
10.0.2.105-200.25.207.173-443-tcp
10.0.2.105-186.31.224.64-443-tcp
10.0.2.105-181.143.71.20-443-tcp
10.0.2.105-37.19.85.9-443-tcp
10.0.2.105-190.0.99.80-443-tcp
10.0.2.105-88.101.108.254-443-tcp
10.0.2.105-88.209.248.135-443-tcp
10.0.2.105-88.209.248.139-443-tcp
10.0.2.105-86.100.25.233-443-tcp
10.0.2.105-86.100.251.174-443-tcp
10.0.2.105-178.249.175.151-443-tcp
10.0.2.105-212.5.207.78-443-tcp
10.0.2.105-213.81.199.121-443-tcp
10.0.2.105-193.201.207.106-443-tcp
10.0.2.105-193.93.216.191-443-tcp
10.0.2.105-195.211.240.166-443-tcp
10.0.2.105-195.113.232.73-80-tcp
10.0.2.105-79.142.203.213-443-tcp
10.0.2.105-91.192.131.229-443-tcp
10.0.2.105-94.153.193.190-443-tcp
10.0.2.105-95.67.79.58-443-tcp
10.0.2.105-87.244.175.114-443-tcp
10.0.2.105-46.37.201.165-443-tcp
10.0.2.105-109.111.100.48-443-tcp
10.0.2.105-77.242.22.182-443-tcp
10.0.2.105-178.219.202.80-443-tcp
10.0.2.105-24.33.131.116-443-tcp
10.0.2.105-72.230.82.80-443-tcp
10.0.2.105-173.248.31.6-443-tcp
10.0.2.105-69.9.204.114-443-tcp
10.0.2.105-69.144.171.44-443-tcp
10.0.2.105-24.148.217.188-443-tcp
10.0.2.105-208.117.68.78-443-tcp
10.0.2.105-203.129.197.90-443-tcp
10.0.2.105-112.133.203.43-443-tcp
10.0.2.105-27.109.20.53-443-tcp

```

Figure 3: Upatre malware group behavior. Upatre generated several identical connections to its C&C servers. The connections where generated in a specific order that revealed its behavior.

```

147.32.83.53-147.32.83.105-arp
147.32.83.53-147.32.80.105-53-udp
147.32.83.53-173.194.113.191-443-tcp
147.32.83.53-195.113.214.241-443-tcp
147.32.83.53-130.15.100.14-80-tcp
147.32.83.53-80.242.138.72-80-tcp
147.32.83.53-195.113.214.230-80-tcp
147.32.83.53-147.32.80.9-53-udp
147.32.83.53-152.3.140.5-80-tcp
147.32.83.53-147.32.80.9-0x1206-icmp
147.32.83.53-81.27.192.20-123-udp
147.32.83.53-37.157.199.158-123-udp
147.32.83.53-178.63.212.146-123-udp
198.143.173.181-147.32.83.53-80-tcp
147.32.83.53-173.194.113.183-443-tcp
147.32.83.53-54.243.152.237-443-tcp
147.32.83.53-195.113.214.219-443-tcp
147.32.83.53-146.95.130.38-80-tcp
147.32.83.53-195.113.214.249-443-tcp
147.32.83.53-80.64.49.11-443-tcp
147.32.83.53-173.194.70.95-443-tcp
147.32.83.53-54.230.8.157-443-tcp
147.32.83.53-173.194.32.222-443-tcp
80.75.105.148-147.32.83.53-13907-tcp
147.32.83.53-80.64.49.12-443-tcp
147.32.83.53-195.113.214.215-443-tcp
147.32.83.53-167.68.20.174-80-tcp
147.32.83.53-84.18.180.87-80-tcp
147.32.83.53-209.167.231.15-80-tcp
147.32.83.53-195.113.214.211-80-tcp
147.32.83.53-167.68.20.181-80-tcp
84.18.164.15-147.32.83.53-33435-udp
147.32.83.53-224.0.0.251-5353-udp
147.32.83.53-173.194.35.79-443-tcp
147.32.83.53-212.70.64.183-80-tcp

```

Figure 4: Behavior of a normal computer in a University network. The connections have different patterns and none is periodic. It is not easy to group connections since there are no clear relationships. This image only shows the significant connections, see the video for a complete recollection.

patterns, they use the same destination port, and they were made almost simultaneously. We could identify two ways of creating the groups: (1) group the connections that have the same individual behavioral pattern of letters, (2) group the connections that have a relationship based on the destination IP address (including the WHOIS information). Both of these techniques would group the malware connections and would also group some normal connections. We use the first way of creating the groups based on the similarity of the individual behavioral patterns.

To analyze the behavior of the group it is necessary to analyze its components. Each group is composed of the connections on it, identified as C_i , and the flows arriving to all the connections, identified as F_j . From each flow it is possible to extract: its arrival time, its duration, its size and the connection to which flow belongs to. Each time a flow arrives there is a *transition* from the last connection that received a flow to the next, and therefore a sequence: $C_1 \rightarrow C_2 \rightarrow C_3 \rightarrow C_1$. Based on C_i , F_j and this last sequence, the following components of a group are computed:

- A directed graph, identified as D computed from the previous sequence.
- The amount of connections.
- The Stratosphere behavioral pattern of each individual connection computed from F_j .

The use of a directed graph gives us a structure to analyze the group behavior. More information is needed to complement it, such as the frequencies of the transitions in the graph. Each directed graph D has its own properties:

1. If D is a cyclic graph, i.e. if it is a loop.
2. The amount of cycles performed by the flows.
3. The frequency of the loop.
4. The frequency of sending a flow to the next connection in D .
5. The frequency of sending a next flow to the *same* connection in D .
6. The delay between the last and first connection in the group.

An example graph of the Geodo group behavior is shown in Figure 5, with its fourteen connections. This graph shows how each new arriving flow goes from one connection to the next connection in the group.

The properties of D can then be analyzed based on trained malware and normal captures to find the behaviors that best separate them. In particular the first restriction for being considered a group behavior is that the graph should be cyclic. If there is no repetitive pattern, then there is no group behavior. Another restriction is that there should be a minimum amount of cycles to be detected as malware, since it indicates that the loops are not by chance. The frequency of the loop and the frequency of sending a flow to the next connection may be a characteristic of each botnet and therefore may be trained. The frequency of sending a flow to the same connection is already incorporated into the Stratosphere behavior of that connection. Finally, the delay

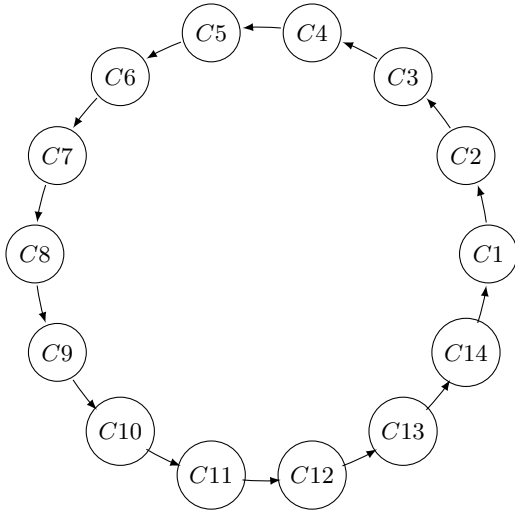


Figure 5: Directed Cycle Graph for the Geodo malware group behavior. There are 18 connections, and the sequence is a round-robin schedule.

to restart the loop is also characteristic and can be changed during the lifetime of the malware.

Based on these properties it is possible to analyze the malware. It should be considered that the Geodo malware showed two distinctive patterns, $b+b+b+$ and $b*b*b*$, and therefore some of the values changed for the group behavior. The properties of Geodo group are:

- Amount of connections: 14
- Amount of cycles performed: 848
- If the graph D is cyclic: Yes.
- Frequency of loop: 3 minutes, 20 seconds ($b+b+b+$ pattern).
- Frequency of loop: 17 minutes 20 seconds ($b*b*b*$ pattern).
- Frequency of sending a flow to the next connection: 10 seconds without considering the restart of the loop.
- Delay between the last and first connection: 1 minute 10 seconds ($b+b+b+$ pattern).
- Delay between the last and first connection: 15 minutes 10 seconds ($b*b*b*$ pattern).
- Frequency of sending a flow to the same connection: 3 minutes, 20 seconds ($b+b+b+$ pattern).
- Frequency of sending a flow to the same connection: 17 minutes, 20 seconds ($b*b*b*$ pattern).

The example of the Upatre malware is more complicated, given that the C&C connection *10.0.2.105-197.149.90.166-12165-tcp* is used in two groups with different connections but the same properties. The properties of both groups are:

- Amount of connections: 24
- Amount of cycles performed: 72

- If the graph D is cyclic: Yes.
- Frequency of loop: 33 minutes.
- Frequency of sending a flow to the next connection: 21 seconds
- Delay between the last and first connection: 17 minutes.
- Frequency of sending a flow to the same connection: 33 minutes.

These values for the properties of the malware groups show how different they are and how characteristic they can be. The case of the Normal capture was easy in this example because it didn't have any group of connections that could be distinguished by our method. This means that if we use the detection method of the group behavior it may be possible to differentiate the malware infections from the normal traffic. We acknowledge that more normal connections are needed for a better testing. The identification of group properties in these examples suggest that they may be a good model of behavior for a detection method.

5. CONCLUSIONS

This position paper proposes that by studying the behavior of a group of connections it may be possible to improve the detection of infected computers. The motivation was that some normal connections look like malware connections and they are not easy to distinguish. We proposed that connections can be grouped together according to how they are related, and those groups can be analyzed to find their properties, sequences and behavioral patterns. The properties of a group seem to have enough differential power to be able to separate normal and botnet computers. Malware may change to avoid any technique, but the behavioral changes may be too costly for them. This research will continue by formally defining the properties of groups, experimenting with larger datasets and using the groups to detect infected computers.

6. REFERENCES

- [1] S. Garcia. Modelling the Network Behaviour of Malware To Block Malicious Patterns . the Stratosphere Project : a Behavioural Ips. In *Virus Bulletin*, number September, pages 1–8, 2015.
- [2] S. Garcia, M. Grill, J. Stiborek, and A. Zunino. An Empirical Comparison of Botnet Detection Methods. *Computers & Security*, 45(0):100 – 123, 2014.
- [3] S. Garcia, V. Uhlir, and M. Rehak. Identifying and Modeling Botnet C&C Behaviors. In ACM, editor, *Proceedings of the 1st International Workshop on Agents and CyberSecurity (ACySE '14)*, New York, 2014. ACM.
- [4] S. Nari and A. A. Ghorbani. Automated Malware Classification based on Network Behavior. pages 642–647, 2013.
- [5] M. Patterson. Security Vendors Teaching Bad Actors How to Get Past Firewalls, 2015.
- [6] R. Sommer and V. Paxson. Outside the Closed World: On Using Machine Learning for Network Intrusion Detection. *2010 IEEE Symposium on Security and Privacy*, 0(May):305–316, 2010.

- [7] H. R. Zeidanloo and A. B. A. Manaf. Botnet Detection by Monitoring Similar Communication Patterns. *International Journal of Computer Science and Information Security*, 7(3):36–45, 2010.