

## ASMATRA: Ranking ASs Providing Transit Service to Malware Hosters

Cynthia Wagner, Jérôme François, Radu State, Alexandre Dulaunoy, Thomas Engel, Gilles Massen

### ► To cite this version:

Cynthia Wagner, Jérôme François, Radu State, Alexandre Dulaunoy, Thomas Engel, et al.. ASMATRA: Ranking ASs Providing Transit Service to Malware Hosters. International Symposium on Integrated Network Management, May 2013, Ghent, Belgium. IEEE, 2013, Proceedings of the 13th IFIP/IEEE International Symposium on Integrated Network Management. <hal-00846082>

HAL Id: hal-00846082

<https://hal.inria.fr/hal-00846082>

Submitted on 18 Jul 2013

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# ASMATRA: Ranking ASs Providing Transit Service to Malware Hosters

Cynthia Wagner\*, Jérôme François†, Radu State†, Alexandre Dulaunoy ‡, Thomas Engel†, Gilles Massen\*

\*RESTENA Foundation, Luxembourg, firstname.lastname@restena.lu

†SnT – University of Luxembourg, Luxembourg, firstname.lastname@uni.lu

‡Computer Incident Response Center Luxembourg, Luxembourg, alexandre.dulaunoy@circl.lu

**Abstract**—The Internet has grown into an enormous network offering a variety of services, which are spread over a multitude of domains. BGP-routing and Autonomous Systems (AS) are the key components for maintaining high connectivity in the Internet. Unfortunately, Internet Service Providers (ISPs) operating ASs do not only host normal users and content, but also malicious content used by attackers for spreading malware, hosting phishing websites or performing any kind of fraudulent activity. Practical analysis shows that such malware-providing ASs prevent themselves from being de-peered by hiding behind other ASs, which do not host the malware themselves but simply provide transit service for malware.

This paper presents a new method for detecting ASs that provide transit service for malware hosters, without being malicious themselves. A formal definition of the problem and the metrics are determined by using the AS graph. The PageRank algorithm is applied to improve the scalability and the completeness of the approach. The method is assessed on real and publicly available datasets, showing promising results.

## I. INTRODUCTION

In general, data is transmitted through multiple Autonomous Systems (ASs), which can be operated by distinct legal entities. Referring to [1], an AS is defined as the unit of routing policy or a collection of links and routes for an operator. AS routing is based on the Border Gateway Protocol (BGP) [2]. This protocol acts at the level of address prefixes as a path vector protocol with aim to guarantee high end-to-end connectivity in the Internet.

Since data travels through more than one AS on its path from sender to receiver, this is also effective for the transmission of malware<sup>1</sup>. While tracking ASs hosting malware is feasible, it is used in this paper in combination with PageRank to show that ASs looking innocent can forward malware from source to an end-point as a transit service. This is the main objective of our work as such ASs guarantee a high connectivity for malware spreading without being necessary present in common blacklists since they are not malware provider themselves.

In this paper, our approach named ASMATRA (1) formally defines a method to assess the capacity of an AS to provide transit to malware hosting ASs, (2) shows how it can be approximated by a new metric over the AS path, (3) analyses the link structure of the entire AS graph in a scalable way and (4) validates the approach on real scenarios.

<sup>1</sup>In this paper, malware is used as a generic term for any kind of malicious activities including worm spreading, spamming, phishing, etc.

The paper is organized as follows: ASMATRA is described from a macroscopic point of view in section II, section III introduces background information. Section IV describes the approach for a single AS and the global approach using PageRank is given in V. Section VI discusses the experimental results. Section VII presents related work and section VIII presents the conclusions and planned future work.

## II. OVERVIEW

Figure 1 highlights information flow and processing in ASMATRA. The initial step is to collect BGP announces which are exchanged between routers of distinct ASs. As highlighted, BGP-ranking is leveraged [3]. This tool relies on blacklists of IP addresses and subnetworks in order to score each AS regarding the proportion of hosted malware. Hence, it is useful to detect ASs hosting malware but not those providing transit service. To do that, ASMATRA uses the output of BGP-ranking but will also create a graph representing the interconnections between ASs (2). Then, a graph analysis using Page is performed where BGP-ranking helps in weighting properly the algorithm. Using also other standard metrics about the graph and nodes, each AS can be individually scored to represent its capacity to provide transit for malware.

Therefore, ASMATRA needs few requirements: BGP announces and blacklists for running BGP-ranking (or access to the web interface).

## III. BACKGROUND

### A. The Border Gateway Protocol (BGP) and Autonomous Systems (AS)

AS is the unit of a routing policy or a collection of IP links/routes for one or more administrative operators [1]. It can be distinguished between three AS types. The first type is *stub* where an AS is only connected to one AS. The next type is *multihomed*, where an AS has connections to many ASs for improving its own connectivity, but does not forward traffic between connected ASs (transit). The last type is the *transit* AS, where an AS provides transit between ASs connected to it. Each AS has an ASN (Autonomous System Number), a unique number to identify its network, but for more about ASs, the reader may refer to [1] and [4].

The BGP was introduced to control the route selection and the transmission of data between ASs. A BGP-router maintains a table with the path (AS path) to reach a given IP-prefix. Since 2001, CIDR (Classless Inter-Domain Routing) is leveraged for route aggregation.

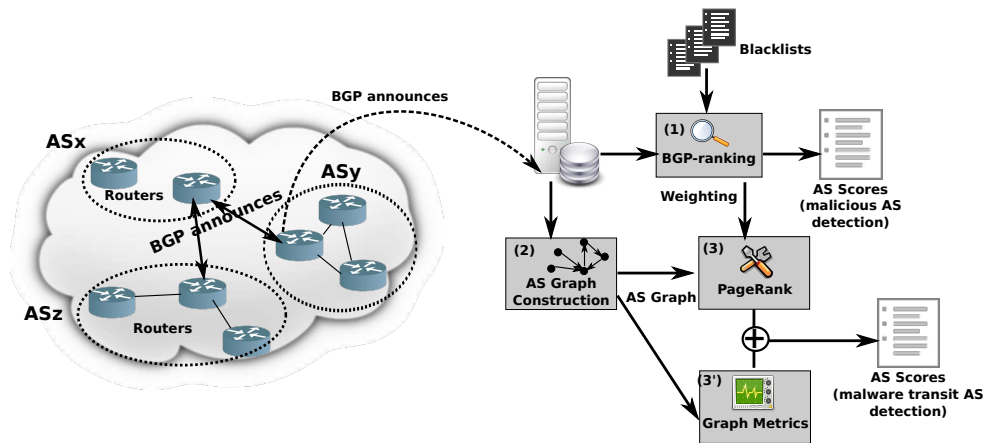


Fig. 1: Information flow and processing of ASMATRA

The information from the IP-table is sufficient to generate AS graphs for representing the connectivity and reachability of ASs. Based on this, routing loops may be pruned or decision policies be enhanced. Details about BGP are given in [2].

### B. Malicious Autonomous Systems

Some ASs are more tolerant regarding their customers' activities than others. Besides this, each AS may be susceptible to host malware, since customer activities are not permanently controlled. Incident report tools may help in the response to malicious activities, since complaints or discovered malware can be reported there and further action be taken, for example by using filtering. In [5] for example, the authors show that there are numerous ASs, which exhibit an abnormal proportion of malware hosting and collateral effects, i.e. launching attacks, etc. For McColo [6] or Atrivo [7], for example, 78% of the servers and domains have been estimated being hostile. In [8], it is claimed that there are others, which are exclusively operated for criminal activities, like for instance ValueDot. In general, once such a behaviour is detected, the AS is disconnected/de-peered from the Internet by a court decision and in most cases all by cooperation among legitimate upstream ASs, which provide connectivity to the rest of the Internet. Even if this works, it acts with delay and usually does not prevent from a resurrection of this AS elsewhere on the Internet, especially in countries where legislations and politics are more permissive. BGP-ranking [3], described in the next section, can help anticipating this reappearance process by combining knowledge from multiple blacklists.

Cyber-criminals go beyond this and operate transit ASs as well. Hence, the upstream transit AS, which was previously able to disconnect a malicious AS once detected, is now controlled by the criminals. Such transit ASs are harder to detect because IP addresses of malware are not listed in this AS. Looking at the well known example of the Russian Business Network (RBN) in Fig.2, it was operated until 2007 and is composed of several ASs involved in malicious activities such as the main RBN or Credolink AS, but connectivity in the Internet is obtained by transiting through SBTtel. Since this one-hop transit does not really obfuscate and is easily discoverable, cyber-criminals integrate more levels of transit

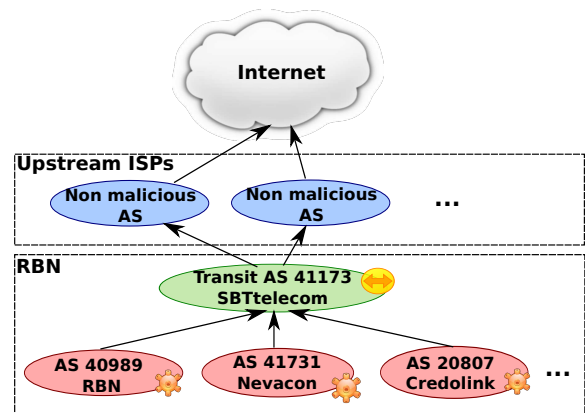


Fig. 2: Overview of the Russian Business Network (RBN)

AS, which allows RBN to escape a de-peering for 3 years. An explicit study of RBN is given in [8].

ASMATRA focuses on detecting AS providing transit to malware hoster but does not provide any counter-measures. However, automatic de-peering is a possibility but it's quite extreme and a more cautious mode, in particular through manual investigation, is advised. In this way, ASMATRA can be seen as a tool to guide manual investigation in order to reduce the necessary time to perform it which, as shown before, can take several years.

### C. BGP-ranking

This section describes BGP-ranking [3] software<sup>2</sup> that aims to score ISPs based on malware they host. BGP-ranking has been operated for more than 2 years, collecting malicious activities per ISP.

The scoring,  $AS_{rank}$ , is a float value starting from 1 to an unbounded value. 1 is the default score for an active ASN (announcing at least 1 network block). The scoring is computed from a set of lists,  $BL$ , holding publicly known

<sup>2</sup>available at <https://github.com/CIRCL/BGP-Ranking>

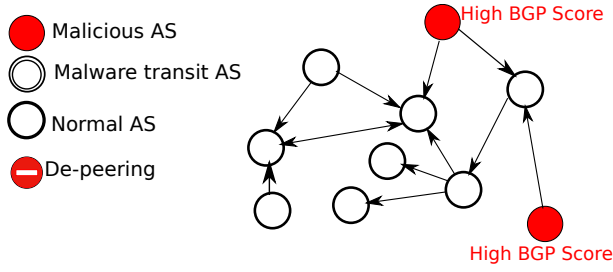


Fig. 3: Trivial example of an AS topology

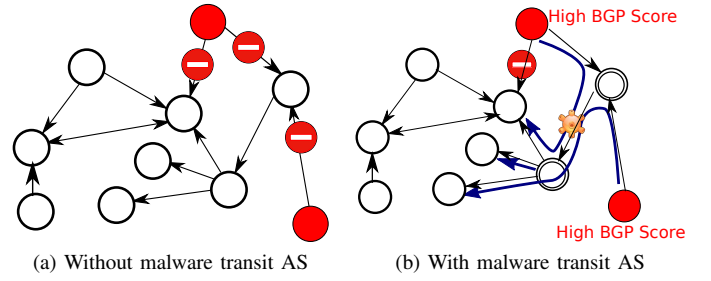


Fig. 4: De-peering known malicious ASs

malicious source IP addresses<sup>3</sup>. A list,  $b \in BL$ , is collected by the BGP-Ranking system for a day. Since the quality of the lists is variable, the BGP-Ranking operator adds an impact factor to each list,  $b_{impact}$ .

Then, for each occurrence of a malicious IP address within  $ASx$ , an occurrence counter is incremented leading to determine  $occ(b, ASx)$ , the number of occurrences of IP addresses included in a prefix announced by  $ASx$  and listed in  $b$ . For all lists, the sum of all occurrences is computed by using the impact factor  $b_{impact}$ . ASNs announce network sizes, therefore the sum is divided by the overall network size,  $ASx_{size}$ , expressed in the number of IP addresses. This weights the overall scoring for each ASN and gives a comparable result that can be defined as:

$$AS_{rank}(ASx) = 1 + \frac{\sum_{b \in BL} occ(b, ASx) b_{impact}}{ASx_{size}} \quad (1)$$

The scoring system is operated by CIRCL<sup>4</sup> and a sample output of BGP-ranking looks like:

```
whois -h pdns.circl.lu 6661
# ASN,Rank,Matched black list,Highest Ranking, Position
6661,1.00001786394817,1/13,3.078125,4379/52883
```

#### IV. PROBLEM DESCRIPTION AND METRICS

##### A. Detecting malware transit service

The objective is to detect ISPs providing transit service to malware hosters (e.g. SBTtel in Fig. 2). A simple example of our approach is given in Fig. 3. The naming of the AS is described below:

- **malicious AS**: hosts malware or from a general perspective related to malicious activities,
- **normal AS**: benign AS, not exhibiting an excessive proportion of malicious activities compared to other ones,
- **malware transit AS**: provides transit for malicious ASs.

Malicious and normal AS are supposed known (assuming a BGP ranking threshold for instance) in the examples to keep the problem focused on malware transit AS detection. However, this clear distinction (malicious/normal) does not

hold in practice and we use the BGP Ranking scores as a mean for weighting our approach as mentioned in next sections.

Based on BGP-ranking, detecting malicious ASs is viable, see Fig. 4(a), where such an AS can be isolated thanks to coordination of normal upstream ASs. However, a malware transit AS is not necessarily highly scored by BGP-ranking but provides Internet connectivity to malicious ASs. Thus, the latter can continue its activity without being filtered like in Fig. 4(b). Our method is conceived to detect such ASs.

##### B. AS graph

The ASs have to exchange data with selected neighbouring ASs. This information is essential for determining AS paths, i.e. successive ASs to go through for reaching an IP address block. The definitions *carry* and *transit* are taken from [9].

*For the inter-domain routing process, if  $ASy$  can reach a given prefix  $N1$  through  $ASx$ , then  $ASx$  carries  $N1$  for  $ASy$ . This is denoted as  $N1 \in ASx \xrightarrow{c} ASy$*

*An  $ASx$  transits a given prefix  $N1$  for another  $ASy$ , if and only if*

- $ASx$ ,  $ASy$ , and  $ASz$  are 3 different ASs,
- $ASy$  carries  $N1$  for  $ASz$
- $ASx$  carries  $N1$  or  $N1$ 's less specific for  $ASz$ .

*This transit can be denoted as:  $N1 \in ASx \xrightarrow{t} ASy$*

Therefore, an AS-path between  $ASx$  and  $ASy$  for a network prefix  $N$  is a sequence of ASs that will forward traffic towards  $N$  from  $ASx$  to  $ASy$ . Thus, these ASs transit the prefix  $N$  for  $ASy$ . For our purpose, We refine the definition of transit AS as follows:  *$ASx$  transits  $ASy$  to  $ASz$ ,  $ASy \xrightarrow{ASx} ASz$ , if and only if there is a network prefix  $N$ , such that  $N \in ASy \xrightarrow{c} ASz$  and  $ASx$  carries  $N$  or  $N$ 's less specific for  $ASz$ .*

Referring to the description of [9], an AS graph is defined as:

*An AS graph is a directed graph, denoted as  $G(V, E)$ , where  $V$  is a set of nodes representing the ASs and  $E$  the set of directed edges, denoting the connections between two end-node ASs. By definition, this graph may contain loops, but not multiple edges.*

The AS graph does not hold multiple edges since the contribution in routing is aggregated. It shows the reachability

<sup>3</sup>see [www.blocklist.de](http://www.blocklist.de), [www.dshield.org](http://www.dshield.org)... Full configuration at: <https://github.com/CIRCL/bgp-ranking/blob/master/etc/bgp-ranking.conf.redis>

<sup>4</sup><http://bgpranking.circl.lu/>

among ASs. There is an edge between two ASs, if they are sequentially represented in at least one AS path. Such a representation summarizes BGP-routing by factorizing AS paths. Fig 3 illustrates these concepts with unnumbered ASs.

### C. Theoretical measures for malware transit ASs

By definition, a malware transit AS provides transit service to other ASs, known to host malware. Assuming,  $MT(ASx)$  measures the capacity of  $ASx$  to support malware transit, it can be computed by summing the BGP-ranking value of the ASs for which transit is provided:

$$MT(ASx) = \sum_{\substack{(ASy, ASz) \\ \in \{(a,b) | a \xrightarrow{ASx} b\}}} AS_{rank}(ASy) \quad (2)$$

$MT(ASx)$  does not measure the impact of a malicious AS on other ones through  $ASx$  in terms of volume of new malware spread. From a theoretical point of view, this needs exact knowledge about hosted malware (type, version,...), which is practically impossible. However BGP ranking may be considered as a good reference for computing the difference between ASs and equation (2) is refined:

$$MT(ASx) = \sum_{\substack{(ASy, ASz) \\ \in \{(a,b) | a \xrightarrow{ASx} b\}}} \frac{(AS_{rank}(ASy) - AS_{rank}(ASz))^+}{card(\{ASu \in V, ASy \xrightarrow{ASu} ASz\})} \quad (3)$$

where  $card(S)$  is the cardinality of  $S$ . The malware impact of  $ASy$  on  $ASz$  evaluated as the BGP Ranking difference (numerator) is distributed among all ASs, which provide an equivalent transit (denominator). By approximation, an AS can transfer malware to another only if its BGP-ranking is higher. Therefore, we use the positive part  $^+$  ( $x^+ = max(x, 0)$ ).

### D. Theoretical measure estimation

Since the previous measure relies on exact AS paths, it raises scalability issues due to the large number of paths to consider. Moreover, while some routes are stable [10], most of them are still unstable [11]. For example, it has been shown that most the average life of a route is around one day. The analysis needs to collect AS paths over long time periods (weeks) to correctly identify ASs providing transit for others. This leads to an increase in complexity, while the AS graph is a summarized view with some advantages:

- the number of links is very low,
- it highlights possible routes between ASs which are not visible in AS paths (yet or anymore). For example in Fig. 5, the  $ASx$  routing table (plain arrows) has a path to  $ASz$  through  $ASy$  and has a direct path to  $ASu$ . However,  $ASu$  is also peered with  $ASz$  and can also be an AS-path from  $ASx$  to  $ASz$ , when the other route is congested or if the policy changes. This is also true for malicious ASs, which regularly have paths filtered and so change their AS-paths frequently.

Based on this reasoning, the AS-graph is now considered as undirected, which is achieved by removing directions of the AS graph as described before. For the sake of clarity, a



Fig. 5: Example for an alternative route in a trivial AS graph

path denotes now a standard path (sequence of nodes) in the AS graph and not an AS-path anymore. The malware transit capacity of an AS,  $ASx$ , is computed by the BGP-ranking differences between all pairs of ASs,  $(ASy, ASz)$ , such that there is a path between  $ASy$  and  $ASz$  through  $ASx$ .

However, there are numerous paths connecting two ASs in Internet. Considering all of them is not realistic for scalability and also from a conceptual point of view (e.g., the probability is very low that routing from  $ASx$  to  $ASy$  contains hundreds intermediate nodes while there are shorter connections). The exploration of the graph can be limited by a predefined radius around the considered ASs.

Assuming  $G'(ASx, k)$ , a graph composed of nodes at a maximum distance  $k$  from  $ASx$ . The latter is considered as a transit AS between  $ASy$  and  $ASz$  if and only if it is in the unique path from  $ASy$  to  $ASz$ . The notion of unique path avoids that too long paths are artificially considered. We propose to find graph partitions (connected components) of  $G'(ASx, k)$  by discarding  $ASx$ . Thus,  $G''(ASx, k)$  is the graph generated from  $G'(ASx, k)$  where  $ASx$  and subsequent edges are removed. The malware transit capacity of  $ASx$  is now evaluated thanks to its capacity to transfer malware between groups of ASs in  $C_k = CC(G''(ASx, k))$ , the set of subgraphs representing the connected components of  $G''(ASx, k)$ :

$$MT'_k(ASx) = \frac{\sum_{(c1, c2) \in pairs(C_k)} \left| \sum_{a \in c1} Rank_a - \sum_{b \in c2} Rank_b \right|}{\#neighs(ASx)} \quad (4)$$

where  $pairs(S)$  is the set of all possible unordered pairs of elements from  $S$  (2-subsets). Assuming a pair  $\{a, b\}$ , the absolute value has to be used to estimate the ability of  $a$  to affect  $b$  or vice-versa. To avoid a bias for highly peered AS like major operators, the value is normalized by the number of peering connections represented by the number of neighbours,  $\#neighs(ASx)$ .

A first metric to asses the capacity of an AS to transit malware (equation (3)) was introduced, based on the natural definition of a malware transit AS. Since its computation raises several issues, an estimated metric was derived that limits the evaluation of an AS only to its neighbourhood. A new parameter,  $k$ , needs to be defined. Next section presents a solution avoiding this additional parameter and employing PageRank. PageRank limits the computational overhead as it allows to compute the malware transit capacity for all ASs within a single execution, while equation (4) calculates it for an individual AS.

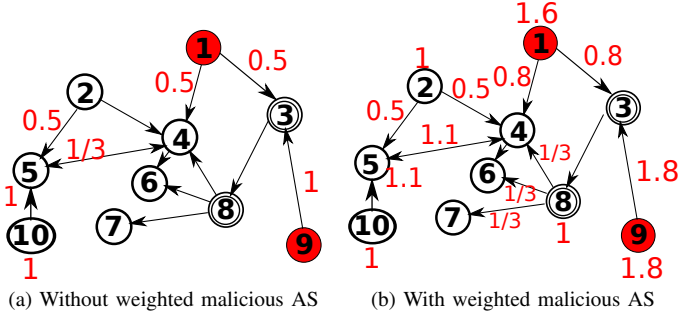


Fig. 6: PageRank on AS topology (Only score values mentioned in text are represented for sake of clarity)

## V. AS RANKING USING PAGERANK

The PageRank algorithm [12], [13] is used by the Google web search engine. In PageRank, a score is estimated for each webpage and is based on the Weblink organization. A web page is referred to by other pages, if other pages contain a hyperlink pointing to it. Intuitively, a frequently-referred page is important and pages referred to by a few important pages, are considered important too. For instance, if a web page is linked by only one web site, but this web site being Yahoo, then this page should be considered as important too.

In this paper, PageRank is applied to an AS graph to detect malware transit ASs because such ASs are well connected to malicious ones. This can have a higher impact during PageRank computation using weighting based on BGP-ranking. PageRank is computed iteratively until scores are stable. At each iterating step, the current score of a node is distributed through its outgoing links. The new score for a node is given by the sum of the incoming edges scores.

Fig. 6(a) and 4(b) show an example using the same node legend than Fig. 3. Assume no prior knowledge, then each node (equivalent to an AS) starts with a score of one. Then, the node 1 distributes its initial score of one through its two outgoing links. Node 3 receives 0.5 as well as 1 from node 9 leading to a global score of 1.5 after one iteration. Node 5 receives a total of 1.83. The iteration process continues, but within a single step, it is already observable that node 5 provides better connectivity to others, but the malware transit AS (node 3 and 8) cannot be detected.

Adding knowledge, such as the BGP-ranking score, is required to propagate bigger values from malicious ASs to transit ASs. In Fig. 6(b), scores close to nodes represent the BGP-ranking for the ASs. Since node 1 and 9 represent the malicious ASs, they receive higher scores, *i.e.* 1.6 and 1.8, which are transmitted to node 3. The latter thus obtains  $0.8 + 1.8 = 2.6$ . However, a normal AS that plays a central role in the topology can have a higher score as this an objective of PageRank. For instance, node 4 in Fig. 6(b) got  $0.8 + 0.5 + 1.1 + 0.33 = 2.73$ . In fact, it also provides transit service to a malicious AS (node 1), but should be less considered as a malware transit AS because it is also the upstream AS of many normal ones. By de-peering the malicious AS similar to Fig.4(b), the scores drop to 1.93, quite lower than the score of node 3. De-peering is not instantaneous because a malicious AS has first to be detected

(IP addresses related to malicious activities for example) before triggering legal and/or technical countermeasures.

In this paper, we consider this case by supposing that an AS providing transit service to normal ASs mainly, cannot be considered as a malware transit AS. The score is thus normalized regarding the number of neighbours. In this trivial example, a simple division by the number of neighbours is enough. Node 4 having three incoming connections will then obtain a score  $2.73/4 = 0.6825$ , while the malware transit node 3 will get  $2.6/2 = 1.3$ .

From a formal point of view, PageRank can be described as follows [14]:

*Let  $P_t(i)$  be the score of a node  $i \in V$  for an iteration  $t$ ,  $(j, i) \in E$  which describes a directed edge from a node  $j$  to a node  $i$ .  $O_j$  is denoted as the number of outgoing links in a node  $j$ .  $I_j$  is denoted as the number of incoming links in a node  $j$ .  $W(i)$  is the initial weight of the node  $i \in V$*

As mentioned before, the weight  $W(k)$  of an AS is given by its BGP-ranking score. While the damping factor  $d$  in the original PageRank represents a user clicking randomly on a webpage, it allows here to balance the impact on the score between the previously computed score and the weight of a node, *i.e.* the BGP-ranking of the AS. The computation efficiency [13] is optimized as follows,

$$P_t(i) = (1 - d) \sum_{k=1}^n W(k) + d \sum_{(j,i) \in E} \frac{P_{t-1}(j)}{O_j} \quad (5)$$

In our case, the graph is undirected, *i.e.*  $(i, j) \in E \Leftrightarrow (j, i) \in E$ . Besides, the iterations stops when  $|P_t - P_{t-1}| < 10^{-12}$ .

More details about PageRank can be read in [12].

As discussed before, PageRank is slightly modified for discarding normal ASs providing minor transit service to malicious AS based on the number of neighbours. To normalize the PageRank score, it can be assumed that most ASs are not providing transit for malware. Hence, from the score of an  $ASx$ , the average score of all ASs having the same number of neighbours is subtracted.

$$P'_t(i) = P_t(i) - \frac{\sum_{j \in V, \#neighs(j) = \#neighs(i)} P_t(j)}{\text{card}(\{j \in V, \#neighs(j) = \#neighs(i)\})} \quad (6)$$

Furthermore, as this method considers potential future AS paths or BGP announcements, distinction between AS announcing an IP block or not in BGP-Ranking computation is not necessary anymore. Thus, the BGP-Ranking is slightly modified by not adding 1 in equation (1)<sup>5</sup>.

Finally, each AS is associated to new score value which can be used to establish another ranking which differs from BGP Ranking and which aims to figure out ASs tolerant to provide transit to malware hoster. Thus, once such a ranking is determined, ASs in the first ranks having an abnormal high score regarding the other ASs are suspects. This process is considered in the following experiments.

<sup>5</sup>this version is used in the rest of the paper

## VI. EXPERIMENTAL RESULTS

### A. Methodology

Some ASs are known to be tolerant with malware [15], [3], but that there is no simple mean to truly identify a malware transit AS, except after huge anti-cybercrime investigations [8]. The more, datasets are not publicly available.

It was never researched in the past and so there is no available public datasets. That is why the paper introduces a metric qualified as an estimation. However, our first definition in equation (3) is built from the natural definition, i.e. an AS is considered to transit malware based on its capacity to forward traffic between malware hoster and benign AS (the differentiation between ASs is done using BGP-ranking). Even if this metric cannot be directly computed, one experiment is dedicated to show that the estimated value based on the same natural definition (equation (4)) is quite stable, after the main parameter  $k$  (number of hops) reaches a certain value. The estimation given in equation (4) is a local metric as it has to be calculated individually for each AS, while PageRank is computed globally for the entire AS graph. The evaluation in following section checks that PageRank's output is compliant with the estimated metric for selected ASs as the latter cannot be computed for all ASs due to scalability issues. This issue is also discussed in last section about complexity.

In the experiment, ASs have been renumbered to mask their real identities. All used datasets are publicly available and our tool is available at <http://lorre.uni.lu/~jerome/files/asmatra.zip> for readers interested in experimenting themselves.

### B. Dataset

Our dataset was built from the Routing Information Service (RIS) of RIPE<sup>6</sup>. It consists of raw BGP-data (packets) collected from the rrc00.ripe.net collector in from 2012 April 17 to 30 in order to avoid any bias due to a particular route change at a certain day. 7243k AS paths have been observed during this period, but only 1028k are unique and correspond to 41k distinct ASs for which BGP-Ranking has been computed from the different blacklists (see section III-C). The theoretical measure presented in equation (3) would have to iterate over all 1028k paths and examine each intermediate AS on the paths. This would give a total number of 4150k iterations, the equivalent to the total number of intermediate links in the AS paths. Using the AS graph, the 41k ASs are interconnected only with about 95k edges only, which shows a higher scalability (with a factor of 44 regarding the number of connections to consider) while using the estimation approach in section IV-D.

The BGP-ranking (version where one is subtracted) varies between zero (no malware hosted) and 0.17565 with a mean of 0.00005. Hence, most ASs are not logically providing malware. Fig. 7 highlights BGP-ranking values where ASs have been indexed in logarithmic scale on the x-axis from high to low BGP-ranking values. Less than 4000 ASs have been identified as hosting malware, i.e. having a strictly positive BGP-ranking value. Thus, the increasing malicious activity in the Internet comes from a minority of ASs, which needs to obtain a good connectivity for being efficient. This strengthens the motivation of our approach to discover ASs providing

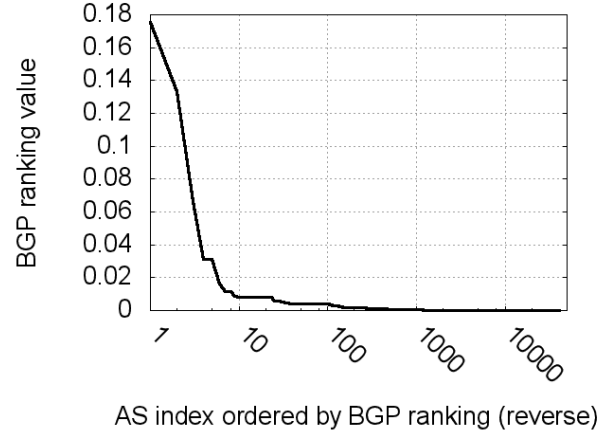


Fig. 7: BGP-ranking (ASs ordered from high to low BGP-ranking values)

such a connectivity. The BGP-ranking file is downloadable at <http://lorre.uni.lu/~jerome/files/bgpranking.txt>.

### C. PageRank Tuning

Explained in section IV, the advantage of using PageRank compared to the estimation in equation (4) is to avoid the evaluation of an AS by exploring a limited subgraph of neighbours reachable in  $k$  hops. PageRank also considers the entire graph by successive iterations and basically propagates BGP-ranking scores along edges. But the damping factor,  $d$ , needs to be defined. It balances the score computation between the connectivity (graph edges) and the BGP-ranking of nodes.

Except for the extrema values of  $d$  (0 and 1), tests have been performed from 0.1 to 0.9 by steps of 0.05. To measure the impact of  $d$ , the dispersion of the PageRank scoring for each AS is evaluated by using the variation coefficient that is the ratio between the standard deviation and the mean. As a relative metric, it is preferred to the standard deviation. The average value over all ASs is around 41% which is quite high.

In Fig. 8, the graph shows the scores in reverse order for different values of  $d$ . It clearly highlights that only a few of them (top 30) are really distinguishable.

### D. Validation

Due to the variability of results regarding  $d$ , an AS is considered to provide transit for malware, if it is always ranked in the top 30, independently of the value of  $d$  when executing PageRank. This gives a set  $T$  of 23 ASs. Fig. 9(a) shows the estimated value for the capacity to transit malware based on equation (4). Each line represents an AS and the first value is the BGP-ranking value. To get a comparative view, Fig. 9(b) represents the ASs always in the top 100, but out of top 30, independently on the damping factor  $d$ . This gives a set  $B$  of 30 ASs which mainly reach a value of zero within few hops.

Compared to Fig. 7, there is no correlation between hosting and providing transit, as some ASs out of the top 30 even have a higher BGP-ranking.

<sup>6</sup><http://www.ripe.net/data-tools/stats/ris/ris-raw-data>

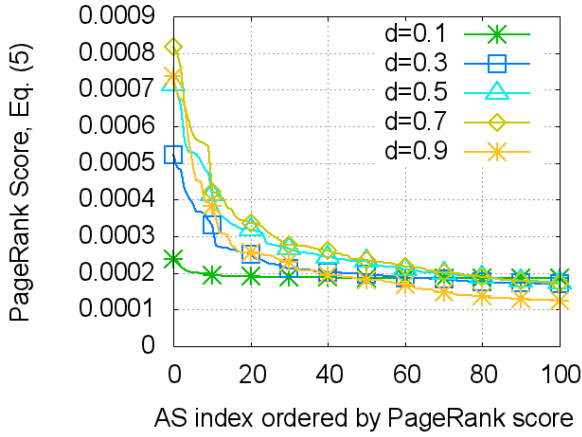


Fig. 8: Pagerank values for top ranked ASs for a range of damping factors  $d$

In Fig. 9, the curves show the same shape with higher values for the first hops with a peak in 1. In fact, by increasing the neighbourhood radius (hop count) in equation (4), it tends to merge connected components and, by this, reduce the potential number of malware transit possibilities. Thus, the malware transit ASs are connected to malware hosters within a few hops. This means that attackers do not need to be necessarily hidden behind multiple ASs for being efficient.

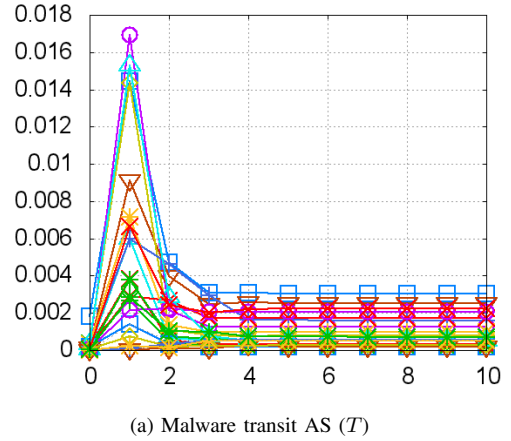
Moreover, Fig. 9 clearly shows that the PageRank approach obeys the properties of the estimated transit value (equation (4)) as the top 30 ASs have higher values (always  $> 0.0002$ ) than ASs out of the top 30 ( $< 0.00015$  except for one AS).

To illustrate our validation from a practical perspective, a concrete example is given in Fig. 10. It shows the neighbourhood for some ASs, T14 and T27, which are thought to provide transit to malware. Both were chosen out of  $T$  because they can easily be represented here while other ASs in  $T$  have more nodes.

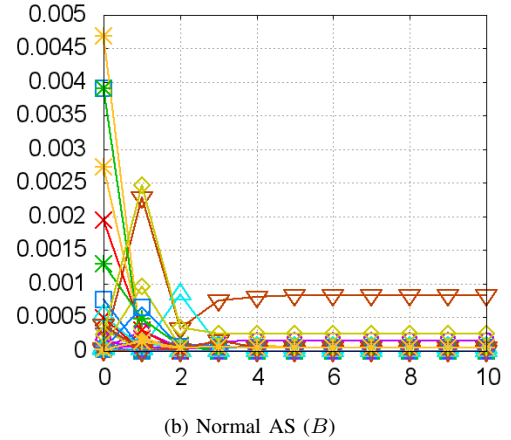
The lighter a node, the higher is its BGP-ranking. This remains also valid for its size. T14 and T27 are not large ASs. They are peered with several ASs and in particular they are close to AS 24 and 15, which are connected together and provide malware, as indicated by their size. Moreover, the AS 1 is directly connected to T14, while it is highly ranked by BGP-ranking. Thus, T14 and T27 play a major role for transiting malware. Because this is a manual investigation, we cannot applied exhaustively to all ASs.

### E. Complexity

The computation of the estimated value for every node, based on the neighbourhood in equation (4), requires to extract the subgraph constructed with nodes at a maximal distance  $k$  of each node. Extracting a subgraph for  $k$  hops is done by traversing the edges of the node to be analysed. Doing this operation for all nodes is achieved in  $O(n)$ , where  $n$  is the number of edges. Once achieved, the second step is the calculation of the connected components. This can be done in parallel with the computation of the aggregated BGP-ranking for each component. The task of determining a connected component is



(a) Malware transit AS ( $T$ )



(b) Normal AS ( $B$ )

Fig. 9: Theoretical estimation value for ASs tagged or untagged as malware transit by the PageRank based approach (each line represents a unique AS)

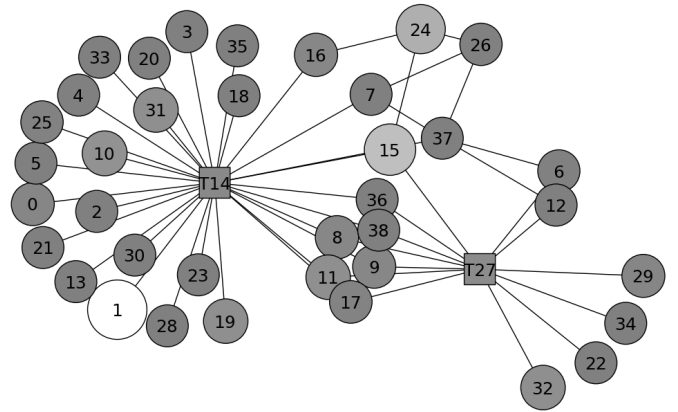


Fig. 10: Sample of an AS topology with malware transit ASs (node size and brightness are proportional to BGP-Ranking)

common in graph analysis and is achieved in linear time  $O(n)$ , where  $n$  is the number of edges. Thus, the overall complexity is  $O(\#nodes \times \#edges(neighborhood(k)))$  where  $\#nodes$  is the total number of AS and  $\#edges(neighborhood(k))$  is the average number of edges in the subgraph including nodes at a maximal distance  $k$  for any node.



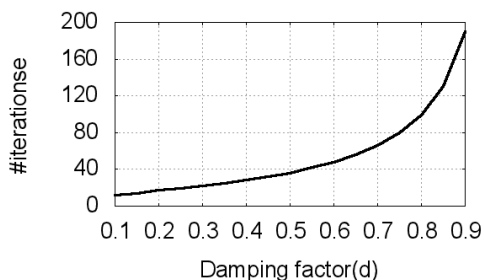


Fig. 11: Number of PageRank iterations



Fig. 12: Average number of nodes in a subgraph of neighbours

In PageRank, an iteration performs in linear time regarding the number of edges (score propagation on each edge). Whereas the total number of nodes and edges are constant in an AS graph and in the same order of magnitude in our experiments, as shown in section VI-B, the complexity differs according to the number of iterations (PageRank) or edges in the neighbourhood of a node (theoretical measure estimation). Fig. 11 shows the number of PageRank iterations depending of the damping factor  $d$ , while Fig. 12 shows the average number of edges in a subgraph of nodes for previously selected nodes for validation. This shows the advantage of the PageRank based approach from complexity view even for small values for  $k$ , since the average number of nodes for  $k = 2$  is 2839. This is why the estimation cannot be calculated for all ASs and only some of them have been selected from the results of the PageRank approach. In addition, a single PageRank iterations requires around 1.143 seconds in previous tests. Computing the estimated metric with  $k$  approaching the graph diameter was even not complete after more than 4 hours computation for only one AS. To conclude, combining both approaches is feasible by selecting some ASs (e.g. top 30 ASs) for a further check with the estimated metric.

## VII. RELATED WORK

Over the last years, a lot of approaches have been presented for the ranking of ASs based on various metrics, which resulted in a variety of ranking criteria. However, we are the first to make a global analysis, not limited to the first hop, for automatically revealing malware transit AS such as those involved in complex rogue network compositions.

Some relevant works done in AS ranking are the methods presented in [16], [17], which analyse the importance of the different ASs in Internet based on standard metrics like the

number of customers or their connectivity. An interesting work is presented in [9], where a new model for the ranking of ASs is presented by using inter-domain access volume, called IDAV. In this model, the contribution of an AS to the Internet routing is analysed by extending it with access volume, as for example carry and transit volume.

In [18], a new approach based on traceroute measurements for rating ASs is presented where quality of service metric are leveraged. A similar work is presented in [19]. The authors present a new model for constructing a network map that includes more data than only information about connectivity. For example, the authors extend their map by using latency and routing information. In a work from [20], the authors analyse the evolution of the Internet connectivity over ten years by analysing ASs and BGP data. In [21], the authors applied a method based on the power-law distribution for AS degrees to analyse the topology of the Internet hierarchy. A more customer oriented work is presented in [22], where the authors define customer-provider relations by applying a novel ranking algorithm, combining collaborative filtering and webpage ranking. In [15], the authors analyse the dynamics in the Internet by studying the properties of ASs by classifying time-scales of events. Considering security perspectives, the authors in [15] analysed the impact of some observable short-lived events.

Identifying abnormal ASs as source of malicious activities is mainly done using blacklist connectivity [3], [23], [5], [24] and some differ from additional features like BGP-behaviour [5] or botnet communications [24]. In [25], the authors show that the major volume of spam comes from few ASs only. Revealing complex structure of ASs involved in cyber-criminal organization require deep manual inspection, as shown in [8]. In [5], the authors study the neighbours at the first hop to figure out some ASs more prone to peer with malicious ASs. A complementary work to ours is [26] which evaluates the proximity of malware hosters in Internet.

## VIII. CONCLUSION

This paper presents a new method, called ASMATRA, that evaluates ASs on their capacity to provide transit for malware hosted in other ASs. ASMATRA leverages PageRank and BGP-ranking. We provide theoretical and practical measures for ASs deriving from the natural and intuitive definition of a malware transit AS. The results show that the PageRank-based method is coherent with these measures. It is more scalable, while applied globally to the entire network, since the other measures need to be applied individually to each AS. To conclude, we are able to track ASs used for malware transit without the AS being necessarily a malicious entity itself and which could not be detected by traditional investigation. In future work, it is planned to test other link analysis algorithms and investigate the evolution of malware transit capacity of ASs over time as well as using other BGP datasets.

**Acknowledgement:** this work was partially funded by IoT6, a European FP7 funded project under the grant agreement 288445, OUTSMART, a European FP7 project under the grant agreement 285038 and MOVE, a CORE project funded by FNR in Luxembourg. **Note:** figures include wikimedia contents (visit <http://commons.wikimedia.org> for license information).

## REFERENCES

- [1] J. Hawkinson and T. Bates, "Rfc 1930 - guidelines for creation, selection, and registration of an autonomous system," 1996.
- [2] Y. Rekhter, T. Li, and S. Hares, "Rfc 4271 - a border gateway protocol 4," 2006.
- [3] A. Dulaunoy, "BGP Ranking Project," in *32nd TF-CSIRT Meeting*. [Online]. Available: <http://www.terena.org/activities/tf-csirt/meeting32/dulaunoy-bgpranking.pdf>
- [4] P. Mockapetris, "Rfc 1034: Domain names - concepts and facilities," 1987.
- [5] C. A. Shue, A. J. Kalafut, and M. Gupta, "Abnormally Malicious Autonomous Systems and Their Internet Connectivity," *IEEE/ACM Trans. Netw.*, vol. 20, no. 1, pp. 220–230, 2012.
- [6] B. Krebs, "Major source of online scams and spams knocked offline," accessed on 04/02/12. [Online]. Available: [http://voices.washingtonpost.com/securityfix/2008/11/major\\_source\\_of\\_online\\_scams\\_a.html](http://voices.washingtonpost.com/securityfix/2008/11/major_source_of_online_scams_a.html)
- [7] J. Hruska, "Bad seed isp atrivo cut off from rest of the internet," accessed on 04/02/12. [Online]. Available: <http://arstechnica.com/security/news/2008/09/bad-seed-isp-atrivo-cut-off-from-rest-of-the-internet>
- [8] R. Howard, *Cyber Fraud: Tactics, Techniques and Procedures*. Auerbach Publications, 2009, ch. 5, The Russian Business Network: the Rise and Fall of a Criminal ISP.
- [9] Y. Wang, Y. Wang, M. Chen, and X. Li, "Information networking. towards ubiquitous networking and services." Springer, 2008, ch. Inter-Domain Access Volume Model: Ranking Autonomous Systems, pp. 482–491.
- [10] J. Rexford, J. Wang, Z. Xiao, and Y. Zhang, "Bgp routing stability of popular destinations," in *SIGCOMM Workshop on Internet measurement*. ACM, 2002.
- [11] J. Li, M. Guidero, Z. Wu, E. Purpus, and T. Ehrenkranz, "Bgp routing dynamics revisited," *SIGCOMM Comput. Commun. Rev.*, vol. 37, no. 2, pp. 5–16, 2007.
- [12] L. Page, S. Brin, R. Motwani, and T. Winograd, "The pagerank citation ranking: Bringing order to the web," 1999.
- [13] T. Haveliwala, "Efficient computation of pagerank," Tech. Rep., 1999.
- [14] J. François, S. Wang, R. State, and T. Engel, "Bottrack: tracking botnets using netflow and pagerank," in *Proceedings of the 10th international IFIP TC 6 conference on Networking - Volume Part I*, ser. NETWORKING'11. Springer, 2011, pp. 1–14.
- [15] M. Gaertler and M. Patriginani, "Dynamic analysis of the autonomous system graph," in *International Workshop on Inter-domain Performance and Simulation (IPS)*, 2004.
- [16] CAIDA, "Introduction to relationship-based as ranking," 2010, [http://www.caida.org/research/topology/rank\\_as/](http://www.caida.org/research/topology/rank_as/), last accessed: 08/03/2012.
- [17] FixOrbit, "Knodes index," <http://www.fixedorbit.com/>, last accessed: 08/03/2012.
- [18] L. Zimmerli, B. Tellenbach, A. Wagner, and B. Plattner, "Rating autonomous systems," in *International Conference on Internet Monitoring and Protection (ICIMP)*. IEEE Computer Society, 2009.
- [19] Z. Wang, J. Cheng, and S. Jamin, "Network maps beyond connectivity," in *Global Telecommunications Conference*. IEEE, 2005.
- [20] N. Alves, M. de Albuquerque, and J. de Assis, "Topology and shortest path length evolution of the internet autonomous systems interconnectivity," 2008, PhD Thesis, Universidade do Estado do Rio de Janeiro.
- [21] G. Siganos, M. Faloutsos, P. Faloutsos, and C. Faloutsos, "Power laws and the as-level internet topology," *IEEE/ACM Trans. Netw.*, vol. 11, no. 4, pp. 514–524, 2003.
- [22] L. Teow and D. Katabi, "Iterative collaborative ranking of customers and providers," 2006, <http://hdl.handle.net/1721.1/33234>.
- [23] A. J. Kalafut, C. A. Shue, and M. Gupta, "Malicious hubs: detecting abnormally malicious autonomous systems," in *Conference on Information communications (INFOCOM)*. IEEE, 2010.
- [24] B. Stone-Gross, C. Kruegel, K. Almeroth, A. Moser, and E. Kirda, "Fire: Finding rogue networks," in *Annual Computer Security Applications Conference*, ser. ACSAC. IEEE, 2009.
- [25] A. Ramachandran and N. Feamster, "Understanding the network-level behavior of spammers," in *Conference on Applications, technologies, architectures, and protocols for computer communications*, ser. SIGCOMM. ACM, 2006.
- [26] G. C. Moreira Moura, R. Sadre, A. Sperotto, and A. Pras, "Internet bad neighborhoods aggregation," in *IEEE/IFIP Network Operations and Management Symposium (NOMS 2012)*, April 2012.