

# Email Shape Analysis for Spam Botnet Detection

Paul Sroufe<sup>†</sup>, Santi Phithakkitnukoon<sup>†</sup>, Ram Dantu<sup>†</sup>, and João Cangussu<sup>‡</sup>

<sup>†</sup>Department of Computer Science & Engineering, University of North Texas, Denton, TX 76207-7102 USA  
{santi, prs0010, rdantu}@unt.edu

<sup>‡</sup>Department of Computer Science, University of Texas at Dallas, Richardson, TX 75083-0688 USA  
cangussu@utdallas.edu

**Abstract**—Botnets have become the major sources of spamming, which generates massive unwanted traffic on networks. An effective detection mechanism can greatly mitigate the problem. In this paper, we present a novel botnet detection mechanism based on the email “shape” analysis that relies on neither content nor reputation analysis. Shape is our new way of characterizing an email by mimicking human visual inspection. A set of email shapes are derived and then used to generate a botnet signature. Our preliminary results show greater than 80% classification accuracy (without considering email content or reputation analysis). This work investigates the discriminatory power of email shape, for which we believe will be a significant complement to other existing techniques such as a network behavior analysis.

## I. INTRODUCTION

Botnets are one of the world’s largest security threats [1, 2, 3] and botnet detecting and countering botnets are a growing field of research for the community. There are many other techniques for detecting spamming botnets. Some use network analysis and traffic generated from IRC commands to detect the presence of attacker control activity [4]. Another is uses a DNS blackhole to discover a botmasters reconnaissance of a particular bot’s “cleanliness” [5]. The Honeynet Project has also done work on collecting data from botnets by setting up honeypots for the sole purpose of being attacked [6].

These techniques rely on the difficult observation of botnets by intercepting network traffic. We have discovered a new way of identifying botnet behavior based on the content of the spam. This method employs a form of human intelligence to derive the shape of an email. Using the shape we can detect the presence of the spamming botnet that sent the email.

The rest of this paper is structured as follows. Section II describes the concept of the email shape, which can be used for the botnet detection. Section III presents some preliminary results of the proposed botnet detector. The paper is concluded in section IV with a summary and an outlook on the future work.

## II. EMAIL SHAPE BASED BOTNET DETECTION

The idea of email shape came during our extensive hand labeling of approximately 1,200 spam emails. These emails were hand labeled into buckets based on content, size, and email type (e.g. Plain, HTML, Multipart). After labeling several hundred of the emails, we started to see patterns emerge. We found evidence to support that botnet spammers used templates to bypass spam filters, and they would fill in the blanks with the links and info they needed to get through.

After seeing several of the templates emerge, we were able to use its shape to classify what bucket the email belonged to (An example of an actual spam template that we received by accident while collecting data for our corpus is shown in Fig. 1). Emails from the same template looked similar without having to read the content. They had different sizes (number of characters) occasionally but they still have the similar shape. Thus, we define the shape of an email as a shape that a human would perceive (e.g. shape of a bottle).

```
-----NextPart_001_2D49_73AC25235E4E77CE
Content-Type: text/plain
Content-Transfer-Encoding: quoted-printable

Company Name
Motto Here
=20
Dear Name,

Run the erranking resultsRun the user-friendly and technology driven Tool P=
rogram.Try the FREE 90 day trial and start achievingoutstanding search engine
e placement and ranking results. Run the user-friendly and technology drive=
n Optimization Tool Program.Try the FREE 90 day trial and outstanding search=
h engine placement and ranking resultsRun the user-friendly and technology=
driven Optimization Tool Program. Try the FREE 90 day trial and start ach=
eivingoutstanding search engine placement and ranking resultsRun the user-f=
riendly and techn ology driven

Sincerely,

John Smith
Manager Accounts
Company Software=20

Tel: your telephone
Fax: your fax
Web: your web site
=20
Copyright©Company Name.com
```

Figure 1. An example of an actual spam template.

Here we present an Email Shape based Botnet Detector (EsBod) that detects email botnet based on email shape analysis. EsBod takes an email spam and feeds it to a *Shape Generator* that extracts its “skeleton” that is simply a set of numbers of character count of each line in the HTML code of the email. Let  $L$  denote the number of lines in the email HTML code, and  $h_k$  denote the total number of characters (including spaces) in line  $k$ . Thus, a skeleton ( $H$ ) of an email can be defined as follows.

$$H = [h_1 \ h_2 \ h_3 \ \dots \ h_L]. \quad (1)$$

The “shape” of the email can then be drawn from its skeleton by applying a Gaussian kernel density estimator [7], which is given by Eq. (2).

$$f(x) = \frac{1}{Lw} \sum_{k=1}^L K\left(\frac{x-h_k}{w}\right), \quad (2)$$

where  $K(u)$  is the kernel function and  $w$  is the bandwidth or smoothing parameter. In our case,  $w$  is obtained by using AMISE optimal bandwidth selection based on Sheather Jones Solve-the-equation plug-in method [8], and  $K(u)$  is an Gaussian of zero mean and unit variance given by Eq. (3).

$$K(u) = \frac{1}{\sqrt{2\pi}} e^{-u^2/2}. \quad (3)$$

This work is supported by the National Science Foundation under grants CNS-0627754, CNS-0619871 and CNS-0551694.

The derived email shape is then fed into the *Classifier* that classifies an email spam to different botnets based on the Hellinger distance [9], which has been widely used for estimating a distance (difference) between two probability measures (e.g., probability density functions (pdf), probability mass functions (pmf)). Hellinger distance ( $d_H^2(P, Q)$ ) between  $P$  and  $Q$  can be computed by Eq. (4), where  $P$  and  $Q$  are two probability measures, which are  $M$ -tuple  $\{p_1, p_2, p_3, \dots, p_M\}$  and  $\{q_1, q_2, q_3, \dots, q_M\}$  respectively.  $P$  and  $Q$  satisfy  $p_m \geq 0$ ,  $\sum_m p_m = 1$ ,  $q_m \geq 0$ , and  $\sum_m q_m = 1$ . Hellinger distance of 0 implies that  $P = Q$  whereas disjoint  $P$  and  $Q$  yields the maximum distance of 1.

$$d_H^2(P, Q) = \frac{1}{2} \sum_{m=1}^M (\sqrt{p_m} - \sqrt{q_m})^2. \quad (4)$$

In our case,  $P$  can be a pmf of an email shape ( $s$ ), which can be obtained by taking  $M$  equal limit integrations over pdf given by Eq. (2) and  $Q$  can be an email signature of a botnet  $i$  ( $b_i$ ). The classification is therefore based on the minimum Hellinger distance as follows.

$$\text{Classified Botnet} = \arg \min_i d_H^2(s, b_i). \quad (5)$$

The system overview of the EsBod is illustrated in Fig. 1, which also shows a graphical example of a received email spam being classified as an email spam coming from botnet #2.

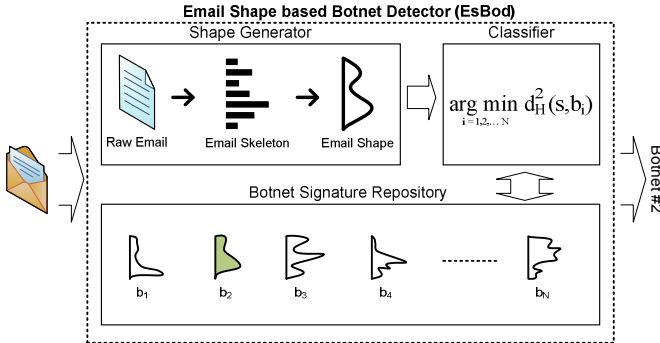


Figure 1. Basic architecture of EsBod.

By not considering the content of the email but focusing on its shape, EsBod is language independent. It is also feasible for matching the shapes of the emails of different sizes (e.g., different number of total lines in the emails) because of the normalization provided by the kernel density function.

### III. PRELIMINARY RESULTS

To evaluate the performance of EsBod, we used a corpus of approximately 1,200 spam emails with 45 hand labeled botnets. EsBod was trained with 1,100 emails to generate 45 botnet signatures (Fig. 2 shows some samples of these botnet signatures). After using 100 emails for testing, we were able to successfully classify testing emails to the correct botnets with high accuracy of 82% (Fig. 3 shows the accumulative accuracy rate as number of testing emails increases).

EsBod may not be the standalone solution but it can certainly complement other existing solutions. Especially, we believe that combining EsBod with a network behavior

analysis of spam botnets would improve the performance of the detector tremendously (further detail will be discussed in our future work).

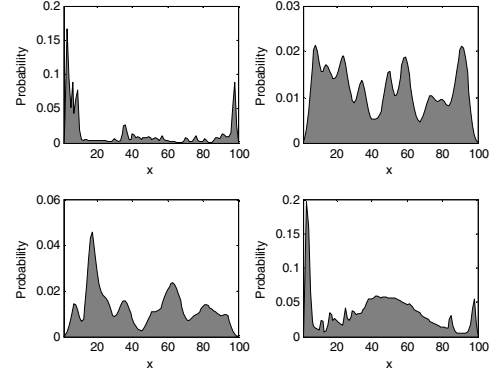


Figure 2. Sample botnet signatures in repository.

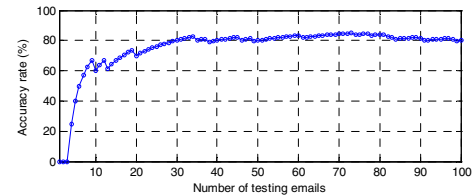


Figure 3. Accumulative accuracy rates of EsBod during the testing as number of testing spam emails increases.

### IV. CONCLUSION AND FUTURE WORK

In this paper, we present a step towards the ultimate goal of eliminating botnets, which are the major sources of malicious activities such as spamming, phishing, click fraud, and copy right violation. We propose here a novel method for detecting botnets based on email shape analysis (EsBod), which does not rely on either email content or reputation analysis. We are able achieve greater than 80% classification accuracy with our current model. Our future works will involve combining our shape analysis with other methods to track botnet activities, refining our shape analysis techniques by including other elements such as content and header information, and collecting a larger botnet email corpus.

### REFERENCES

- [1] M. Overton, *Bots and Boetnets: Reisks, Issues, and Prevention*, Virus Bulletin Conference, Dublin, Ireland, October 2005.
- [2] B. Schneier, *How Bot Those NETs?* Wired Magazine, July 27, 2006.
- [3] R. Narasine, Money Bots: Hackers Cash In on Hijacked PCs, eWeek, Sept. 2006.
- [4] W. Timothy Strayer, D. Lapsely, R. Walsh, and C. Livadas, "Botnet Detection Based on Network Behavior," in *Botnet Detection Countering the Largest Security Threat*, pp. 1-24, 2008.
- [5] A. Ramachandran, N. Feamster, and D. Dagon, "Detecting Botnet Membership with DNSBL Counterintelligence," in *Botnet Detection Countering the Largest Security Threat*, pp. 131-142, 2008.
- [6] P. Bäcker, T. Holz, M. Köster, G. Wicherski, "Know your Enemy: Tracking Botnets", *The Honeynet Project*, 13 March 2005. [Online]. Available: <http://www.honeynet.org/papers/bots/>.
- [7] E. Parzen, "On estimation of a probability density function and mode," *Annual Mathematic Statistics*, 33, vol. 3, pp. 1065-1076, 1962.
- [8] S. J. Sheather and M. C. Jones, "A reliable data-based bandwidth selection method for kernel density estimation," *Journal of the Royal Statistical Society, Series B*, 53, pp. 683-690, 1991.
- [9] G. L. Yang, and L. M. Le Cam, *Asymptotics in Statistics: Some Basic Concepts*, Berlin, Springer, 2000.