

Detection of Host Name Harvesting Attack in PTR Resource Record Based DNS Query Packet Traffic

YASUO MUSASHI,^{†1} DENNIS ARTURO LUDEÑA ROMANA,^{†2}
 SHINICHIRO KUBOTA^{†1} and KENICHI SUGITANI^{†1}

We statistically investigated the total inbound PTR resource record (RR) based DNS query request packet traffic to the top domain DNS server in a university campus network through January 1st to December 31st, 2009. The obtained results are: (1) We observed fourteen host name harvesting (HnH) attacks that we can observe rapid decreases in the unique source IP address based entropy of the inbound PTR RR based the DNS query packet traffic and significant increases in the unique DNS query keyword based one. (2) We found the consecutive and random IP addresses in the PTR RR based DNS query request packet traffic through the days of January 8th and 21st, 2009, respectively. Also (3), we calculated Euclidian distances between the observed IP address and the last observed IP address as the DNS query keywords and we detected two kinds of HnH attacks by a range of thresholds for 1.0-2.0 and 150.2-210.4. Therefore, these results show that we can detect more easily the inbound HnH attacks by calculating the Euclidian distances among the observed IP addresses in the inbound PTR RR based DNS query request packet traffic.

1. Introduction

It is of considerable importance to raise up a detection rate of bots, since they become components of the bot clustered networks that are used to transmit a lot of unsolicited mails including like spam, phishing, and mass mailing activities and to execute distributed denial of service attacks¹⁾⁻⁴⁾.

Wagner *et al.* reported that entropy based analysis was very useful for anomaly detection of the random IP search activity of Internet worms (IW) like an W32/Blaster or an W32/Witty worm, since the both worms drastically change entropy after starting their activity⁵⁾. Then, we reported previously that in the

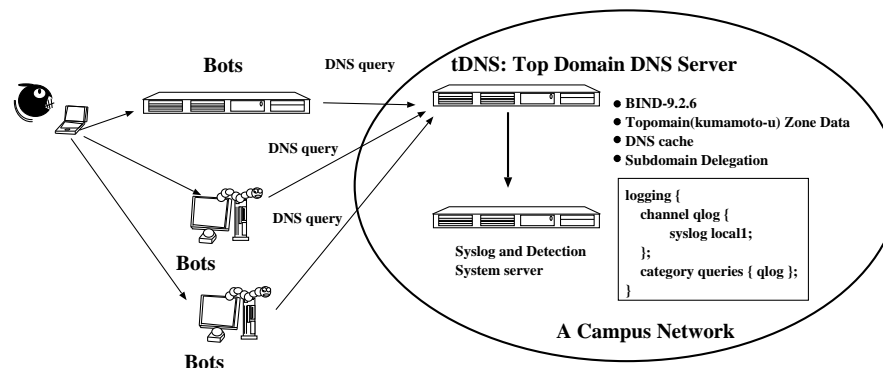


Fig. 1 A schematic diagram of a network observed in the present study. inbound PTR resource record (RR) based DNS query request packet traffic, the unique source IP address based entropy decreases considerably while the unique DNS query keyword based one increases when the host search (HS) activity is high^{6),7)}. The HS activity is recognized to be a pre-investigation activity or a harvesting activity of fully qualified domain names (FQDNs) of the university campus and/or enterprise networks *i.e.* hereafter, we call the HS activity as host name harvesting (HnH) attack. Probably, the attacker can concentrate to check out the vulnerabilities in the targeted servers or hosts by employing the HnH attack.

In this paper, (1) we carried out entropy and Euclidian distance based analyses on the total PTR resource record (RR) based DNS query request packet traffic from the Internet through January 1st to December 31st, 2009, and (2) we assessed the both results for entropy and Euclidian distance based analyses on the IP addresses as the query keywords in the PTR-RR based DNS query packet traffic.

2. Observations

2.1 Network Systems and DNS Query Packet Capturing

We investigated on the DNS query request packet traffic between the top domain (tDNS) DNS server and the DNS clients. Figure 1 shows an observed network system in the present study, which consists of the tDNS server and the

^{†1} Center for Multimedia and Information Technologies, Kumamoto University

^{†2} Graduate School of Science and Technology, Kumamoto University

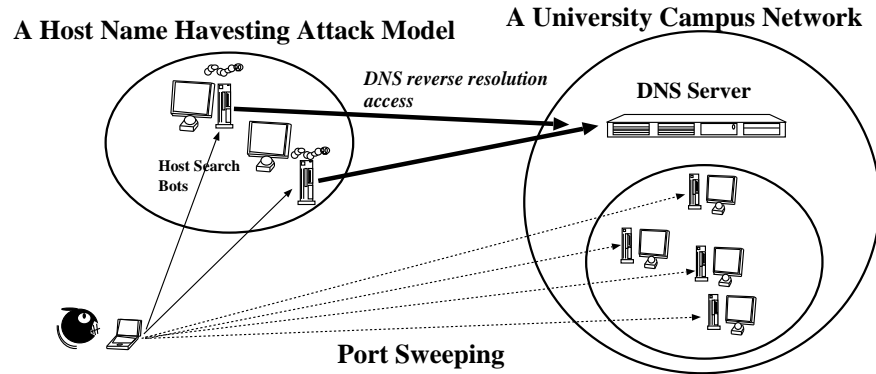


Fig. 2 A host name harvesting (HnH) attack model.

PC clients as bots like a host search bot or a spam bot in the campus or enterprise network, and the victim hosts like the DNS servers on the campus network. The **tDNS** server is one of the top level domain name (kumamoto-u) system servers and plays an important role of domain name resolution including DNS cache function and subdomain name delegation services for many PC clients and the subdomain network servers, respectively, and the operating system is Linux OS (CentOS 4.3 Final) in which the kernel-2.6.9 is currently employed with the Intel Xeon 3.20 GHz Quadruple SMP system, the 2GB core memory, and Intel 1000Mbps EthernetPro Network Interface Card.

In the **tDNS** server, the BIND-9.2.6 program package has been employed as a DNS server daemon⁸⁾. The DNS query packet and their query keywords have been captured and decoded by a query logging option (see Figure 1 and the named.conf manual of the BIND program in more detail). The log of DNS query packet access has been recorded in the syslog files. All of the syslog files are daily updated by the cron system. The line of syslog message consists of the contents of the DNS query packet like a time, a source IP address of the DNS client, a fully qualified domain name (A and AAAA resource record (RR) for IPv4 and IPv6 addresses, respectively) type, an IP address (PTR RR) type, or a mail exchange (MX RR) type.

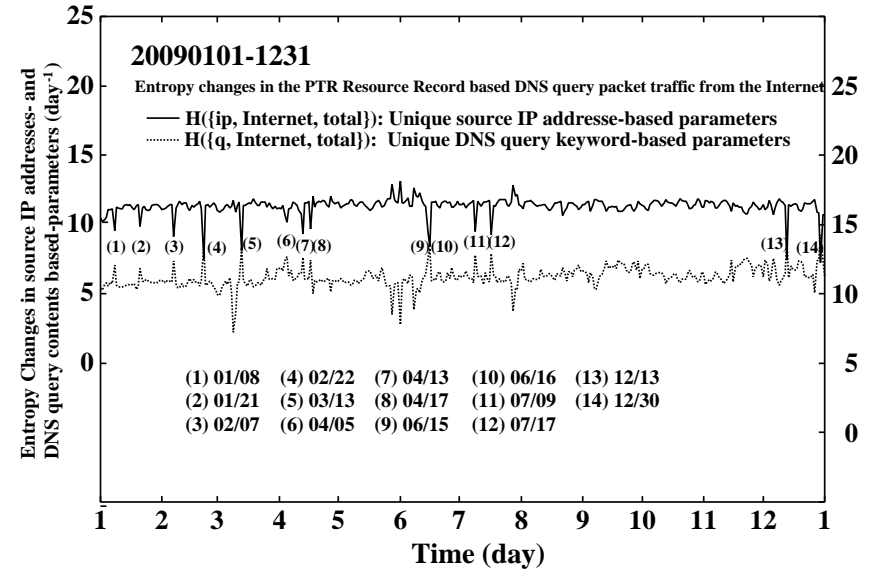


Fig. 3 Entropy changes in the total PTR-resource records (RR) based DNS query request packet traffic from the Internet to the top domain DNS (**tDNS**) server through January 1st to December 31st, 2009. The solid and dotted lines show the unique source IP addresses and unique DNS query keywords based entropies, respectively (day^{-1} unit).

2.2 Estimation of DNS Query Traffic Entropy

We employed Shannon's function in order to calculate entropy $\mathbf{H}(\mathbf{X})$, as

$$\mathbf{H}(\mathbf{X}) = - \sum_{i \in \mathbf{X}} \mathbf{P}(i) \log_2 \mathbf{P}(i) \quad (1)$$

where \mathbf{X} is the data set of the frequency $\mathbf{freq}(\mathbf{j})$ of a unique IP address or that of a unique DNS query keyword in the DNS query packet traffic from the Internet, and the probability $\mathbf{P}(i)$ is defined, as

$$\mathbf{P}(i) = \mathbf{freq}(i) / \left(\sum_j \mathbf{freq}(j) \right) \quad (2)$$

where i and j ($i, j \in \mathbf{X}$) represent the unique source IP address or the unique

```

Jan 8 14:53:15 kun named[10018]: client 66.***.***.#14981: query: 133.95.199.106 IN PTR
Jan 8 14:53:15 kun named[10018]: client 66.***.***.#36073: query: 133.95.199.107 IN PTR
Jan 8 14:53:16 kun named[10018]: client 66.***.***.#38238: query: 133.95.199.108 IN PTR
Jan 8 14:53:16 kun named[10018]: client 66.***.***.#21441: query: 133.95.199.114 IN PTR
Jan 8 14:53:16 kun named[10018]: client 66.***.***.#46887: query: 133.95.199.115 IN PTR
Jan 8 14:53:17 kun named[10018]: client 66.***.***.#22561: query: 133.95.199.116 IN PTR
Jan 8 14:53:17 kun named[10018]: client 66.***.***.#26389: query: 133.95.199.117 IN PTR
Jan 8 14:53:17 kun named[10018]: client 66.***.***.#45189: query: 133.95.199.118 IN PTR
Jan 8 14:53:17 kun named[10018]: client 66.***.***.#37530: query: 133.95.199.119 IN PTR
Jan 8 14:53:17 kun named[10018]: client 66.***.***.#12411: query: 133.95.199.128 IN PTR
Jan 8 14:53:17 kun named[10018]: client 66.***.***.#53097: query: 133.95.199.129 IN PTR
Jan 8 14:53:17 kun named[10018]: client 66.***.***.#14727: query: 133.95.199.130 IN PTR
Jan 8 14:53:17 kun named[10018]: client 66.***.***.#46036: query: 133.95.199.131 IN PTR
Jan 8 14:53:17 kun named[10018]: client 66.***.***.#32521: query: 133.95.199.132 IN PTR
Jan 8 14:53:17 kun named[10018]: client 66.***.***.#12967: query: 133.95.199.133 IN PTR
    
```

Fig. 4 Changes in the IP address as the DNS query keywords in the total PTR-resource records (RR) based DNS query request packet traffic from the Internet to the top domain DNS (tDNS) server at January 8th, 2009.

DNS query keyword in the DNS query packet, and the frequency $\text{freq}(\mathbf{i})$ are estimated with the script program, as reported in our previous work⁹⁾.

2.3 Host Name Harvesting Attack Model

We define here a host name harvesting (HnH) attack model (See Figure 2). *A host name harvesting attack (HnH) model* – the host name harvesting attack can be mainly carried out by a small number of IP hosts on the Internet or in the campus network like bot compromised PCs or like a directory harvesting attack. Since these IP hosts send a lot of the DNS reverse name resolution (the PTR RR based DNS query) request packets to the tDNS server, the unique IP addresses- and the unique DNS query-keywords based entropies decrease and increase, simultaneously.

Here, we should also define thresholds for detecting the HnH attack, as setting to 1,000 packets day⁻¹ for the frequencies of the top-ten unique source IP addresses or the DNS query keywords. The evaluation for threshold was previously reported⁹⁾.

2.4 Estimation of Euclidian Distances Among IP addresses as DNS Query Keywords

The Euclidian distances, $d(\mathbf{IP}_i, \mathbf{IP}_{i-1})$, are calculated, as

$$d(\mathbf{IP}_i, \mathbf{IP}_{i-1}) = \sqrt{\sum_{j=1}^4 (x_{i,j} - x_{i-1,j})^2} \quad (3)$$

where both \mathbf{IP}_i and \mathbf{IP}_{i-1} are the current IP address (i) and the last IP address ($i - 1$) of the DNS query keywords, respectively, and where $x_{i,1}$, $x_{i,2}$, $x_{i,3}$, and $x_{i,4}$ correspond to an IPv4 address like A.B.C.D, respectively. For instance, if an IP address is 192.168.1.1, the vector $(x_{i,1}, x_{i,2}, x_{i,3}, x_{i,4})$ can be represented as (192.0, 168.0, 1.0, 1.0). The detection is decided by thresholds $d_{\min} = 1.0$ and $d_{\max} = 2.0$, as

$$d_{\min} \leq d(\mathbf{IP}_i, \mathbf{IP}_{i-1}) \leq d_{\max} \quad (4)$$

3. Results and Discussion

3.1 Entropy Changes in Total PTR-RRs DNS Query Packet Traffic from the Internet

We demonstrate the calculated unique source IP address and unique DNS query keyword based entropies for the PTR-resource record (RR) based DNS query request packet traffic from the Internet to the top domain DNS (tDNS) server through January 1st to December 31st, 2009, as shown in Figure 3.

In Figure 3, we can find fourteen significant peaks and these peaks of (1)-(14) correspond to January 8th, 21st, February 7th, 22nd, March 13th, April 5th, 13th, 17th, June 15th, 16th, July 9th, 17th, December 13th, and 30th, 2009, respectively, in which all the peaks show significant increase and decrease in the unique source IP address- and the unique DNS query keyword based entropies, respectively. This result indicates that all the peaks (1)-(12) can be assigned to the HnH attack.

In the peak (1), at January 8th, 2009, we investigated the DNS query keywords in the total inbound PTR RR based DNS query packet traffic and the results are shown Figure 4. In Figure 4, we can view scenery that the IP address as DNS query keyword is consecutively incremented.

Therefore, it has a possibility that this consecutive increment of the IP address

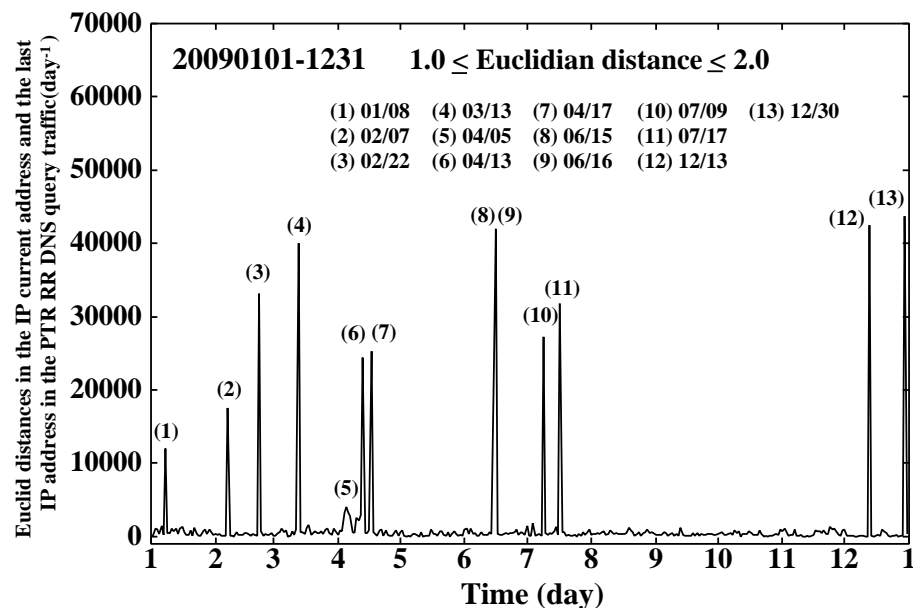


Fig. 5 Changes in Euclidian distance between the current IP address and the last IP address, as the unique DNS query keywords in the total PTR-resource records (RR) based DNS query request packet traffic from the Internet to the top domain DNS (tDNS) server through January 1st to December 31st, 2009 (day⁻¹ unit).

can be useful to detect the HnH attack in the PTR RR based DNS query request packet traffic.

3.2 Euclidian Distances in Consecutive Incremental DNS Reverse Queries

We illustrate the calculated Euclidian distance ($1.0 \leq d(\mathbf{IP}_i, \mathbf{IP}_{i-1}) \leq 2.0$) between the current IP address and the last IP address, as the unique DNS query keywords in the PTR resource record (RR) based DNS query request packet traffic from the Internet to the top domain DNS (tDNS) server through January 1st to December 31st, 2009, as shown in Figure 5.

In Figure 5, we can observe thirteen significant peaks (1)-(13) being allocated

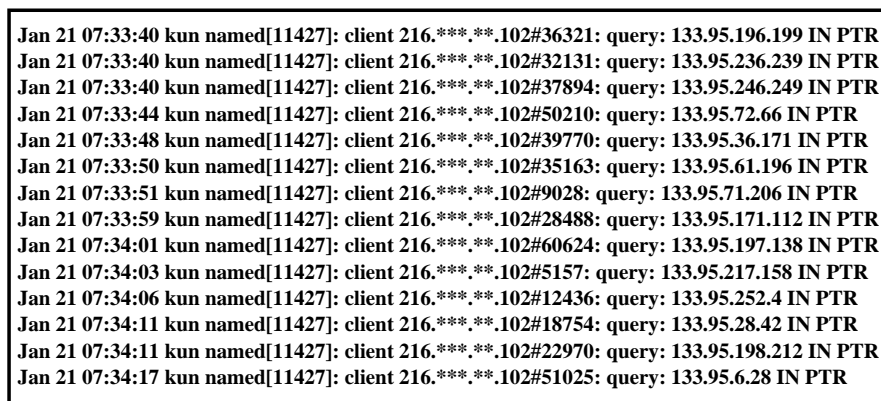


Fig. 6 Changes in the IP address as the DNS query keywords in the total PTR-resource records (RR) based DNS query request packet traffic from the Internet to the top domain DNS (tDNS) server at January 21st, 2009.

to January 8th, February 7th, 22nd, March 13th, April 5th, 13th, 17th, June 15th, 16th, July 9th, and 17th, 2009, respectively.

In Figure 3, we can observe the peak (2), corresponding to January 21st, 2009, however, we can find no peak for January 21st, 2009, in Figure 5. Thus, we investigated in more detail the DNS query keywords in the total inbound PTR RR based DNS query packet traffic at January 21st, 2009, and the results are shown Figure 6. In Figure 6, we can watch a scene that the IP address as DNS query keyword is discontinuously or randomly changed.

As a result, the discontinuous or randomized DNS query keywords cause a factor for disappearing of the peak at January 21st, 2009, in Figure 5.

3.3 Euclidian Distances in Random Sequential DNS Reverse Queries

The campus IP addresses are represented as $133.95.x_i.y_i$ in which both x_i and y_i can take numbers from 0 to 255, as: $0 \leq x_i \leq 255$ and $0 \leq y_i \leq 255$ *i.e.* the following eq 5 is obtained employing the both newly defined variables (x_i y_i) and eq 3, as:

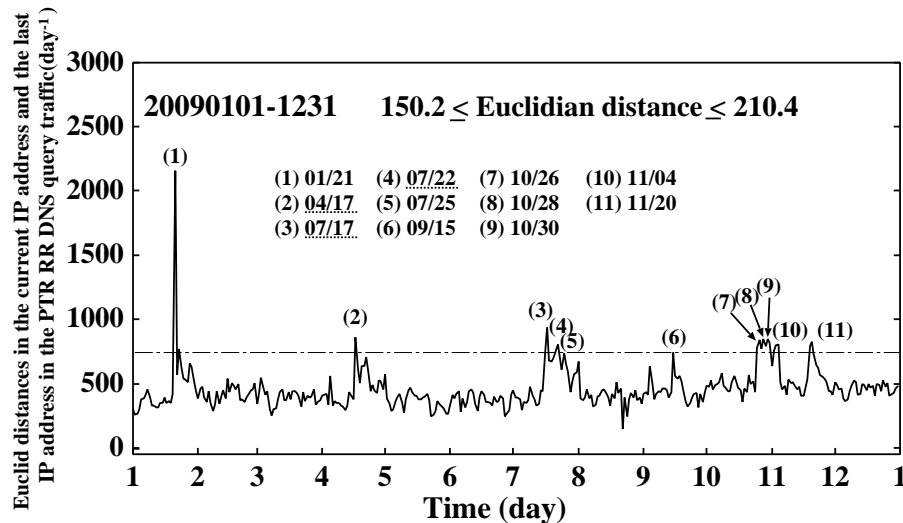


Fig. 7 Changes in Euclidian distance between the current IP address and the last IP address, as the unique DNS query keywords in the total PTR-resource records (RR) based DNS query request packet traffic from the Internet to the top domain DNS (tDNS) server through January 1st to December 31st, 2009 (day⁻¹ unit).

$$d(\mathbf{IP}_i, \mathbf{IP}_{i-1}) = \sqrt{(x_i - x_{i-1})^2 + (y_i - y_{i-1})^2} \quad (5)$$

where $(x_i - x_{i-1})^2$ or $(y_i - y_{i-1})^2$ takes a range from 0 to 255^2 *i.e.* the range of the Euclid distance, $d(\mathbf{IP}_i, \mathbf{IP}_{i-1})$, should be between 0.0 to $\sqrt{255^2 + 255^2}$ (~ 360.6).

If the both variables are random sequences, the Euclid distance, $d(\mathbf{IP}_i, \mathbf{IP}_{i-1})$, can also take a random sequence. Also, if the random sequence follows the Gaussian distribution, the probability for the Euclid distance takes a maximum value between at 180.3 ($\sim 360.6/2$) with a standard deviation of 30.1 ($\sim 360.6/12$) because of the central limit theorem *i.e.* the d_{\min} and d_{\max} should take values of 150.2 ($\sim 180.3 - 31.1$) and 210.4 ($\sim 180.3 + 31.1$).

We demonstrate the calculated Euclidian distance ($150.2 \leq d(\mathbf{IP}_i, \mathbf{IP}_{i-1}) \leq 210.4$) between the current IP address and the last IP address, as the unique DNS query keywords in the PTR resource record (RR) based DNS query request packet traffic from the Internet to the top domain DNS (tDNS) server through January 1st to December 31st, 2009, as shown in Figure 7.

In Figure 7, we can observe eleven significant peaks (1)-(11) being allocated to January 21st, April 17th, July 17th, 22nd, 25th, October, 26th, 28th, 30th, November 4th, and 20th, 2009, respectively. Especially, the peak (1) is very sharp and it is also assigned to January 21st, 2009, which is observed in Figure 3 but disappeared in Figure 5.

4. Conclusions

We investigated entropy and Euclidian distance based analyses on the total inbound PTR resource record (RR) based DNS query request packet traffic through January 1st to December 31st, 2009. The following interesting results are found: (1) we observed fourteen considerable host name harvesting (HnH) attacks in the entropy change in the PTR RR based DNS query request packet traffic, and (2) we found the thirteen consecutive incremental IP address- and the eleven random sequence of the IP address-based HnH attacks in the Euclidian distances between the current IP address and the last IP address in the PTR RR based DNS query request packet traffic.

From these results, it is concluded that we can detect the two typical HnH attacks by observing the Euclidian distance between the current and the last IP addresses in the total inbound PTR RR based DNS query packet traffic to the campus top level domain DNS server, from the Internet.

5. Acknowledgment

All the studies were carried out in CMIT of Kumamoto University and this study is supported by the Grant aid of Graduate School Action Scheme for Internationalization of University Students (GRASIUS) No. 165240040213 in Kumamoto University.

References

- 1) Barford, P. and Yegneswaran, V.: An Inside Look at Botnets, Special Workshop on Malware Detection, *Advances in Information Security*, Springer Verlag, 2006.
- 2) Nazario, J.: Defense and Detection Strategies against Internet Worms, 1 Edition; *Computer Security Series*, Artech House, 2004.
- 3) Kristoff, J.: Botnets, *North American Network Operators Group (NANOG32)*, Reston, Virginia (2004), <http://www.nanog.org/mtg-0410/kristoff.html>
- 4) McCarty, B.: Botnets: Big and Bigger, *IEEE Security and Privacy*, No.1, pp.87-90 (2003).
- 5) Wagner, A. and Plattner, B.: Entropy Based Worm and Anomaly Detection in Fast IP Networks, *Proceedings of 14th IEEE Workshop on Enabling Technologies: Infrastructure for Collaborative Enterprises (WETICE 2005)*, Linköping, Sweden, 2005, pp.172-177
- 6) Ludeña Romaña, D. A. and Musashi, Y.: DNS Based Analysis of DNS Query Traffic in the Campus Network, *Journal of Systemics, Cybernetics and Informatics*, Vol. 6, No.5, pp.42-44 (2008).
- 7) Ludeña Romaña, D. A., Kubota, S., Sugitani, K., and Musashi, Y.: Entropy Study on A and PTR Resource Record-Based DNS Query Traffic, *IPSJ Symposium Series*, Vol. 2008, No.13, pp.55-61 (2008).
- 8) BIND-9.2.6:
<http://www.isc.org/products/BIND/>
- 9) Ludeña Romaña, D. A., Musashi, Y., Matsuba, R., and Sugitani, K.: Detection of Bot Worm-Infected PC Terminals, *Information*, Vol. 10, No.5, pp.673-686 (2007).