

Traffic Analysis on a Domain Name System Server. SMTP Access Generates Many Name-Resolving Packets to a Greater Extent than Does POP3 Access

Yasuo Musashi,^{*} Ryuichi Matsuba,^{*,†} and Kenichi Sugitani[‡]

*Center for Multimedia and Information Technologies, Kumamoto University,
Kumamoto 860-8555 Japan, E-mail: musashi@cc.kumamoto-u.ac.jp*

[†]*Center for Multimedia and Information Technologies, Kumamoto University,
Kumamoto 860-8555 Japan, E-mail: matsuba@kumamoto-u.ac.jp*

[‡]*Center for Multimedia and Information Technologies, Kumamoto University,
Kumamoto 860-8555 Japan, E-mail: sugitani@cc.kumamoto-u.ac.jp*

Abstract: A traffic analysis on the domain name system (DNS) server of Kumamoto University was carried out with the multivariate statistical analysis. It is found that the total number of DNS packets, D_q , are generated from an electronic mail (E-mail) server, as represented: $D_q = m_{\text{SMTP}} N_{\text{SMTP}} + m_{\text{POP3}} N_{\text{POP3}}$, where N_{SMTP} and N_{POP3} represent the number of the simple mail transfer protocol (SMTP) access and that of the post office protocol version 3 (POP3) access, respectively. The linear coefficients m_{SMTP} and m_{POP3} are calculated to be 8.6 and 1.0. From these results, it is clearly concluded that the DNS access from the E-mail server is mainly driven by the SMTP access to a greater extent than the POP3 access. Also, it is found that m_{SMTP} is represented, as follows: $m_{\text{SMTP}} = 2 + 4n(1 - q)$, where q is a mail-receiving rate and n is a number of different domain hosts.

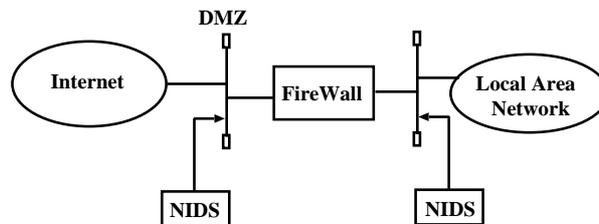
Keywords: multivariate statistical analysis, DNS access, SMTP access, POP3 access

1. Introduction

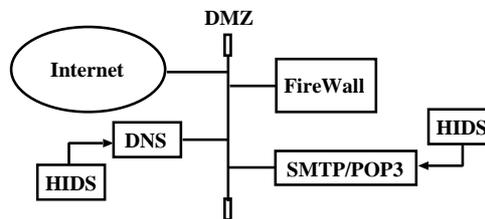
The domain name system (DNS)[1] is one of the most important services on the Internet/local area network (LAN) like a campus LAN system, since the recent network applications, such as simple mail transfer protocol (SMTP)[2] service, post office protocol version 3 (POP3)[3] service, file transfer protocol (FTP) [4] service, depend upon the DNS service, deeply. It is known that the DNS service provides us a host domain name, an internet protocol (IP) address[5], and how to exchange the electronic mail (E-mail). If the DNS service stops, a lot of network applications freeze or crash. Furthermore, remote attack and intrusion take place through an attack on the DNS server. From these points, we need to protect the DNS server, firmly.

Intrusion Detection System (IDS) is one of attractive solutions to protect security of the network server like the DNS server[6,7]. There are two kinds of IDS systems: (1) a network based IDS (NIDS) which directly checks the IP packets on a LAN like the demilitarized zone (DMZ) (Scheme 1A), and (2) a host-based IDS (HIDS) that checks only a network

(A) Network Based Intrusion Detection System



(B) Host Based Intrusion Detection System



Scheme 1

host servers (Scheme 1B). The latter IDS is suitable for the DNS server, since it is easily installed into the DNS server by use of a packet-logging program such as iplog[8,9].

There are two methods of detection of the intrusion: One is the scanning a signature file which is a database file of the remote attack pattern, and the other is de-

tection of the abnormality of the network server by checking packet logs with the statistical analysis. The former method needs to update frequently the signature file because new remote attack patterns or new cracking technologies have been developed quickly. However, the latter method does not need the signature file but can detect an unknown pattern of the remote attack. To develop a new effective IDS against remote attacks on the DNS server, it is of considerable of importance to get further detailed information of traffic of the DNS query packets (UDP packets) between the DNS server and the DNS clients[1].

In the present paper, we can present the statistical investigation on traffic of the DNS query packets between the DNS server (**1DNS**)[10] and the E-mail server (**1MX**)[11]. The traffic is schematically drawn in Scheme 2. Our purposes are (1) to compare both logs of SMTP and POP3 accesses with that of DNS query access, (2) to show how the DNS query packets are generated by the SMTP and POP3 accesses, and (3) to find out methods to detect abnormality of the DNS and the E-mail servers.

2. Computations

2.1 Multivariate Statistical Analysis

Multivariate statistical analysis[12] was carried out on traffic of the DNS query access between the **1DNS** and **1MX**. We set the total number of the DNS query access D_q as

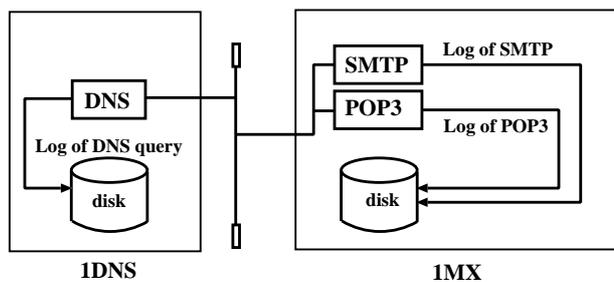
$$D_q = R_{\text{SMTP}} + R_{\text{POP3}} + R_{\text{FTP}} + \dots \quad (1)$$

where R_i is the access numbers from the DNS clients. The index i represents, here and hereafter, a network protocol, such as SMTP, POP3, FTP, and so on. We assume that R_i is a linear function of N_i

$$R_i = m_i N_i \quad (2)$$

Here, N_i is the access counts of a network application, *i.e.*, **1MX**. m_i is a linear coefficient. In **1MX**, $R_{\text{SMTP}} + R_{\text{POP3}} \gg R_{\text{FTP}} + \dots$, since the SMTP and POP3 daemons in the E-mail server are considered to be main DNS clients. Thus, D_q is given by

$$D_q = m_{\text{SMTP}} N_{\text{SMTP}} + m_{\text{POP3}} N_{\text{POP3}} \quad (3)$$



Scheme 2

As a result, the normal equation for multivariate regression analysis[12] is represented, as follows:

$$\mathbf{A}_{\text{SMTP,POP3}} \mathbf{x}_{\text{SMTP,POP3}} = \mathbf{d}_{\text{SMTP,POP3}} \quad (4)$$

where

$$\mathbf{A}_{\text{SMTP,POP3}} = \begin{bmatrix} \sum_{j=1}^n N_{\text{SMTP},j}^2 & \sum_{j=1}^n N_{\text{SMTP},j} N_{\text{POP3},j} \\ \sum_{j=1}^n N_{\text{SMTP},j} N_{\text{POP3},j} & \sum_{j=1}^n N_{\text{POP3},j}^2 \end{bmatrix}$$

$$(j = 1, 2, 3, \dots, n; \text{days})$$

$$\mathbf{x}_{\text{SMTP,POP3}} = \begin{bmatrix} m_{\text{SMTP}} \\ m_{\text{POP3}} \end{bmatrix}$$

$$\mathbf{d}_{\text{SMTP,POP3}} = \begin{bmatrix} \sum_{j=1}^n N_{\text{SMTP},j} D_{q,j} \\ \sum_{j=1}^n N_{\text{POP3},j} D_{q,j} \end{bmatrix}$$

2.2 Used Server Daemon Programs and Estimation of D_q and N_i

In **1DNS**, the BIND-9.1.3 program package have been employed as a DNS server daemon[13]. The log of DNS query packet have been recorded by the iplog-1.2 program[8,9] with the UNIX syslog system.

In **1MX**, the sendmail-8.9.3 program package[14] and the Qualcomm qpopper-4.0 program package[15] were installed as SMTP and POP3 server daemons, respectively. The log of SMTP and POP3 accesses have been observed in the syslog file. All of the syslog files are daily updated by UNIX crond.

The D_q , N_{SMTP} , and N_{POP3} values are estimated, as follows;

1. We connect to the DNS server (**1DNS**) by a ssh client, and then change into the “/var/log” directory. We enter the following commands:

```
% grep domain messages.1 >/tmp/1dns
```

After writing its output into a file at the “/tmp” directory, we count lines of the file by a “wc” command:

```
% grep "133.95.xx.yy:" /tmp/1dns | wc
```

The D_q value is given as an output of the wc command.

2. We connect to the E-mail (**1MX**), and then we change into the “/var/log” directory. We enter the following commands:

```
% grep "sendmail" syslog.0 >/tmp/1smtp
```

After using this command, we enter the next commands:

```
% grep "from=" /tmp/1smtp | wc
```

The N_{SMTP} value is given as an output of the wc command.

3. We enter the following commands to estimate the N_{POP3} value:

```
% grep "poppe\[\" syslog.0 | wc
```

3. Results and Discussion

3.1 Multivariate Statistical Analysis

The observed data in February, 2002 are shown in Table 1. From Table 1, $\mathbf{A}_{SMTP,POP3}$ and $d_{SMTP,POP3}$ of the normal equation are calculated, as follows:

$$\mathbf{A}_{SMTP,POP3} = \begin{bmatrix} 3.120 \times 10^8 & 9.084 \times 10^8 \\ 9.084 \times 10^8 & 2.652 \times 10^9 \end{bmatrix} \quad (5)$$

$$d_{SMTP,POP3} = \begin{bmatrix} 3.612 \times 10^9 \\ 1.052 \times 10^{10} \end{bmatrix} \quad (6)$$

From eqs (4), (5), and (6), $x_{SMTP,POP3}$ is obtained:

$$x_{SMTP,POP3} = \begin{bmatrix} 8.6 \\ 1.0 \end{bmatrix}$$

Therefore, eq (3) is rewritten as,

$$D_q = 8.6N_{SMTP} + N_{POP3} \quad (7)$$

Table 1. Observed data of N_{SMTP} , N_{POP3} , and D_q (day^{-1}).

j	N_{SMTP}	N_{POP3}	D_q
2002/02/11	1878	4480	26845
02/13	6010	17701	70327
02/14	5647	17663	68574
02/15	5744	16469	65849
02/17	1487	4004	18370
02/18	5973	16959	67262
02/19	5594	16118	62489
02/20	5666	17178	66718
02/21	5701	15851	63614
02/23	2363	6451	27540
02/24	1749	3814	20199
02/25	5731	16020	63626
02/26	5675	17688	68612

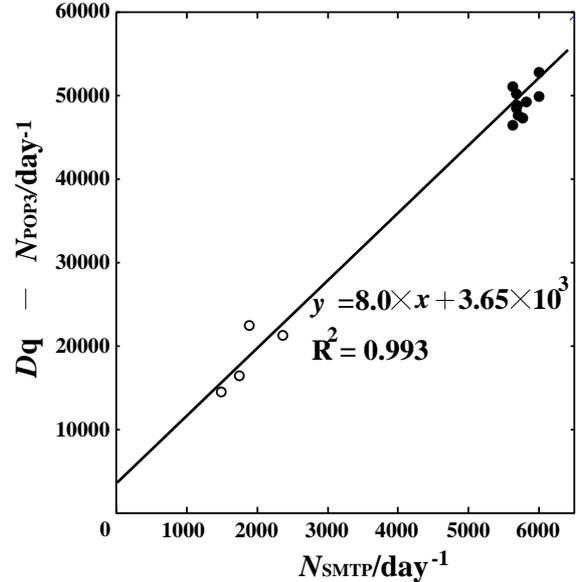


Figure 1. $D_q - N_{POP3}$ vs N_{SMTP} plot. Correlation coefficient (R^2) is 0.996.

The m_{SMTP} value of 8.6 is about 9 times greater than m_{POP} value of 1.0. This feature indicates that the SMTP access generates the DNS query access, rather than that of the POP3 access.

It is worthwhile to investigate graphically the relation between $D_q - N_{POP3}$ and N_{SMTP} , since it is likely that if the m_{POP3} value always is 1.0, the regression curve fits a linear function. Figure 1 illustrates regression analysis between $D_q - N_{POP3}$ and N_{SMTP} . The opened and closed circles show observed values in a holiday and a weekday, respectively. Expectedly, the correlation coefficient (R^2) is 0.993. These results indicate that m_{POP3} is 1.0 and that m_{SMTP} is about 8 ~ 9.

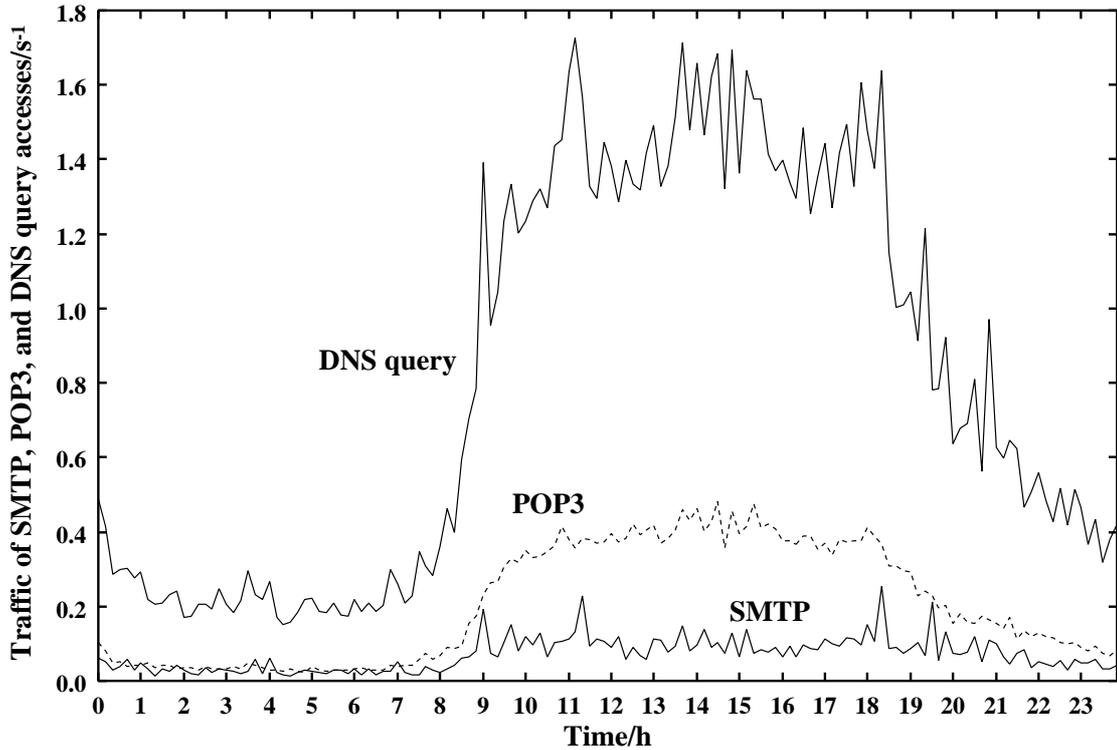


Figure 2. Traffic of the SMTP, POP3, and DNS query accesses in February 13th, 2002. The upper real line shows the DNS query access, the middle broken line means the POP3 access, and the bottom real line indicates SMTP access (s^{-1} unit).

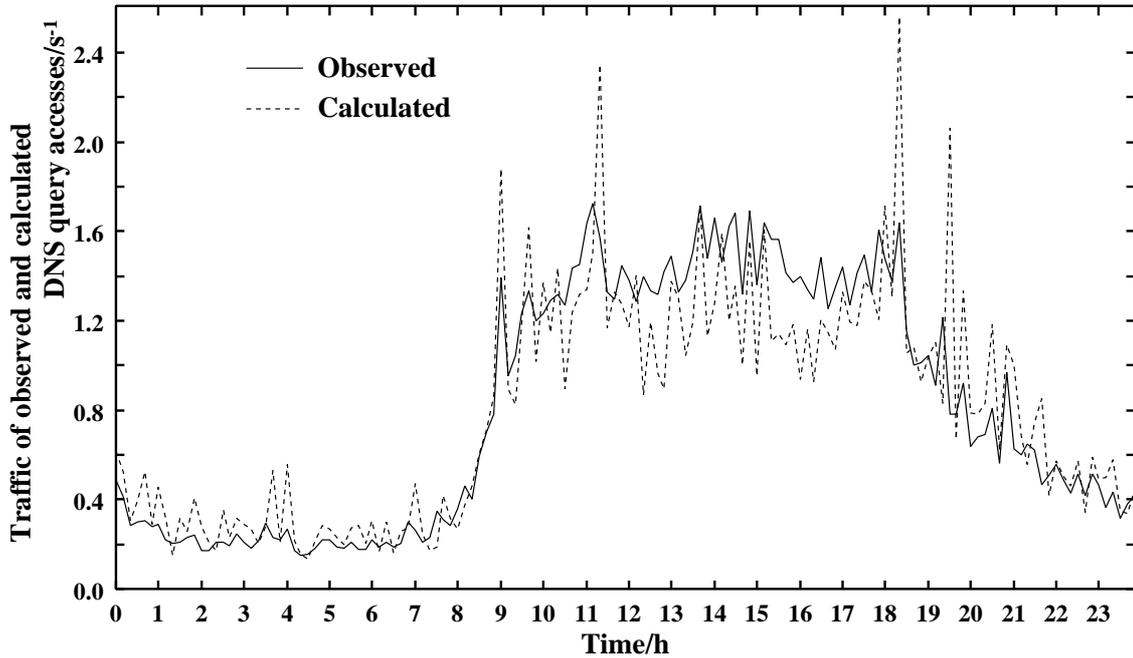


Figure 3. Traffic of the observed and calculated DNS query accesses in February 13th, 2002. The real line indicates the observed DNS query access, and the broken line means the calculated one (s^{-1} unit).

3.2 Analysis and Prediction of Traffic between the DNS server and the E-mail Server

We plot observed traffic curves of the DNS query access D_q (1DNS), the SMTP access (1MX), and the POP3 access (1MX) in Figure 2. The observation

was performed at February 13th, 2002.

In Figure 2, the traffic curve of D_q rises straight upon going from 08:00 to 09:00, considerably increases up to 11:00 with small fluctuations, slightly decreases to a local minimum at 12:00, repeats a local maximum twice, and decreases gradually to 08:00. These

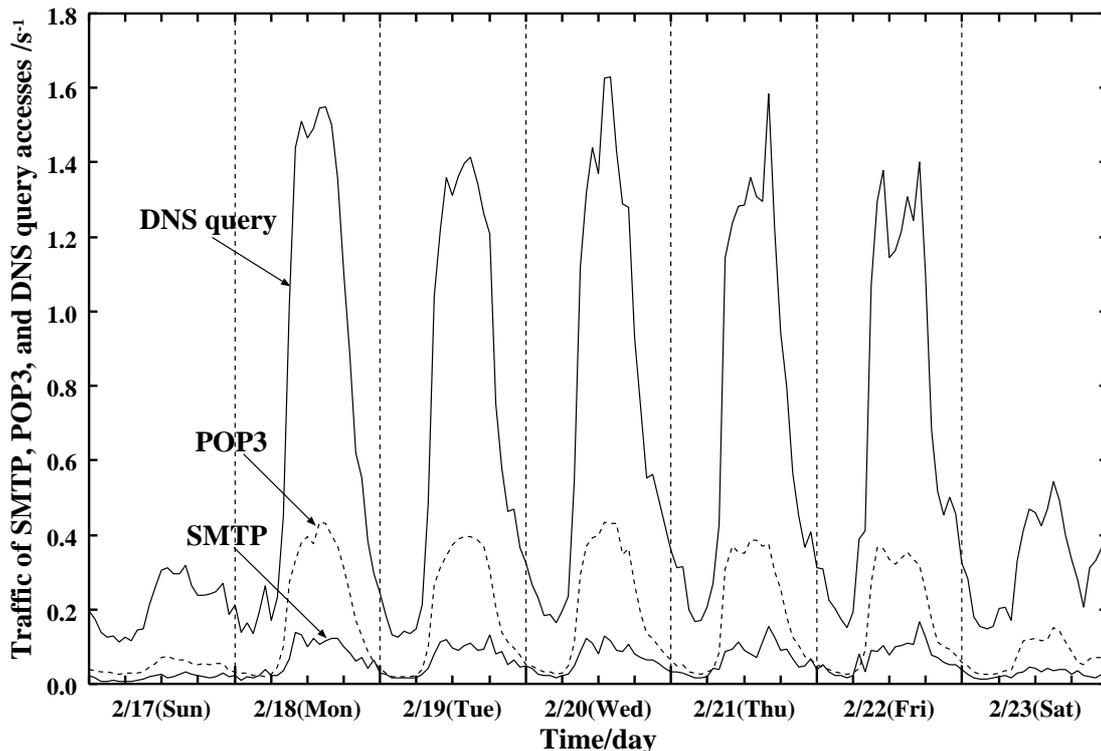


Figure 4. Traffic of the SMTP, POP3, and DNS query accesses in the third week of February, 2002. The upper real line shows the DNS query access, the middle broken line means the POP3 access, and the bottom real line indicates SMTP access (s^{-1} unit).

features indicate that almost users of **1MX** start to use an E-mail application in the morning, have a lunch at noon, and start to return back to home from 18:00.

The traffic curve of N_{POP3} changes in a mild manner and slightly resembles that of D_q but it is not so clear. This is probably because the N_{POP3} gives only small contribution to D_q (see eq (7)).

On the other hand, the traffic curve of N_{SMTP} resembles well that of D_q in a small scale manner. This is because the contribution of N_{SMTP} to D_q is a much greater extent than that of N_{POP3} .

It is worthwhile to compare the observed D_q with a calculated one. As shown in Figure 3, both observed and calculated curves of D_q resemble well each other. Especially, the calculated curve almost fits the observed one through 08:00-09:00. This result clearly shows that eq (7) can predict approximately D_q from N_{SMTP} and N_{POP3} .

In addition, the calculated curve through the night is slightly larger than the observed one. The calculated curve at the day time except large peaks is somewhat less than the observed one.

Figure 4 demonstrates traffic of D_q , N_{SMTP} , and N_{POP3} through a week of February 17th - 23th, 2002. The D_q and N_{SMTP} curves are similar to each other. Peaks of D_q curves emerge between 12:00 and 18:00 everyday. This feature is probably interpreted in terms of the following reasons; in the morning, almost

the users of **1MX** start to send E-mails, and in the afternoon, they receive the reply mail of the morning and send E-mails, again. In fact, these features appear in Figures 3 and 4: In the morning, the D_q curve fits well a linear function of the time. In the afternoon, on the other hand, the D_q curve fits a parabolic function of the time. Yoshida *et al.* reported that the correlation between input and output speed of traffic indicated a linear function[16], and the input function was a parabolic function of the time.

The traffic curves are classified into two kinds of curve: One is a weekday curve, and the other is a holiday curve. The curve of holiday D_q is a lesser scale than that of the weekday one. This is common situation because all the users of **1MX** work well through weekdays, and several users work on the holiday.

3.3 Why do both m_{SMTP} and m_{POP3} values become 8.6 and 1.0 ?

We need to know experimentally how many DNS query accesses, R_{SMTP} and R_{POP3} , are generated by one SMTP access and one POP3 access. The Scheme 3 sketches how many DNS query packets are requested by one SMTP access and one POP3 access.

Simply, one POP3 access generates only one DNS query packet. This result corresponds with the m_{POP3} value of 1.0, This feature is interpreted in terms that

the POP3 program calls only a resolver API such as `gethostbyaddr()` and `gethostbyname()` because of the reverse domain name resolving look-up. From this result, R_{POP3} is written as,

$$R_{\text{POP3}} = N_{\text{POP3}} \quad (8)$$

In the SMTP access, there are two kinds of SMTP access, as follows; one is a receiving SMTP access, and the other is a transmitting SMTP access. The former SMTP access yields two DNS query packets in an access (see Scheme 3A). Therefore, the DNS query packet of the receiving SMTP access, $R_{\text{SMTP}}^{\text{rec}}$ is represented,

$$R_{\text{SMTP}}^{\text{rec}} = 2N_{\text{SMTP}}^{\text{rec}} \quad (9)$$

The latter SMTP access generates two DNS packets when a SMTP client connects to the SMTP server and this SMTP server requests of the DNS server to get the canonical and reverse domain name resolving look-ups. In general, at least one E-mail destination address is required to transmit an E-mail. The E-mail destination address consists of fully qualified domain name (FQDN) or domain name (DN). As shown in Scheme 3B, the transmitting SMTP access needs to check four times: the first is to get a DN (checking the MX record), the second is to get a FQDN of the IP address of the SMTP destination, the third is to get a FQDN of the SMTP source to show the IP address of the SMTP source, and the fourth is to get the DN again. As a result, one E-mail destination address needs 4 DNS query packets, two E-mail destination addresses needs 8 DNS query packets, and n different E-mail destination addresses need $4n$ DNS query packets. From these reasons, the DNS query access of the transmitting SMTP access, $R_{\text{SMTP}}^{\text{tr}}$ is written, as

$$R_{\text{SMTP}}^{\text{tr}} = (2 + 4n)N_{\text{SMTP}}^{\text{tr}} \quad (10)$$

R_{SMTP} is

$$R_{\text{SMTP}} = R_{\text{SMTP}}^{\text{rec}} + R_{\text{SMTP}}^{\text{tr}} \quad (11)$$

Here, we set E-mail receiving rate,

$$q = \frac{N_{\text{SMTP}}^{\text{rec}}}{N_{\text{SMTP}}^{\text{rec}} + N_{\text{SMTP}}^{\text{tr}}} \quad (12)$$

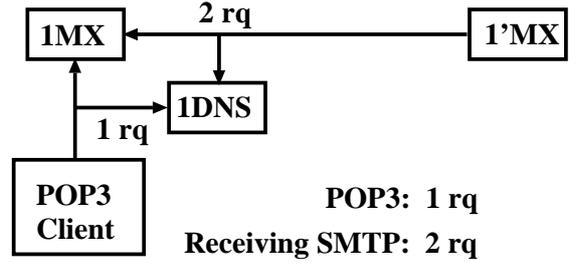
From this q and N_{SMTP} , $N_{\text{SMTP}}^{\text{rec}}$ and $N_{\text{SMTP}}^{\text{tr}}$ are rewritten as qN_{SMTP} and $(1 - q)N_{\text{SMTP}}$, respectively. From eqs (2) and eq (9)-(10), eq (11) is written as,

$$m_{\text{SMTP}}N_{\text{SMTP}} = 2qN_{\text{SMTP}} + (1 - q)(2 + 4n)N_{\text{SMTP}}$$

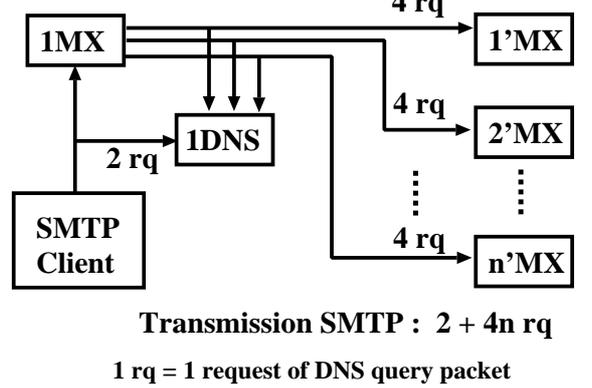
Since N_{SMTP} is an observed number *i.e.*, $N_{\text{SMTP}} > 0$, m_{SMTP} is rewritten by

$$\begin{aligned} m_{\text{SMTP}} &= 2q + (1 - q)(2 + 4n) \\ &= 2 + 4n(1 - q) \end{aligned} \quad (13)$$

(A) POP3 access and Receiving SMTP access



(B) Transmission SMTP access



Scheme 3

Therefore, D_q is represented, as follows:

$$D_q = (2 + 4n(1 - q))N_{\text{SMTP}} + N_{\text{POP3}} \quad (14)$$

If the q value is 0.50 ~ 0.75 and m_{SMTP} is 8.6, the n value is calculated to be 3.3 ~ 6.6. This means that the user of the E-mail server simultaneously sends a same mail to at least three different mail servers and that the user of 1MX sends to at least 3 ~ 7 persons by one E-mailing.

3.4 DNS Cache Effects on the DNS query access between the DNS and E-mail servers

As is well-known that DNS cache acts effectively to suppress heavy traffic of the DNS query from the E-mail server. However, quantitative analysis on the DNS cache effects have not been tried yet. We present the DNS cache effects of the DNS query access between 1DNS and 1MX with eq (7). The DNS cache server program is installed into 1MX with the BIND program package[13]. The DNS resolver configuration looks up only BIND in 1MX, in which root.cache is set to 1DNS (see Scheme 4). At 00:00 of March 11th, 2002, the DNS cache server program started. The

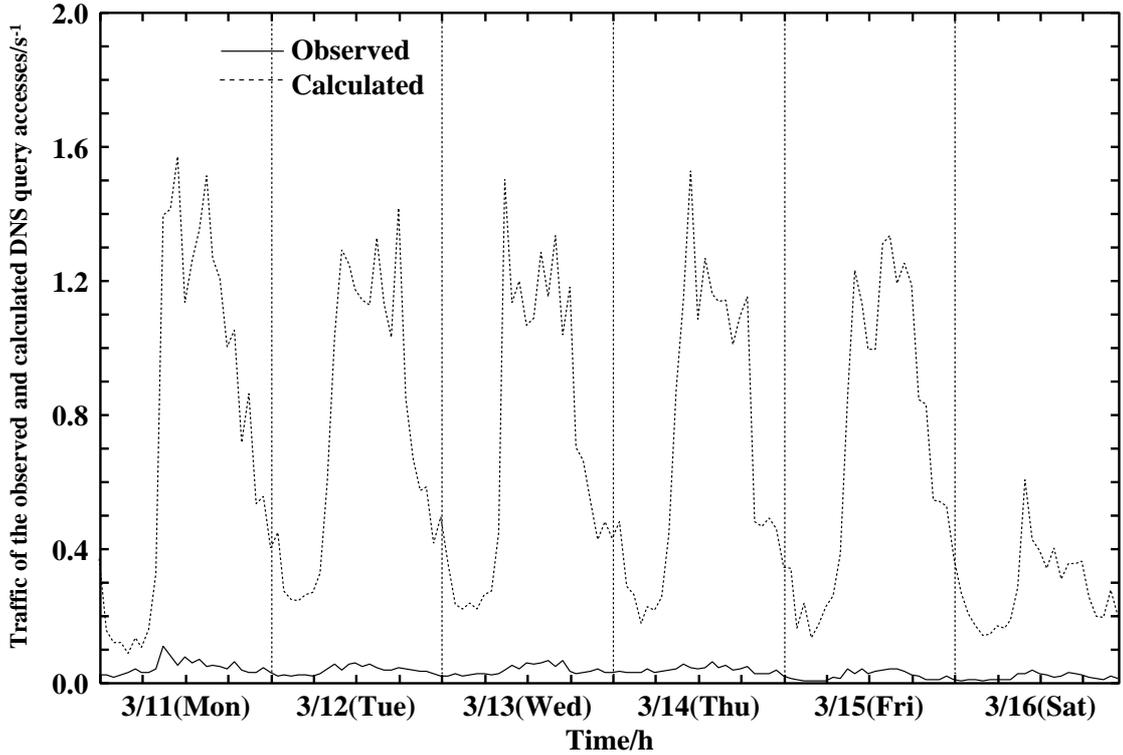


Figure 5. Traffic of the observed and calculated DNS query accesses upon going from March 11th to 16th (2002). The real line indicates the observed DNS query access, and the broken line means the calculated one (s^{-1} unit).

/etc/resolv.conf file was modified to change from an IP address of **1DNS** into an IP address of 127.0.0.1 (localhost).

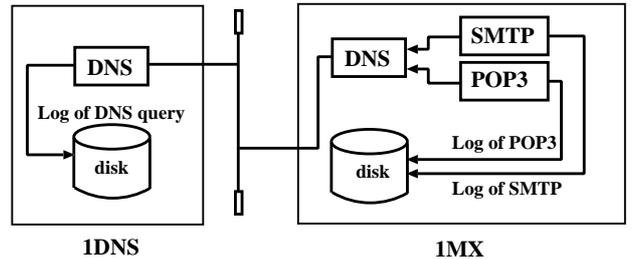
Figure 5 shows the observed (cached) and calculated traffic of D_q through from March 11th to 16th, 2002. Surprisingly, the observed traffic of D_q drastically decreases to a lesser extent than the calculated one. The DNS cache is effective to reduce the traffic between **1DNS** and **1MX**.

In March 11th, both observed and calculated traffic curves resembles in a small scale manner. In the other five days, the observed traffic curve is not similar to the calculated one. This is because the cache effects are so extremely high that the contribution of R_{SMTP} (N_{SMTP}) becomes obscure.

Here, we describe the DNS cache efficiency (DCE) as

$$DCE = 1 - \frac{D_q^{obs}}{D_q^{calc}} \quad (15)$$

In Figure 6, we plot DCE for same days in Figure 5. As shown in Figure 6, the DCE curve slightly moves to a local minimum at early morning in March 11th. After 06:00 in March 11th, the DCE curve slightly fluctuates about 0.85-0.99. From these results, it is clearly concluded that DCE is nearly 0.9 and that almost the traffic of the DNS query access between the DNS server and the E-mail server is cached by the DNS cache server at the E-mail server.



Scheme 4

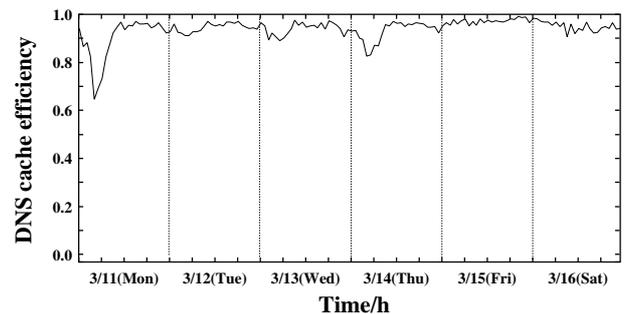


Figure 6. Changes in DNS cache efficiency upon going from March 11th-16th (2002).

4. Conclusions

Conclusions presented in this work are summarized as follows: (1) The total number of DNS packets, D_q , are represented, $D_q = m_{SMTP} N_{SMTP} + m_{POP3} N_{POP3}$, where N_{SMTP} and N_{POP3} represent the number of the

SMTP access and that of the POP3 access, respectively. (2) The linear coefficients m_{SMTP} and m_{POP3} are calculated to be 8.0-8.6 and 1.0. (3) Traffic between DNS server and E-mail server is mainly driven by traffic by the SMTP and the POP3 access, in which the traffic by the DNS query of the SMTP access ($m_{\text{SMTP}} = 8 - 9$) is about eight or nine times greater than that of the POP3 access ($m_{\text{POP3}} = 1.0$). (4) m_{SMTP} is represented, $m_{\text{SMTP}} = 2 + 4n(1 - q)$, where q is a mail-receiving rate and n is a number of different domain hosts. (5) The DNS cache sufficiently affects on the traffic between the DNS server and the E-mail server, and the cache efficiency is about 0.85-0.99. The DNS cache on the E-mail server reduces the traffic between the DNS server and the E-mail server, drastically.

It is of considerable importance to investigate the DNS query traffic generated by the E-mail server. Especially, it should be noted that the DNS query traffic is mainly generated by the SMTP access in the E-mail server. This information is essential to detect statistically a mass mailing internet worm (MMIW)[17] and to develop a new statistical IDS. This is because the MMIW has been diffused by way of the attachment file of the E-mail. Furthermore, several MMIW have a private SMTP engine and can attack on the DNS server, directly. If MMIW infection is expanded widely, the SMTP access becomes increases and a lot of DNS query packets are generated simultaneously. These features indicate that we can statistically detect infection of the MMIW and can know quickly a location of the MMIW-infected PC by only watching traffic between the DNS server and the E-mail server/the PC DNS clients. To get further information to develop a new statistics-based IDS (SIDS), a direct/indirect traffic between the DNS server and the MMIW is under further statistical investigation.

Acknowledgement. All the calculations were carried out with AMD Athlon, Intel Pentium III, and Sun Microsystems Ultra-Sparc machines in our center.

References and Notes

[1] Z. S. Su and J. B. Postel, "The Domain Naming Convention for Internet User Applications", RFC819, Network Information Center, SRI International, Menlo Park, California, 1982.

[2] J. B. Postel, "Simple Mail Transfer Protocol", RFC821, Network Information Center, SRI International, Menlo Park, California, 1982.

[3] M. T. Rose, "Post Office Protocol - Version 3", RFC1081, The Wollongong Group, Palo Alto, California, 1982.

[4] J. B. Postel and J. K. Reynolds "FILE TRANSFER PROTOCOL (FTP)", RFC959, University of the Southern California/ Information Sciences Institute, California, 1985.

[5] J. B. Postel, "DoD Standard Internet Protocol", RFC760, University of the Southern California/ Information Sciences Institute, California, 1980.

[6] S. Northcutt and J. Novak, "Network Intrusion Detection", 2nd ed; New Riders Publishing: Indianapolis, 2001.

[7] K. Yamamori, "An Improvement of Network Security Using an Intrusion Detection Software", *Journal for Academic Computing and Networking*, No.4, pp.3-13, 2000.

[8] eric@ojnk.nu, <http://tower.zot.nu/%7Eric/>

[9] <http://www.st.ryukoku.ac.jp/~kjm/security/-memo/1999/07.html>

[10] **1DNS** is the secondary DNS server of the Kumamoto University (kumamoto-u) which is run by our center. The OS is Linux OS (kernel-2.4.16), and the AMD Athlon 1.4 GHz.

[11] **1MX** is our mail server of the generic domain name of the Kumamoto University (kumamoto-u). The OS is Solaris 2.6 (Ultra-SPARC 300MHz, Sun Microsystems Inc.).

[12] P. R. Bevington, "Data Reduction and Error Analysis for the Physical Science"; McGraw-Hill: New York, 1967.

[13] <http://www.isc.org/products/BIND/>

[14] <http://www.sendmail.org/>

[15] <http://www.eudora.com/qpopper/>

[16] H. Yoshida, T. Suzuzki, and T. Iyama, "Traffic Analysis on a Campus Information Network", *Journal for Academic Computing and Networking*, No.4, pp.75-78, 2000.

[17] <http://www.symantec.com/region/jp/sarcj/re-fa.html>