

PC

-

The Research on Web-mail Server System
using a Cluster of Linux PCs

2001 2

PC

-

The Research on Web-mail Server System
using a Cluster of Linux PCs

2000 10

2000 12



'/etc/passwd'

'procmail'

MDA

가

SMTP

MDA

가

PC

-

가

SMTP

가

가

:

, , ,

,

: 99419 - 542

1.	1
2.	4
2.1	4
2.1.1	4
2.1.2	Sendmail program.....	7
2.1.2.1	'sendmail'.....	7
2.1.2.2	'sendmail'.....	8
2.1.2.3	'sendmail'.....	8
2.1.3	9
2.1.3.1	Cluster mailhub.....	10
2.1.3.2	Duke.....	10
2.1.3.3	Earthlink.....	11
2.1.3.4	Porcupine Scalable Mail Server.....	13
2.2	15
2.2.1	Round Robin DNS (RRDNS).....	15
2.2.2	Network Address Translation (NAT).....	16
2.2.3	Tunneling.....	17
2.2.4	Direct Routing (DR).....	17
2.2.5	Broadcasting & Filtering (BF).....	18
3.	19
3.1	19
3.1.1	MTA (Mail Transfer Program).....	20
3.1.2	MDA (Mail Delivery Agent).....	20
3.1.3	mail server program.....	21
3.1.3.1	node thread.....	21
3.1.3.2	operation thread.....	22
3.2	-.....	24
3.2.1	-.....	24
3.2.2	mail client program.....	24
3.2.3	authtake thread, authgive thread.....	25
3.2.4	datatake thread.....	26
3.2.5	operation thread.....	26

4.	27
4.1	27
4.2	28
4.3	30
4.3.1	30
4.3.2	30
5.	33
	34
Abstract	36
	37

1.	5
2. -	19
3.	21
4. mail client program	25
5.	28
6.	32

132

1.

가
가 [1],
Earthlink
가 1996 1 25,000 1997 9 350,000
가 1,000,000 [2],
Hotmail 75,000,000 .1)
AOL, Hotmail 1998
10,000,000
[3], 가
가
PC 가 PC
, network 10Mbps 100Mbps fast ethernet,
Gigabit ethernet
PC
가
가
가
가

1) 2000 Hotmail

가 , 가 , PC
가 , PC
가 PC
가 가 .
가

[3], 가 가

가 [4].
PC -

MDA (Mail Delivery Agent)²⁾ (mailbox)
(mailbox)

MDA (Mail Delivery Agent)

- 가 - .
-

가 MDA (Mail Delivery Agent)
SMTP connection
PC

. 2

2) '2.1.1 , .

가

. 3

. 4

, 5

.

2.

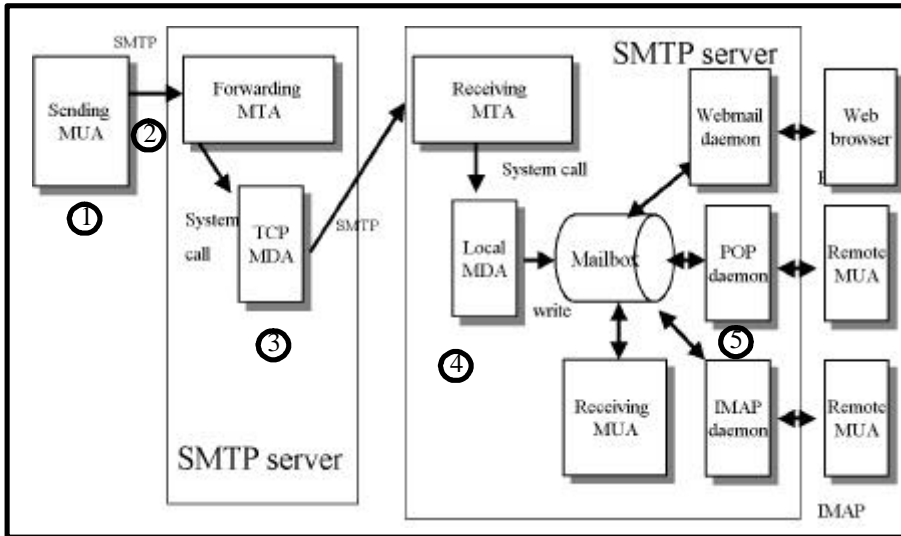
linux
MDA (Mail Delivery Agent)가
(mailbox)
MDA
MDA
MTA (Mail Transfer
Agent) 'sendmail'
MDA MTA
MTA
가

2.1

가
MTA (Mail Transfer Agent) 'sendmail' [5]
가
'sendmail'

2.1.1

가 가



1

[6]

가 'Outlook Express'

MUA(Mail User Agent) (1 'sending MUA':)

MUA(Mail User Agent) 가

Microsoft 'Outlook Express', Netscape 'Netscape mail', UNIX '/usr/ucb/mail' [6].

MUA

가 SMTP(Simple Mail Transfer Protocol)

'sendmail' MTA(Mail Transfer Agent)

SMTP (

). MTA(Mail Transfer Agent)

MDA(Mail Delivery Agent) MTA

[6],

'sendmail' [5], 'qmail' SMTP (Simple

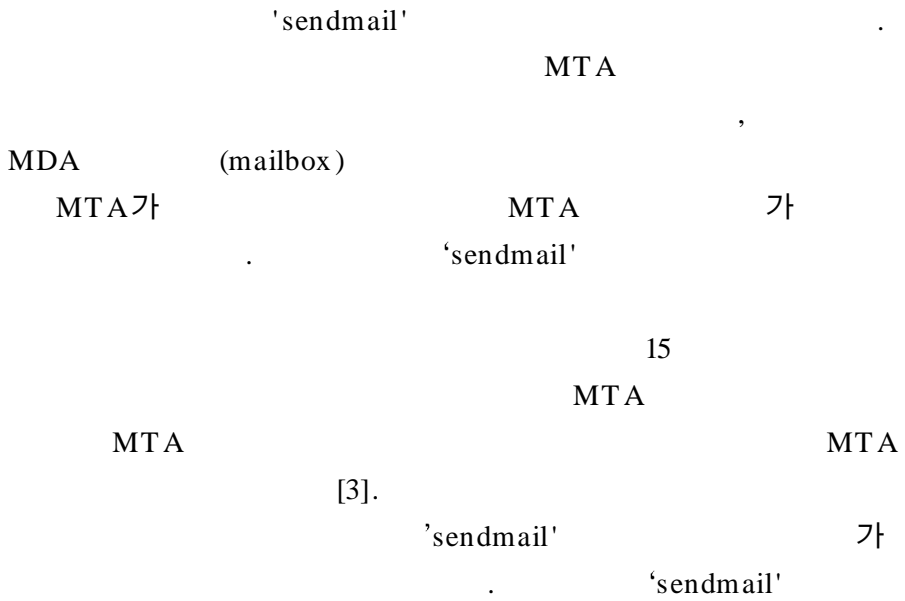
Mail Transfer Protocol)[7] MTA

MTA
 가 ,
 TCP MDA (Mail Delivery Agent) (TCP
 network) . TCP MDA 가
 SMTP MTA SMTP
 (). MDA (Mail
 Delivery Agent) MTA가
 (TCP, uucp)
 MDA가 '/bin/sh' [6].
 가 SMTP MTA
 local MDA
 (). local MDA
 MDA , 'procmail', '/bin/mail'
 . (mailbox)
 (mbox)
 (MH,
 maildir) 가 가 [6].
 가
 SMTP 'mail' MUA
 , POP (Post Office Protocol), IMAP
 (Internet Message Access Protocol)
 (). POP (Post
 Office Protocol)
 ,
 POP3[8] 가
 [9][10]. IMAP (Internet Message
 Access Protocol)
 ,
 IMAP4[11] POP 가
 [9][10].

1

가

2.1.2 Sendmail



2.1.2.1 'sendmail'

sendmail[5][12]

가

MTA 'sendmail' 가 .
 , 'sendmail'
 network
 UNIX, linux .
 PC
 , sendmail 'sendmail.cf'
 가 .
 sendmail sendmail
 , sendmail MDA
 MDA
 가 가 .

2.1.2.2 'sendmail'

'sendmail' temporary file,
 process forking, lock-file 가
 [4]. temporary file sendmail
 queue 가 qf-, df-
 overhead [13]. process forking
 sendmail SMTP process fork
 thread process overhead가
 . lock-file
 sendmail queue qf- lf-
 lock lock linux

sendmail.net [13]

queue 가

DNS .

2.1.2.3 'sendmail' [14]

```

SMTP      'sendmail'      . 'sendmail'
        fork
        'sendmail'
        queue
        load average      1
        process      sendmail      load
        average(la)      QueueLA (x)
        'msgpri > q / ( la - x + 1 )'
        sendmail
        queueing
        refuseLA      load average(la)가
        SMTP      .      QueueFactor (q)
        default 600,000      QueueLA (x)      default 8
        msgpri      message priority      'msgpri = size - ( class
        * z ) + ( recipients * y )'
        , recipients      class
        default 0, ClassFactor(z)      default 1800, RecipientFactor(y)
        default 30,000      . refuseLA      default 12

```

2.1.3

'sendmail'

```

        'sendmail'
        local file system
        가
        'Cluster mailhub      '[15], 'Duke

```


'[1], 'EarthLink Server' [3]

'[2], 'Porcupine Scalable Mail

2.1.3.1 Cluster mailhub [15]

Cluster mailhub NFS
 mail spool scalable
 가 가
 spool file 가 가 server directory
 name cache mailhub nfsd
 process 가 가 NFS
 [1]. 가 SMTP
 sendmail fork
 nfsd process load 가 ,
 SMTP 가 가 가
 process 가
 load average가³⁾ sendmail load average
 SMTP connection 가
 process load average가
 sendmail SMTP
 SMTP
 가 [1].

2.1.3.2 Duke [1]

Duke [15] mailhub mail
 gateway system .
 DNS MX record mail gateway
 SMTP mail
 gateway mailhub .

3) '2.1.2.3 sendmail

. dataless POP
 POP daemon mailbox location . File
 NFS
 '/var/mail#' mount , balanced hash
 mail queue
 가 . Authentication DB passwd
 file
 SQL database passwd file MDA
 getpwnam() SQL query
 mailbox , 가
 . 가
 , OS , dataless POP
 가 . switched FDDI
 , dataless POP server , 가
 . 가 RAID
 . 가
 DNS round robin .
 file locking open system call O_EXCL flag
 lock 가 15
 synchronous NFS operation 가
 . DB
 , 가
 . 가
 가 DB
 . SMTP POP
 lock file

[3].

35 mailbox
 20 . POP
 20 POP connection 가 600
 daemon 1000 POP connection
 .
 DB
 spool hash
 NFS 가 .
 ,
 .

2.1.3.4 Porcupine Scalable Mail Server [3]

Porcupine Scalable Mail Server DNS A, MX record
 SMTP, POP, IMAP socket client .
 SMTP, POP, IMAP session
 DB .
 .
 가 POP IMAP
 .
 user directory
 . 가
 distributed membership protocol user
 directory .
 DB 가

network
 DB 가
 가 DB
 parallel RPC
 가 .
 user directory
 NFS gateway daemon .
 directory
 OS .
 가
 가
 ,
 mailing list DNS
 link
 .
 Myrinet linux pthread
 prototype simulation
 scalable 가
 10 650 가
 network 가
 1Gb/sec bandwidth가 650
 node Disk가 가
 bottleneck disk 가가
 .
 NFS NFS
 broadcast
 network traffic .

2.2

()
IP
(High Performance)
(Scalability) , IP
가 (High Availability)
(Scalability) IP

HTTP

Round Robin DNS (RRDNS)[17],
Network Address Translation (NAT)[18], Tunneling[19], Direct
Routing (DR)[19], Broadcast & Filtering (BF)

2.2.1 Round Robin DNS (RRDNS)

Round Robin DNS (Domain Name Service) NCST
가 . DNS
가 IP
. RRDNS IP
, IP
IP
BIND (Berkeley Internet Name Domain) DNS

, RRDNS
 .
 , 가 DNS
 , DNS 가
 , 가 IP
 가
 . , DNS RRDNS
 DNS 가 , DNS
 가 ,

2.2.2 Network Address Translation (NAT)

Network Address Translation (NAT) IP
 . 가
 IP , (Load
 Balancer) IP
 , IP
 , check sum IP
 가 ,
 ,
 IP
 (Fail-over)
 ,
 RRDNS 가
 IP
 IP IP ,

가
가 , 가
가 .

2.2.3 Tunneling

Tunneling IP IP
IP 가 IP
가
IP
가
 , CRC ,

3.2.4 Direct Routing (DR)

IP 가 IP
 . Direct Routing(DR)
가 IP , 가
 , 가 가 IP
MAC Ethernet
가 IP

가 가 NAT
 가 IP 가 ,
 .

3.2.5 Broadcasting & Filtering (BF)

Broadcasting & Filtering(BF) NAT,
 tunneling, DR

. , 가 IP ,
 (FF:FF:FF:FF:FF:FF)
 . ,
 . 가 IP IP 가 , ID
 가 .
 ID
 가 ,
 가 가 ,
 가 가 ,
 .

3.

PC

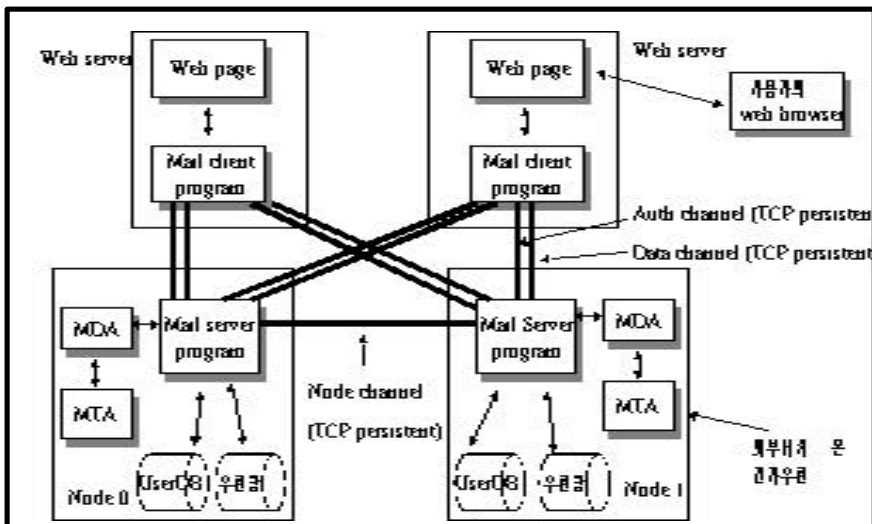
가

(mailbox)

3.1

2

2



2

MTA

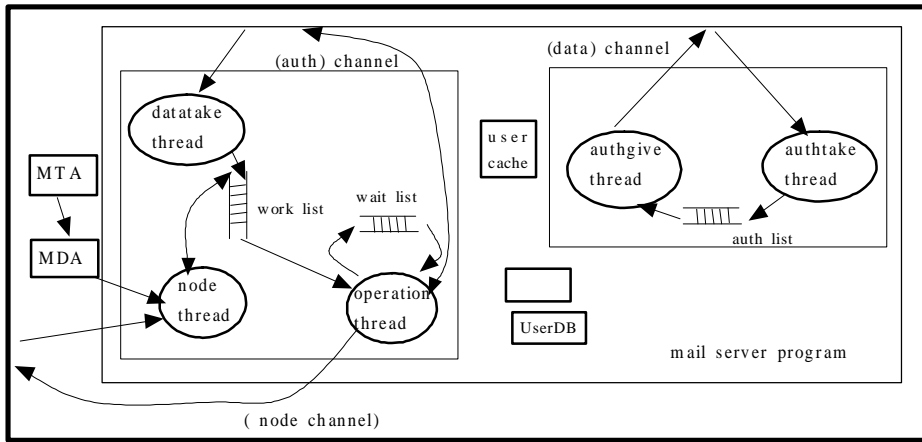
3.1.1 MTA (Mail Transfer Agent)

MTA 'sendmail' MTA
'sendmail' . MTA
sendmail
sendmail 가 MTA
[2].
'sendmail' forwarding
'/etc/passwd'
forwarding 'sendmail.cf'
'Mlocal' 'F=w' flag 가
[2].

3.1.2 MDA (Mail Delivery Agent)

가 'sendmail'
fork local MDA . linux
local MDA 'procmail' -
'procmail'
'/etc/passwd'
30,000 65,000 가
UID 가 system 가
[2].
'procmail' [2]
'sendmail' [18] 'mail.local'
가

'mail server program' TCP 'sendmail'



3

3.1.3 mail server program

'mail server program'

POSIX pthread

fork

overhead가

'authtake thread',

'authgive

thread'

'datatake thread',

'operation

thread',

MDA

'node

thread'

3.1.3.1 node thread

'mail server program'가 local MDA

'node thread'가

'node thread'

'node channel'

TCP persistent , local
 MDA . 'node
 thread'
 'operation thread'가
 queue . queue 'work list'
 'operation thread'가
 . 3 2 'mail server
 program' . 'work list'
 - - queue
 .

3.1.3.2 operation thread

'work list' 'operation thread'
 가 .
 '/etc/passwd' . Berkeley-style
 NEWDB hashed passwd file
 hash rebuild [2], SQL DB
 [2] DB
 . 3 user cache
 ID, ,
 primary node number, secondary node number

Porcupine[3][4]

'operation thread'가 가
 primary node number . primary
 node number -
 secondary node number

primary node 가

SQL DB

primary node가

node channel

'wait list' queue

mbox

UNIX linux

ID'

가

MUA POP

directory

가 32,000

가 Redhat 6.2 linux

directory sub-directory 32,000

ID.sent'

가 mbox

header 'ID.info'

()

가 Porcupine[3][4] 가

network overhead가

3.2 -

가 -
-
가
ID - script
'mail client program' TCP

3.2.1 -

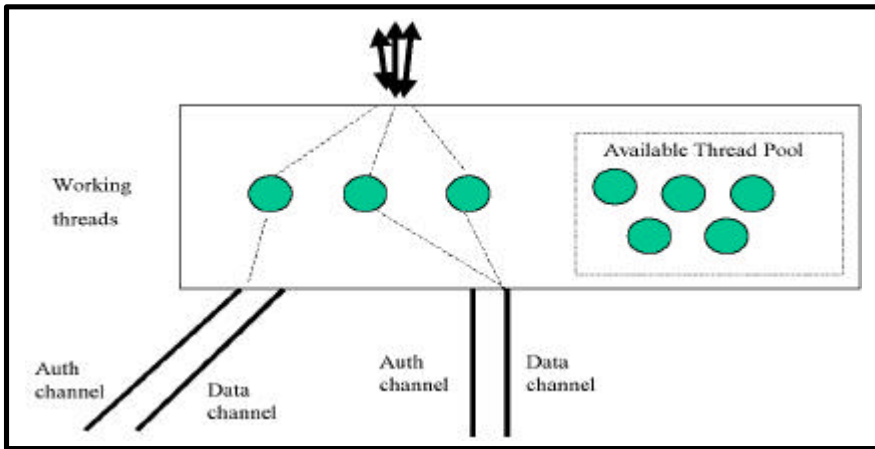
- - -
'mail client
program'
- script
'mail client program'

3.2.2 mail client program

'mail client program'
(thread pool)
가
4 2 'mail client
program'

channel) (auth channel) (data channel)
 TCP TCP overhead
 - 가 script -
 'mail client program' TCP

primary node data channel



4 mail client program

3.2.3 authtake thread, authgive thread

'mail client program' auth channel
 'authtake thread' 'auth list' queue
 'authgive thread'
 'auth list' user cache
 ID 'mail
 client program' auth channel
 primary node
 'mail client program'
 primary node

3.2.4 datatake thread

.
primary node data channel 'mail client program'
가 'operation thread'가 'datatake thread'
queue 'work list' .

3.2.5 operation thread

'work list' 'operation thread'
queue
 . mbox -
script가 .
가 script .

4.

가

가

MDA(Mail Delivery Agent)

(mailbox)

MTA

4.1

round robin DNS [14]

가

가

가

round robin DNS

DNS record

DNS [13]

100Mbps Fast Ethernet switch

switch

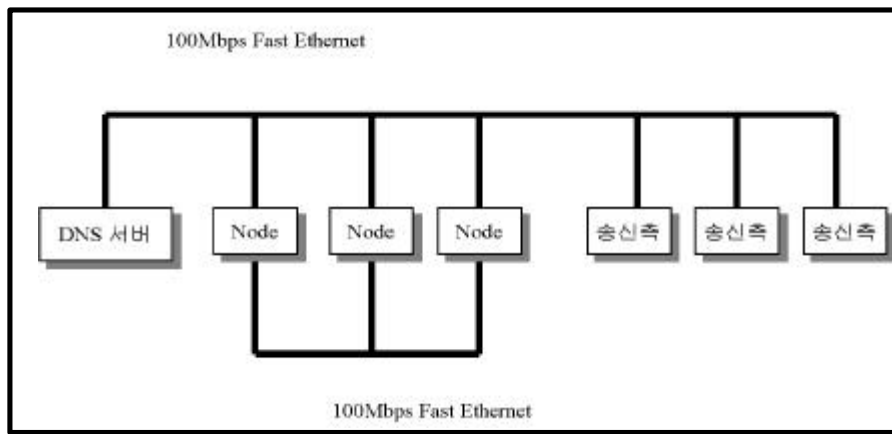
switch

card 2

5

PentiumIII 550 MHz dual CPU, 256MB
RAM, UDMA66 EIDE 5GB HDD Redhat

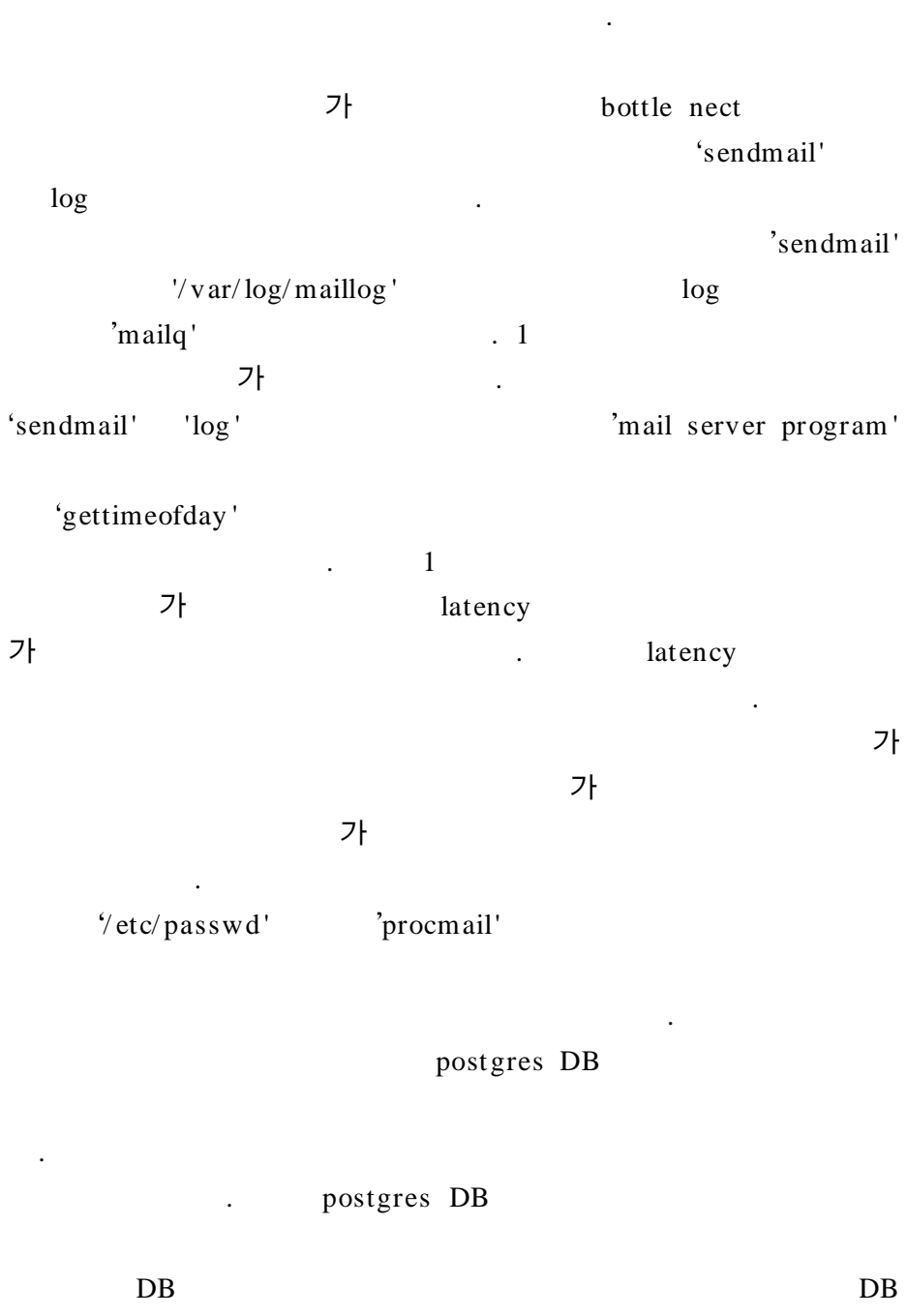
Linux 6.2 . sendmail sendmail
 8.11.1 sendmail 'Mlocal'
 . local MDA 'sendmail' 'mail.local'
 forwarding 'F=w' flag
 round robin DNS 'Cw'
 . default .



5

30,000 .
 '/etc/passwd'
 'etc/passwd' 가 35,000
 60,000 가 [2]
 512 Byte, 4 KByte, 32 KByte, 256 KByte
 가 , , ,

4.2



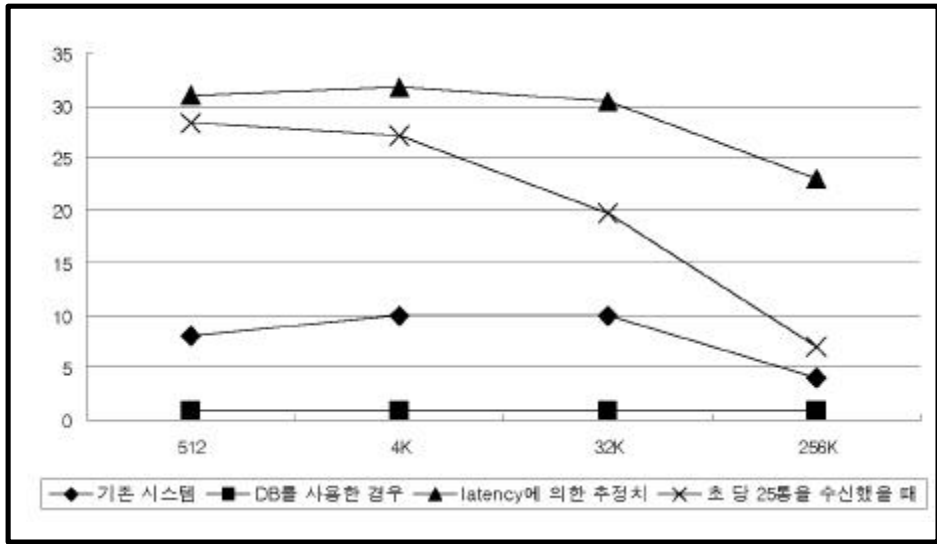
4.3

4.3.1

'/etc/passwd' 'procmail' 가
MDA가 .
dummy
MDA bottle neck 가
가 MDA 가
queue process 가 load
average가 가 . 가
msgpri 가 'msgpri > q / (la - x + 1)' 4)
가 . message 가 가 sendmail
deliver queueing 가 .
MDA
SMTP [1].
512 Byte
27 28 SMTP
35 40 SMTP .
SMTP 가
delay SMTP
가

4.3.2

4) '2.1.2.3 sendmail'



6

	512 B	4 KB	32 KB	256 KB
	8	10	10	4
DB	0.929	0.931	0.928	0.911
latency	30.988	31.763	30.403	22.971
25	28.261	27.167	19.695	7.033

1

5.

'/etc/passwd' 'procmail'
 가 , 가 MDA procmail
 SMTP connection
 .
 MDA
 - . 가
 - ,
 가
 가
 MDA
 가
 가
 PC 가
 SMTP connection 가 가
 -
 round robin DNS virtual host
 load balancing
 ,
 가
 가
 , 가 가
 .

- [1] Michael Grubb - Duke University, How to Get There From Here: Scaling the Enterprise-Wide Mail Infrastructure, Proceedings of the Tenth USENIX Systems Administration Conference (LISA '96), Chicago, IL, 1996, pp. 131-138.
- [2] Nick Christenson, Tim Bosserman, David Beckemeyer, EarthLink Network, Inc. A Highly Scalable Electronic Mail Service Using Open Systems, Proceedings of the USENIX Symposium on Internet Technologies and Systems, Monterey, California, December 1997
- [3] Yasushi Saito, Brian Bershad, Henry Levy, and Eric Hoffman, The Porcupine Scalable Mail Server, Department of Computer Science and Engineering, University of Washington
- [4] Yasushi Saito, Brian N. Bershad, and Henry M. Levy, Manageability, availability and performance in Porcupine: a highly scalable, cluster-based mail service, 17th ACM symposium on Operating System Principles (SOSP'99) Published as Operating Systems Review 34(5):1-15, Dec 1999
- [5] Eric Allman, SENDMAIL - An Internetwork Mail Router, BSD UNIX Documentation Set, University of California, Berkeley, CA, 1986.
- [6] David Wood, Programming Internet Email, Sebastopol, CA: O'REILLY & Associates, Inc., 1992
- [7] Postel, Jonathan B. RFC 821: Simple Mail Transfer Protocol, 1982.
- [8] Myers, J. and M.Rose. RFC 1725: Post Office Protocol - Version 3, 1994.
- [9] Terry Gray, Message Access Paradigms and Protocols, <http://www.imap.org/imap.vs.pop.html>, 1995.08.28
- [10] Terry Gray, Comparing Two Approaches to Remote Mailbox Access: IMAP vs. POP, <http://www.imap.org/imap.vs.pop.brief.html>, 1993.11.05
- [11] Crispin, M. RFC 1730: Internet Message Access Protocol - Version 4, 1994.
- [12] Sendmail, Home page, <http://www.sendmail.org>
- [13] Sendmail, Home page, <http://www.sendmail.net>

- [14] Bryan Costales with Eric Allman, sendmail ; Second Edition, O'REILLY, 1997.
- [15] Cuccia, Nichlos H. "The Design and Implementation of a Multihub Electronic Mail Environment". San Die , CA: USENIX Proceedings -- Lisa V; October 3, 1991.
- [16] Albiz, Paul, and Cricket Liu. DNS and BIND. Sebastopol, CA: O'REILLY & Associates, Inc., 1992
- [17] Thomas T. Kwan and Robert E. McGrath, NCSA ' s World Wide Web Server: Design and Performance, IEEE Computer, pp. 68-74, Nov. 1995.
- [18] D. Dias, W. Kish, R. Mukherjee, and R. Tewari, A Scalable and Highly Available Server. COMPCON pp. 85-92, 1996.
- [19] Wensong Zhang, Shiyao Jin, and Quanyuan Wu, Creating Linux Virtual Servers, In the proceedings of LinuxExpo, 1999.

Abstract

This paper addresses a web-mail server system using a cluster of Linux-based PCs.

The experiment shows that the low performance in general email server systems using '/etc/passwd' file and 'procmail' results from a slow MDA program and the limitation of concurrent SMTP connections.

In this paper, we modify the delivering procedure from MDA to user's mailbox in order to overcome such a problem. It is proved that the DB for a web-mail service and a wrong mailbox management approach can cause the low performance.

As a result, the implemented prototype system based on this observation, has about twice receiving performance as much as general mail server systems.

Because of an implemented web-mail system using linux PC cluster system, it has the advantage as if concurrent SMTP connections increases.

Keywords: email server system, cluster, mailbox, MDA, Linux PC

Student Number: 99419 - 542

