

2003 3

3 1994 12 1992 1989 3 1980

UNIX SVR4 Solaris Linux Windows 2000/XP

CH1

; Windows 2000/XP / ;

CH2

; Window2000/XP Solaris

Linux UNIX SVR4 Windows2000/XP Linux
CH3

]

x86/Pentium : Windows 2000/XP Linux Intel

CH5 I/O I/O I/O I/O ;
RAID /
I/O Windows2000/XP I/O Linux

CH6 Windows2000/XP

Linux
CH7 ;

CH8 Windows 2000/XP

Windows2000/XP 1999

98 99 00 ppt

ppt

ppt !

feixl@nju.edu.cn luobin@nju.edu.cn email

2003 3

CH1	1
1.1	1
1.1.1	1
1.1.2	2
1.1.3	4
1.2	6
1.2.1	6
1.2.2	6
1.2.3	8
1.2.4	11
1.3	18
1.3.1	18
1.3.2	19
1.3.3	20
1.3.4	22
1.4	26
1.4.1	26
1.4.2	29
1.4.3	29
1.4.4	31
1.4.5	/	32
1.4.6	35
1.4.7	<i>Windows 2000/XP</i> /	37
1.5	41
1.5.1	<i>DOS</i>	41
1.5.2	<i>Windows</i>	42
1.5.3	<i>UNIX</i>	44
1.5.4	<i>Linux</i>	45
1.5.5	<i>IBM</i>	47
1.5.6	49
1.6	50
CH2	56
2.1	56
2.1.1	56
2.1.2	57
2.1.3	58
2.1.4	58
2.1.5	59

2.2	60
2.2.1	60
2.2.2	60
2.2.3	61
2.2.4	63
2.2.5	63
2.2.6	68
2.2.7	<i>Windows 2000/XP</i>	69
2.2.8	<i>Solaris</i>	75
2.2.9	<i>Linux</i>	76
2.3	79
2.3.1	79
2.3.2	80
2.3.3	83
2.3.4	87
2.3.5	88
2.3.6	<i>UNIX SVR4</i>	91
2.3.7	<i>Linux</i>	94
2.4	97
2.4.1	97
2.4.2	98
2.4.3	103
2.4.4	<i>Solaris</i>	106
2.4.5	<i>Windows 2000/XP</i>	110
2.5	116
2.5.1	117
2.5.2	117
2.5.3	118
2.5.4	118
2.5.5	119
2.6	119
2.6.1	119
2.6.2	120
2.6.3	121
2.6.4	121
2.7	124
2.7.1	124
2.7.2	124
2.7.3	127
2.7.4	128
2.7.5	<i>UNIX SVR4</i>	132
2.7.6	<i>Windows 2000/XP</i>	133

2.7.7	<i>Linux</i>	140
2.8		143
CH3		151
3.1		151
3.1.1		151
3.1.2		151
3.1.3		153
3.1.4	<i>Interaction Among Processes</i>	155
3.2		156
3.2.1		156
3.2.2		157
3.2.3		158
3.2.4		160
3.3	PV	162
3.3.1		162
3.3.2	<i>PV</i>	163
3.3.3		166
3.3.4	-	168
3.3.5	-	170
3.3.6		172
3.4		172
3.4.1		172
3.4.2	<i>Hoare</i>	175
3.4.3	<i>Hanson</i>	178
3.5		184
3.5.1		184
3.5.2		185
3.5.3		188
3.5.4		189
3.5.5		192
3.6		195
3.6.1		195
3.6.2		196
3.6.3		197
3.6.4		198
3.6.5		206
3.7	WINDOWS 2000/XP	209
3.7.1	<i>Windows 2000/XP</i>	209
3.7.2	<i>Windows2000/XP</i>	210
3.8	LINUX	211
3.9		212

CH4	225
4.1	225
4.1.1	225
4.1.2	<i> caching</i>	226
4.1.3	227
4.2	227
4.2.1	227
4.2.2	229
4.2.3	230
4.3	235
4.3.1	235
4.3.2	236
4.3.3	237
4.3.4	238
4.3.5	238
4.3.6	240
4.4	241
4.4.1	241
4.4.2	241
4.4.3	243
4.4.4	243
4.5	243
4.5.1	243
4.5.2	245
4.5.3	261
4.5.4	262
4.6	INTEL X86/PENTIUM	263
4.6.1	<i> Intel x86/Pentium</i> — —	264
4.6.2	<i> Intel x86/Pentium</i>	265
4.6.3	<i> Intel x86/Pentium</i>	265
4.6.4	<i> Intel x86/Pentium</i>	266
4.7	WINDOWS 2000/XP	268
4.7.1	268
4.7.2	269
4.7.3	273
4.8	LINUX	279
4.8.1	<i> Linux</i>	279
4.8.2	279
4.8.3	280
4.8.4	281
4.8.5	282

4.8.6	283
4.8.7	284
4.9	285
CH5	292
5.1 I/O	292
5.1.1 I/O	292
5.1.2 I/O	293
5.1.3	297
5.2 I/O	298
5.2.1 I/O	298
5.2.2 I/O	299
5.2.3	300
5.2.4	I/O	300
5.2.5	I/O	302
5.3	I/O	302
5.3.1	302
5.3.2 I/O	I/O	304
5.3.3	I/O	305
5.4	306
5.4.1	306
5.4.2	307
5.4.3	307
5.5	308
5.5.1	308
5.5.2	309
5.5.3	310
5.5.4	311
5.5.5	311
5.5.6	313
5.5.7	I/O	316
5.6	316
5.6.1	316
5.6.2	317
5.7	318
5.7.1	318
5.7.2 SPOOLING	319
5.7.3 SPOOLING	321
5.8	WINDOWS 2000/XP I/O	321
5.8.1 Windows 2000/XP I/O	321
5.8.2 Windows 2000/XP I/O	325
5.8.3 Windows2000/XP	328

5.8.4 Windows 2000/XP I/O	331
5.8.5 Windows 2000 XP	333
5.9 LINUX	344
5.9.1 Linux	344
5.9.2 Linux	345
5.9.3 Linux	346
5.9.4 Linux	347
5.10	347
CH6	352
6.1	352
6.1.1	352
6.1.2	353
6.1.3	353
6.1.4	354
6.1.5	355
6.1.6	356
6.2	357
6.2.1	357
6.2.2	358
6.2.3	358
6.2.4	359
6.3	361
6.3.1	361
6.3.2 685.7460#56534a71d58fTT8 1 T#0.5 0 0 10.5 205.92 e...1422c067TD0 Tc(TjTT8 1 T#38.4457 -1.5543 TD-0.	

6.6	WINDOWS 2000/XP	397
6.6.1	<i>Windows 2000/XP</i>	397
6.6.2	<i>Windows2000/XP</i> FSD	398
6.6.3	NTFS	401
6.6.4	NTFS	402
6.6.5	NTFS	405
6.6.6	NTFS	405
6.7		406
CH7		411
7.1		411
7.2		411
7.3		415
7.3.1		415
7.3.2		415
7.3.3		415
7.3.4		416
7.4		416
7.4.1		416
7.4.2		417
7.4.3		420
7.5		421
7.5.1		421
7.5.2		423
7.5.3		423
7.5.4		423
7.5.5		425
7.6		426
7.6.1		426
7.6.2		438
7.6.3		442
7.6.4		452
7.6.5		456
7.7	WINDOWS 2000/XP	457
7.7.1	<i>Windows 2000/XP</i>	457
7.7.2	<i>Windows2000/XP</i>	457
7.7.3	<i>Windows2000/XP</i>	458
7.7.4		458
7.7.5		459
7.7.6		459
7.8		461
CH8		465

8.1	465
8.1.1	465
8.1.2	467
8.1.3	468
8.2	474
8.2.1	474
8.2.2	475
8.2.3	475
8.3	477
8.3.1	477
8.3.2	478
8.3.3	484
8.3.4	485
8.3.5	495
8.3.6	497
8.3.7	502
8.4	WINDOWS2000	504
8.4.1	<i>Windows 2000</i>	504
8.4.2	<i>WindowS 2000</i>	518
8.5	522

CH1

1.1

1.1.1

Operating System OS

- OS
- OS
- OS
- OS
- OS

1-1

/

1.1.2

Virtual Machine

I/O I/O

Machine

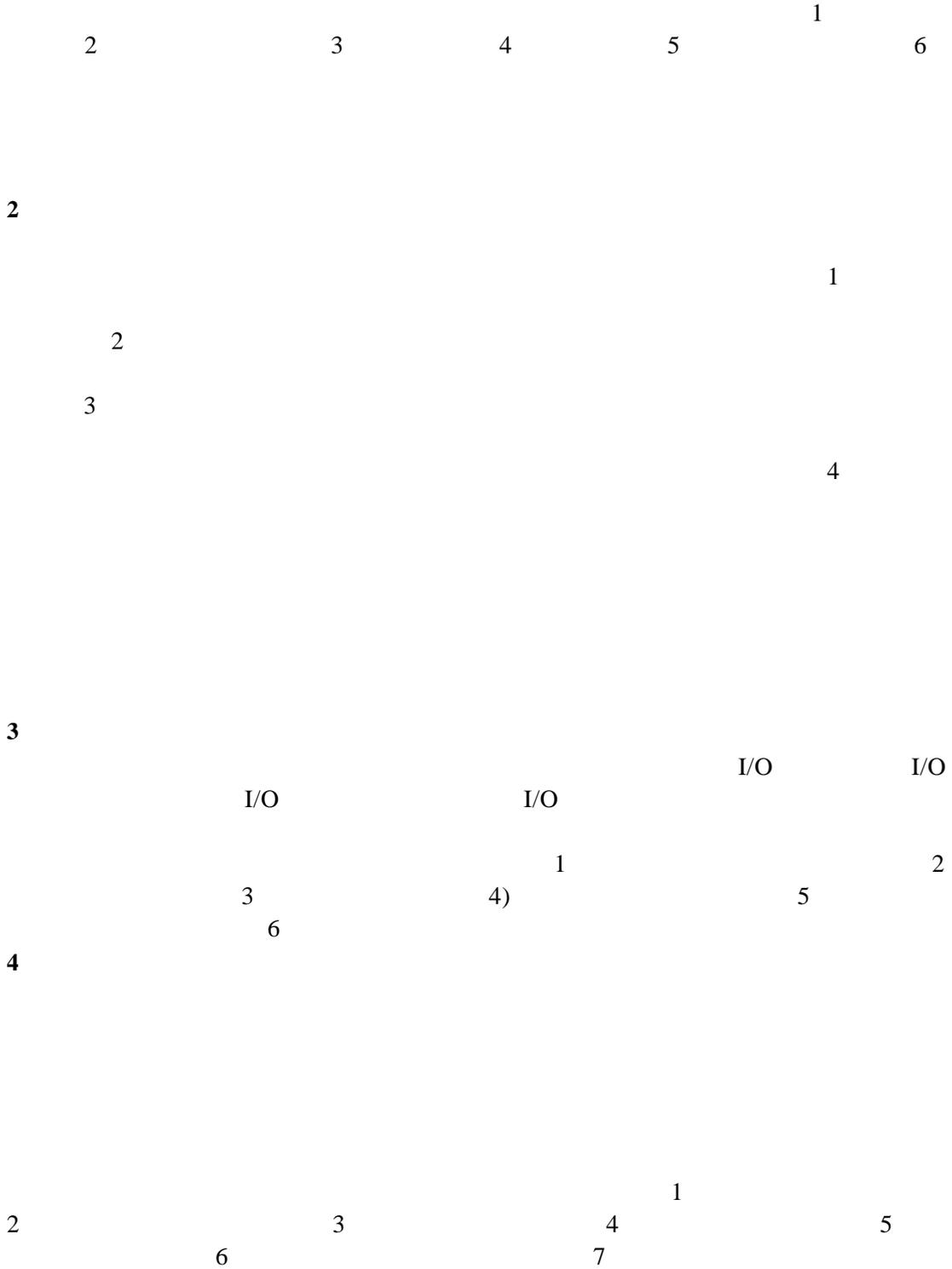
Virtual

I/O I/O

1

process

thread



5

1

2

3

6

1.1.3

1

concurrency

I/O

I/O CPU

I/O
CPU

CPU

CPU

I/O

?

multitasking system

CPU
CPU I/O I/O

I/O

CPU

parallelism

CPU

I/O

CPU

2 (sharing)

" "

3 (asynchronism)

" " CPU CPU
CPU

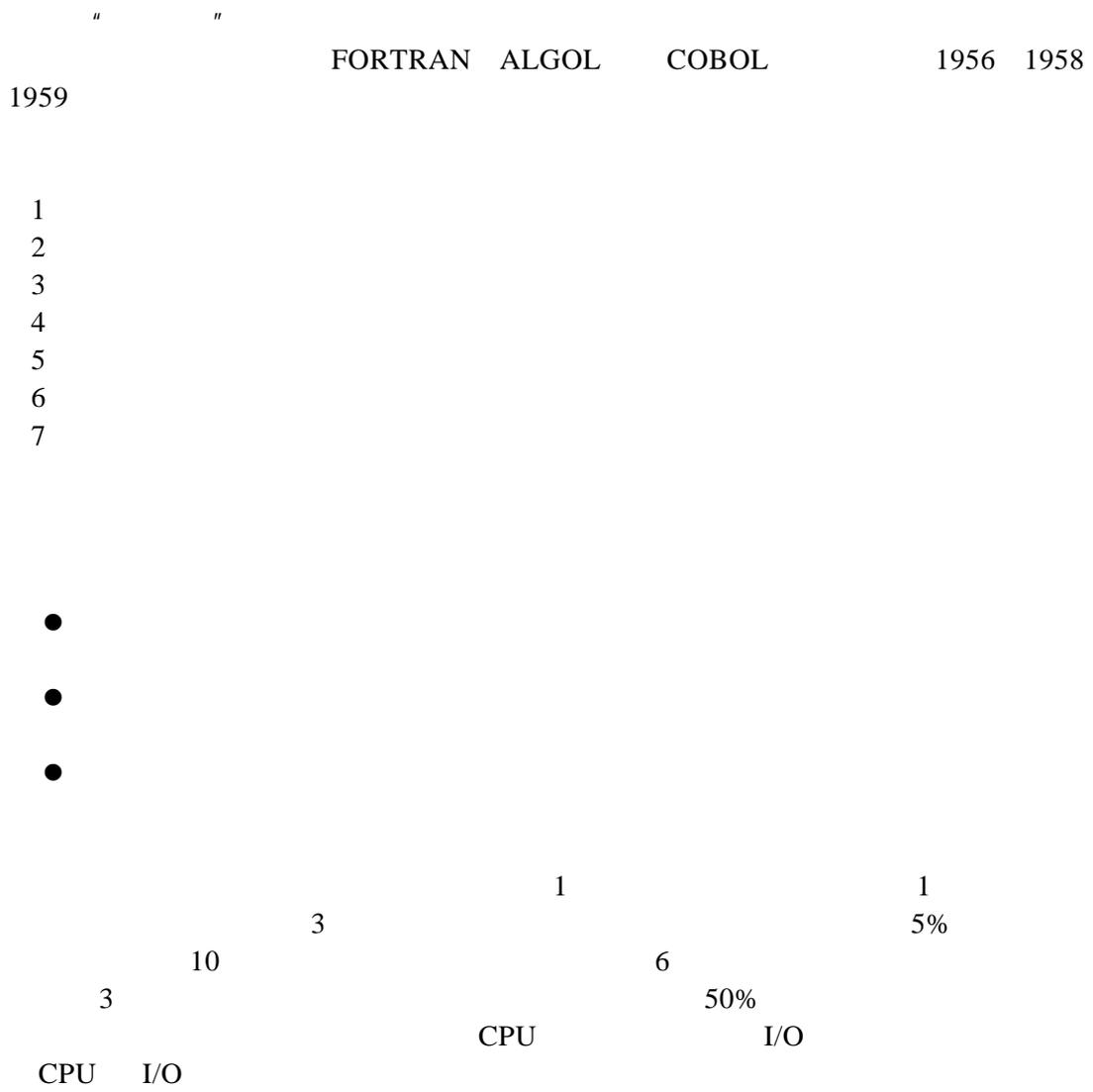
4 (virtual)

Spooling CPU CPU CPU CPU
CPU CPU CPU CPU
IBM VM

) () (

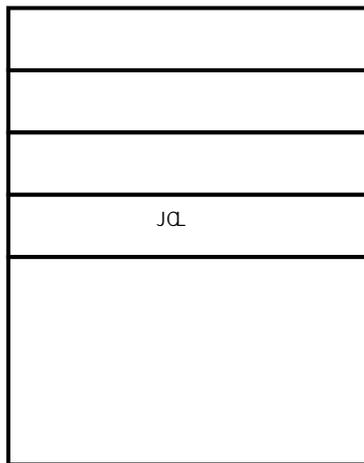
1.2

1.2.1



1.2.2

I/O CPU CPU I/O I/O I/O I/O



1-2

Monitor System IBSYS IBM 7094 Monitor System resident monitor FMS FORTRAN
1-2



Job Control Language

\$FTN FORTRAN \$JOB \$LOAD
\$DATA \$RUN \$SEND

•

•

•

•

I/O I/O

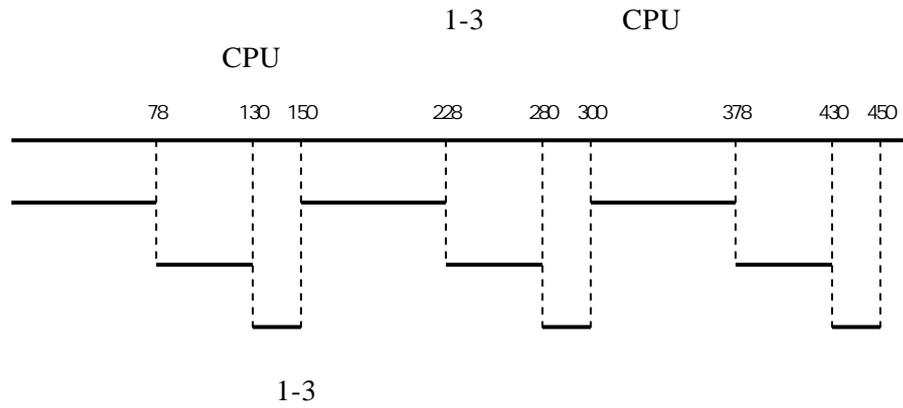
I/O

FORTRAN

I/O

1.2.3

1



20 CPU 60 I/O

multiprogramming

()

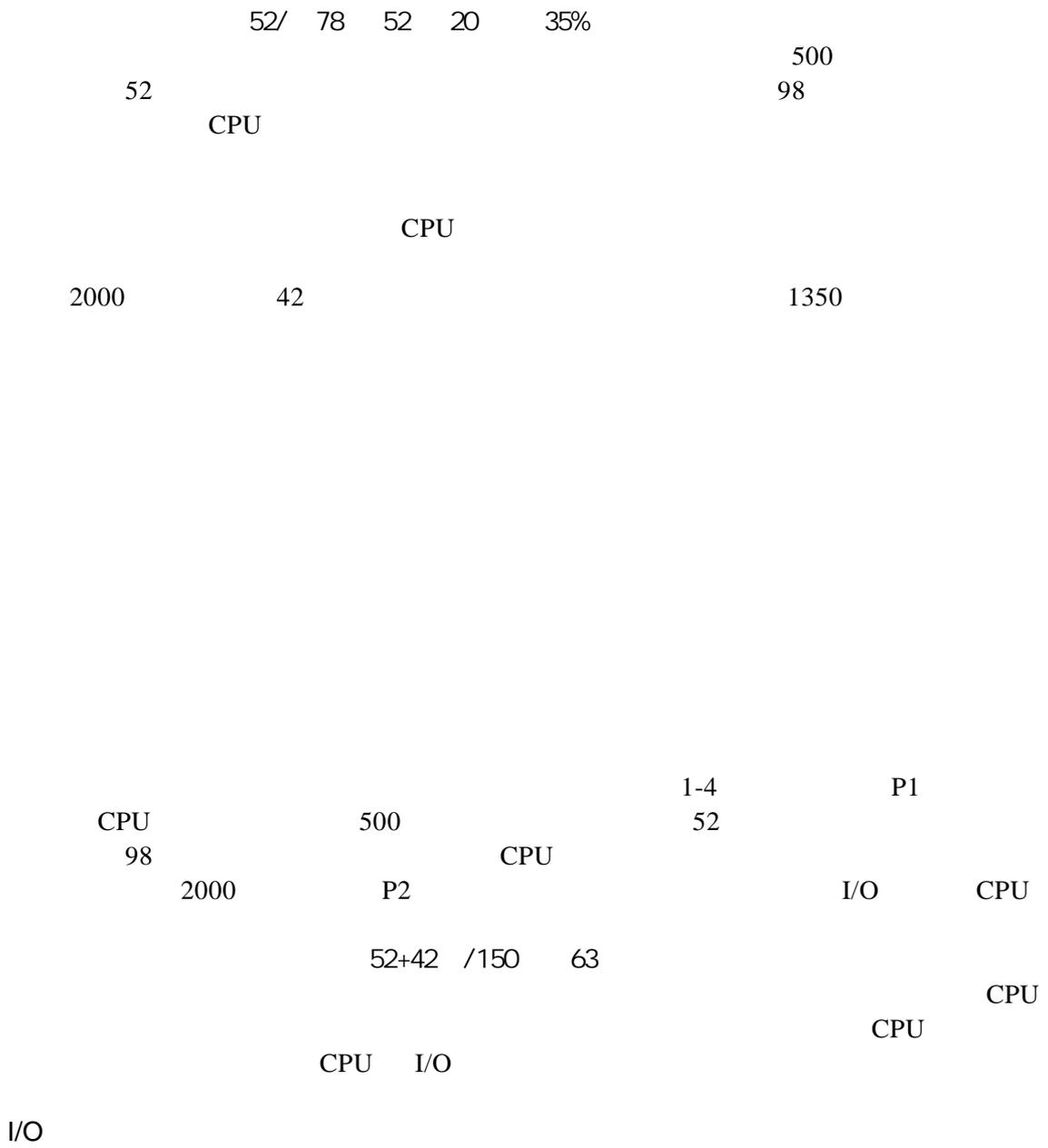
()

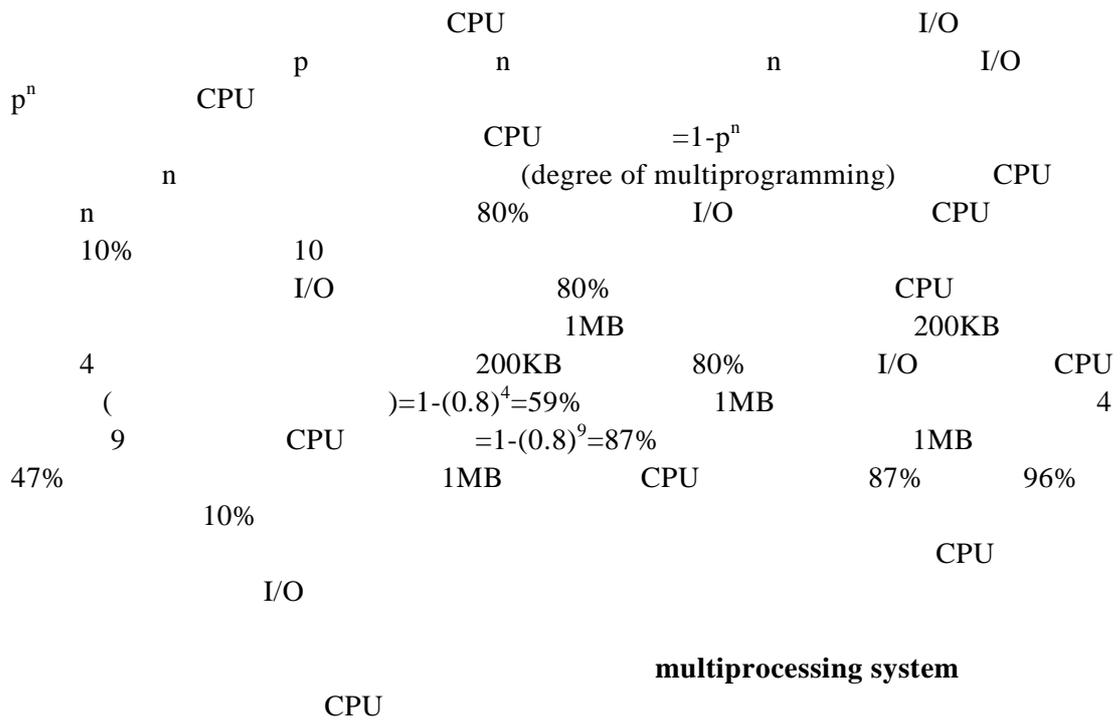
CPU

CPU

500 2000 6400 500 10 52

1-3





CPU

2

I/O

-
-
-
-
-
-

SPOOLING

1.2.4

18	1			8	16	"	"	32
		64		8				16
	32		64					
	2							

3

4

GUI

5

1

System

Batch Operating

•

•

• /

2

Time Sharing Operating System

CPU

	1959	MIT		1962		
CTSS	Compatible Time Sharing System				IBM 7094	32
	1965	8	IBM	360		TSS/360(Time
Sharing System/360)						
1965			MIT	BELL	GE	

/ /

3

Real Time Operating System

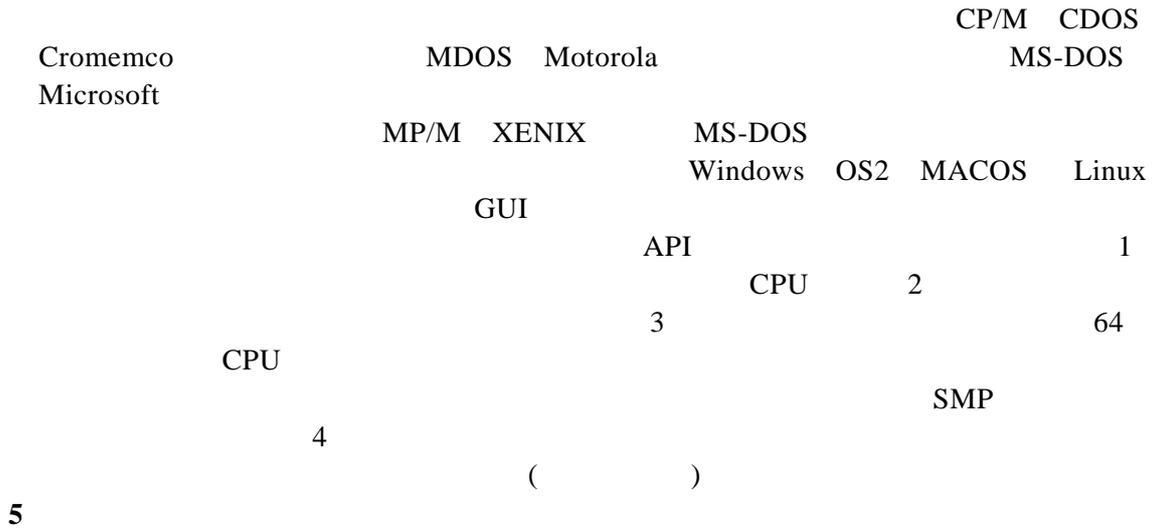
2 1

3 4

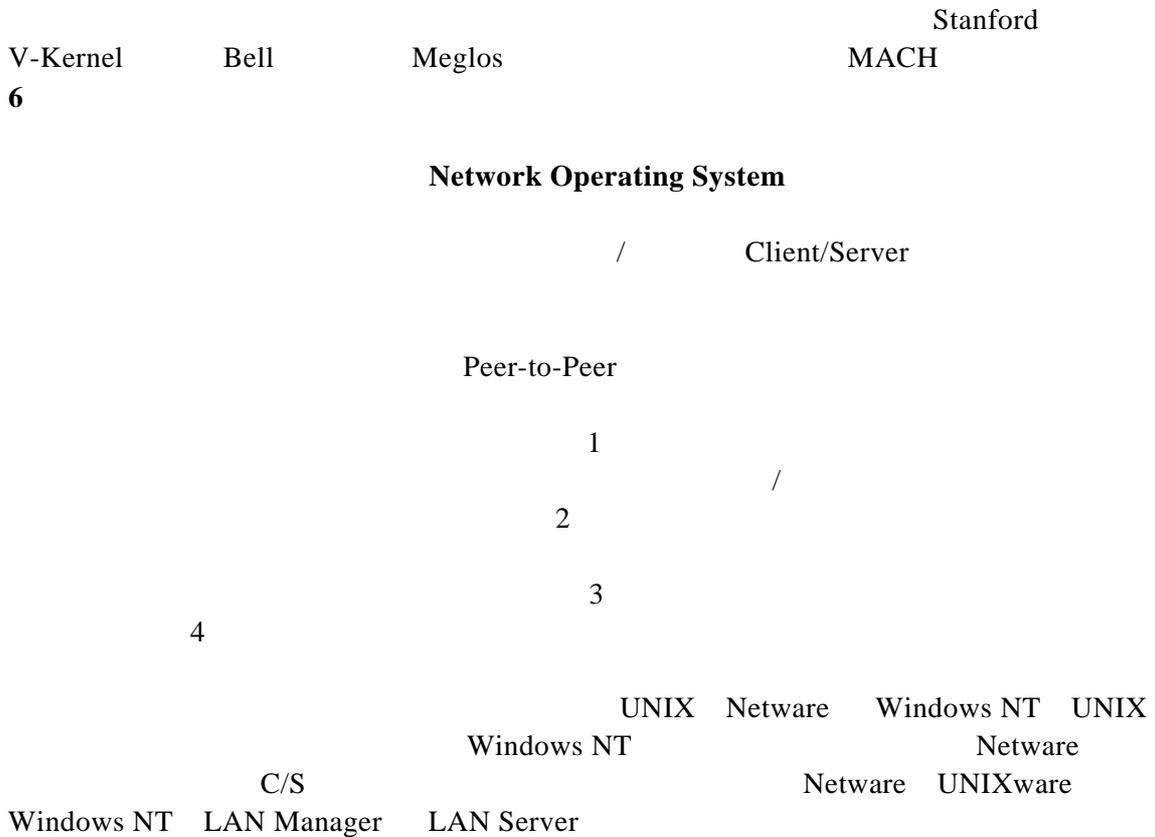
"

" " "

"



parallel processing

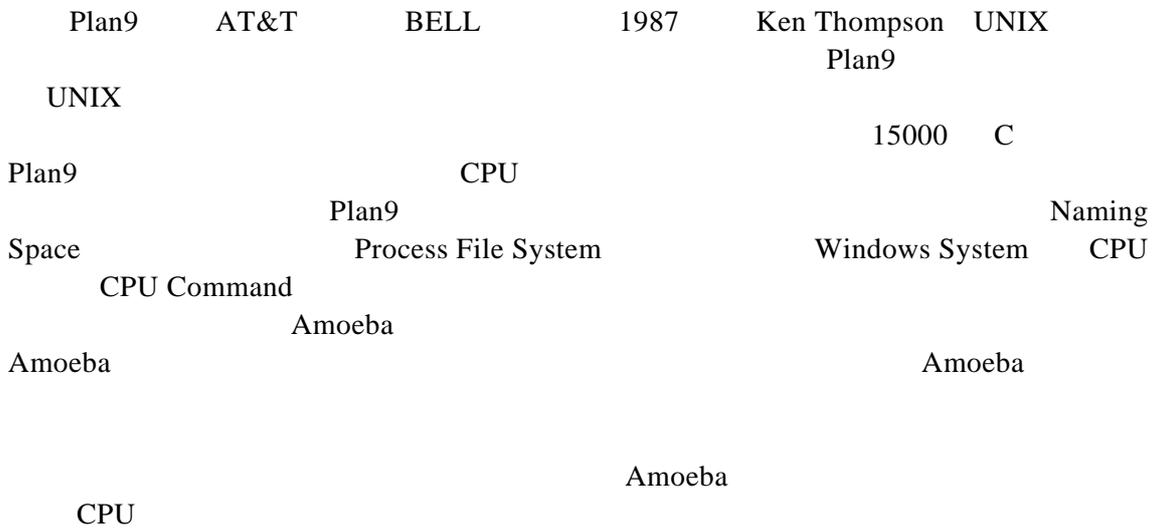


-
-
-
-

7

System

Distributed Operating



-
-
-
-
-

8

3C Computer Communication Consumer Electronics Internet

3C

()

()
()

(embedded software)

()

(8 16) (1) ()

84 ! ã

Abstraction level)	BSP(Board Support Package)	HAL(Hardware HAL)
		BSP
Windows CE	VxWorks()	Linux Symbian
Chorus	Diba Sun	chorus
Psos ISI	QNX QSSL	OS-9 Microsoft
Microsoft	HOPEN	Windows CE
Windows CE	" "	32
	Personal Java	SUN
	Java	Java
		SUN
		Java OS for Consumers
	Embedded Java	Hopen
	" "	Hopen
	10kB	C
Gb2312-80	Hopen	Win32API Personal Java

1.3

1.3.1

1

2

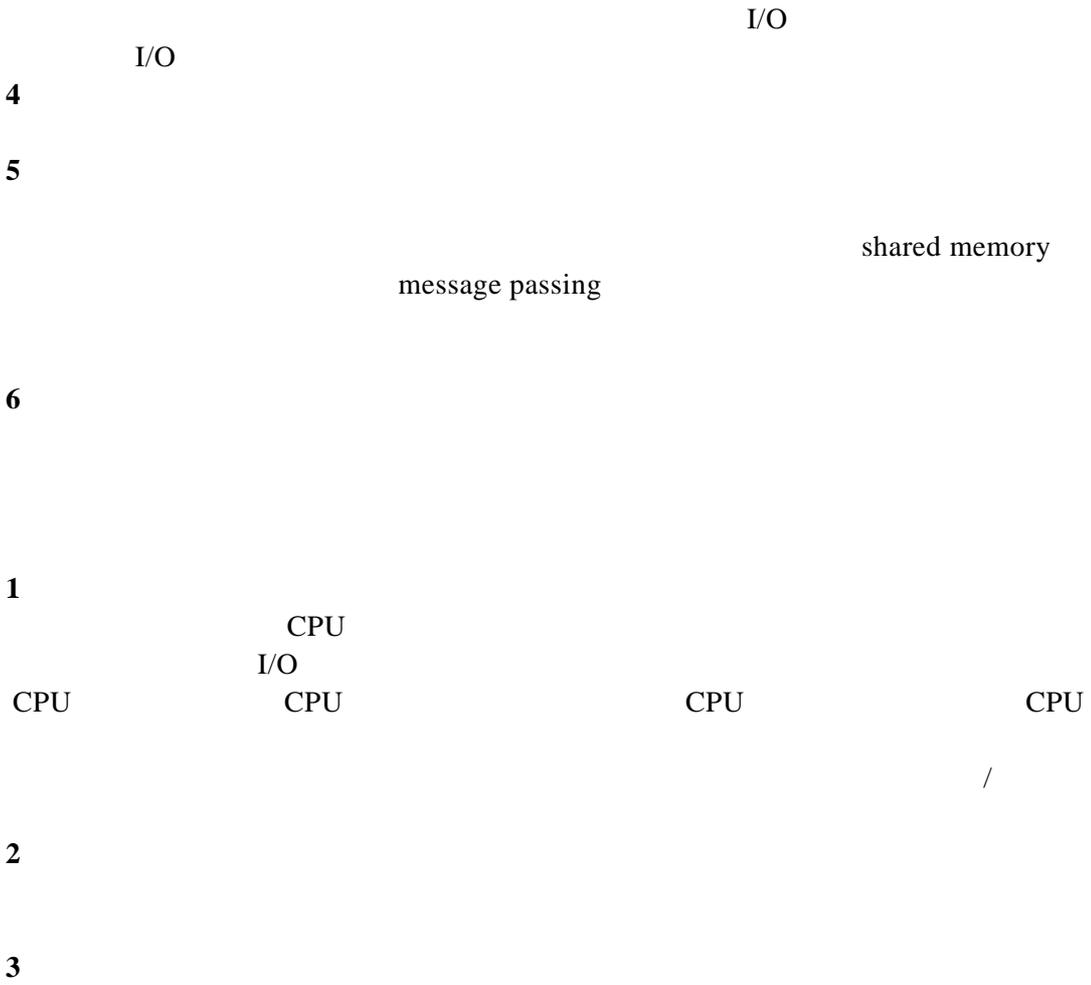
3

I/O

I/O

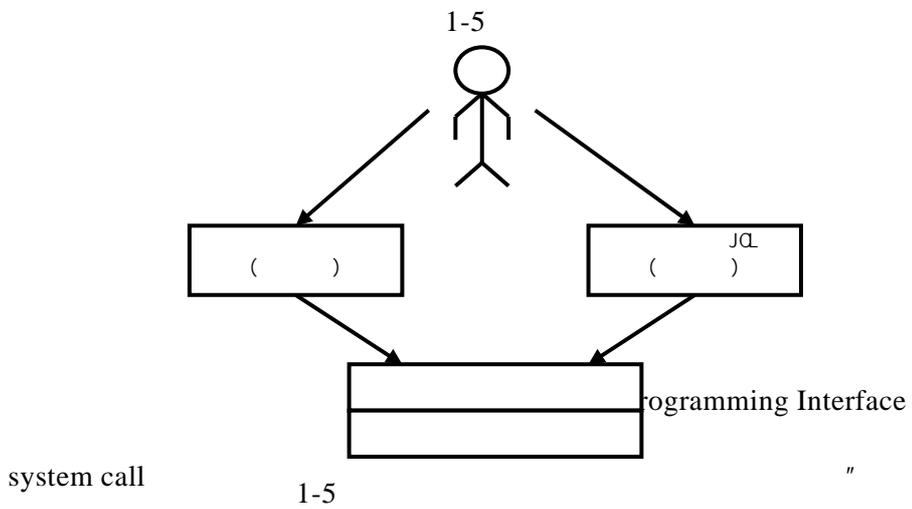
I/O

I/O



shell

1.3.2



() () ()

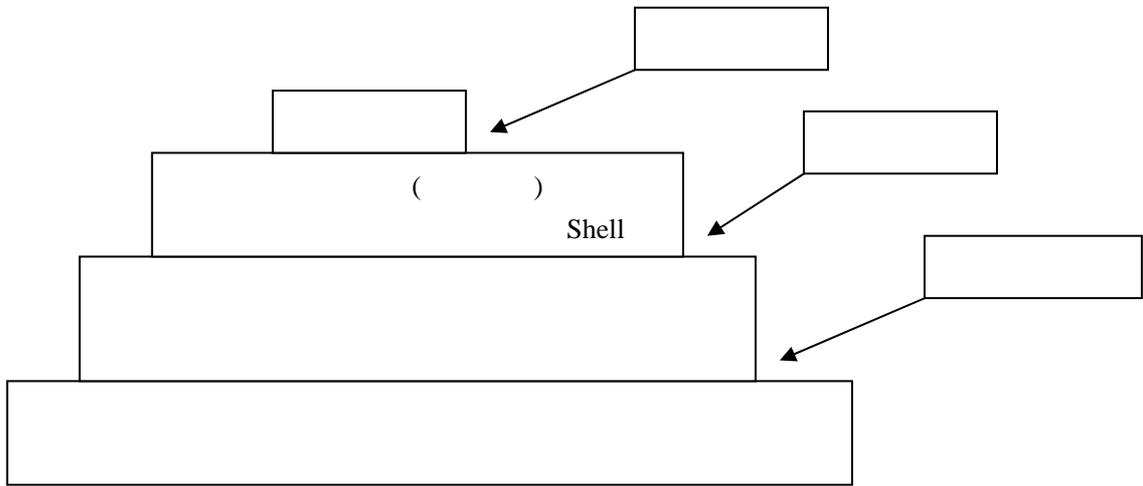
1.3.3

1

) ()
 " UNIX " ()
 POSIX1003.1 Portable Operating System IX
 UNIX

(C)
 ()

UNIX Linux Windows OS 2 C
 C 1-6
 UNIX/Linux



1-6 Unix/Linux

1 2
 3 4 I/O 5)

Windows		API	Kernel	User	GDI	Kernel
			User			
		GDI				
	Windows	API				

Windows

DLL Dynamic Link Library

2

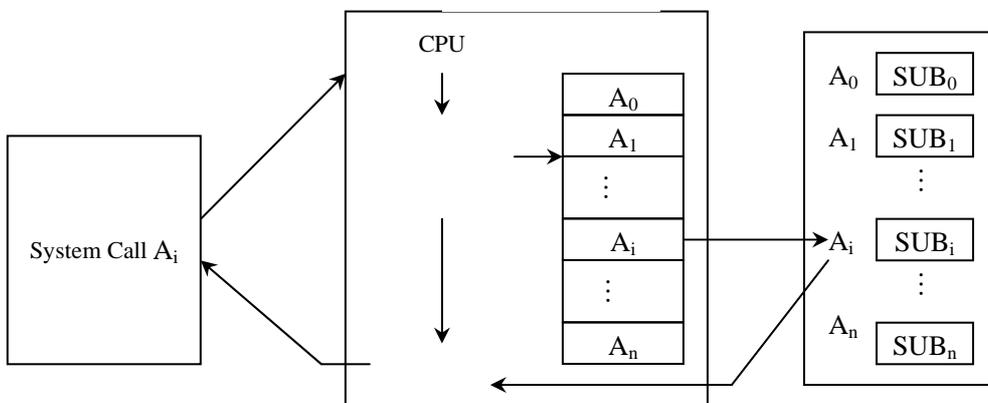
supervisor

trap

interrupt

1-7

CPU



1-7

3

(1)

(2)

(3)

(4)

4 Linux

Linux	UNIX			
Linux	Linux	190	Shell	
Linux			1	C
		2		API
Linux	lib.a			

Linux entry(sys_call_table)

include/asm/unistd.h

ENTRY(sys-call-table)

long SYMBOL_NAME(sys_ni_syscall) 0

long SYMBOL_NAME (sys_exit) 1

long SYMBOL_NAME (sys_fork) 2

long SYMBOL_NAME (sys_read) 3

long SYMBOL_NAME (sys_write) 4

long SYMBOL_NAME (sys_open) 5

long SYMBOL_NAME(sys_close) 6

... ..

long SYMBOL_NAME (sys_ni_syscall)

long SYMBOL_NAME (sys_ni_syscall)

long SYMBOL_NAME (sys-vfork) 190

Linux

Linux 0x80 int80h

1.3.4

1

()

(1) ---

1)

Command	Command	arg1	arg2	argn
Linux		1	cd	chmod chown chgrp
comm cp	crypt diff file find ln ls mkdir mv od pr pwd rm rmdir	2		
	at kill mail nice nohup ps time write mesg	3		
cat crypt grep norff uniq wc sort spell tail troff	4		cc	f77
login logout size yacc vi emacs dbs lex make lint ld	5			date
man passwd stty tty who				

MS-DOS	BAT	BAT
--------	-----	-----

UNIX	Linux	Shell	Shell	Shell	Shell
------	-------	-------	-------	-------	-------

2

		GUI	Graphics User Interface
GUI		WIMP	(Window Icon
Menu	Pointing device)		

Aito Research Center	1981	Star8010	GUI	GUI	Xerox	Palo
Apple	Apple Lisa	Macintosh				1983
Microsoft	Windows	IBM	OS/2	UNIX	Linux	GUI
GUI		GUI			GUI	X-Window
					MIT	X-Windows
			Windows NT	Visual C++	Visual Basic	
			GUI			

3)

(2)

()

JCL Job Control Language

UNIX/Linux

Shell

JCL

JCL

JCL

JCL

2

command interpreter

()

" "

" "

delete G
delete

delete

G
delete

1.4

IBM OS/360	MIT	1963		CTSS	32000	
1975	MIT	Bell		4000		5000
Windows 2000		2500	Multics	2000		Brooks
	OS/360				3200	
		20	60		OS/360	
			10	20		
				CPU		

POSIX

1.4.1

1

(microkernel)

(monolithic kernel)

Linux

(module)

Linux

Linux

1

2

3

1

2

Solaris

3

Linux

top half bottom half

Windows 2000/XP

ISR

Linux

ISR

DPC

DPC

DPC

IRQL

Dispatch/DPC IRQL

DPC

4

1)

2

3

2

1

2

3

3

CPU

multithreaded process

4

() 70E0N,86wj@,€

1À@?Z@GQ•Á&>ZM

-
-
-
-

5

1975
Solo

PDP 11/45

1.4.2

IBM S/360

1.4.3

A_0
 A_1 A_2 A_1
 --- " ---

(1)

(2) () CPU
 ()
 (3)

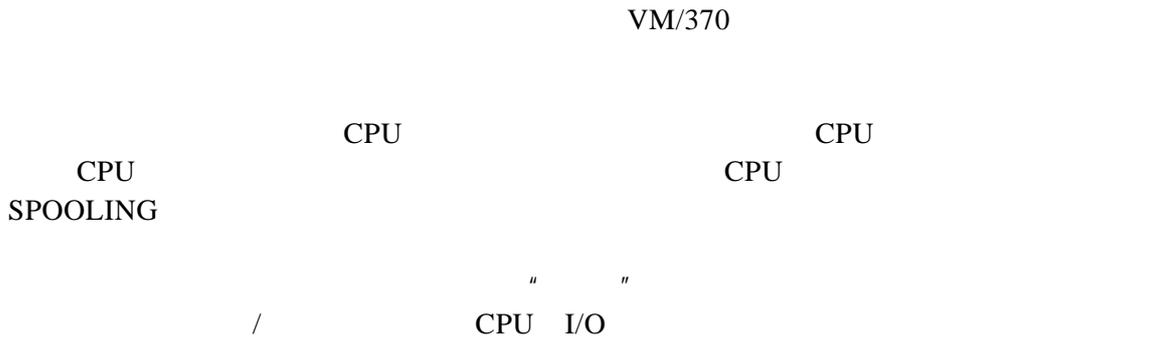
(4)

CPU
 Edsger.W.Dijkstra 1968 THE Electrologica X8
 THE
 6 0
 1
 2
 3 I/O I/O
 4
 I/O 5

1.4.4

MacKinnon 1979
 1-8(a)

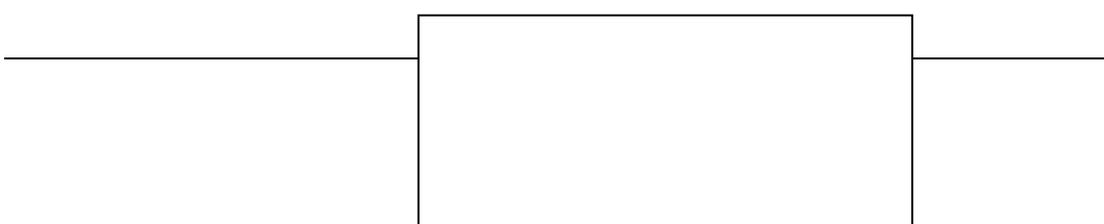
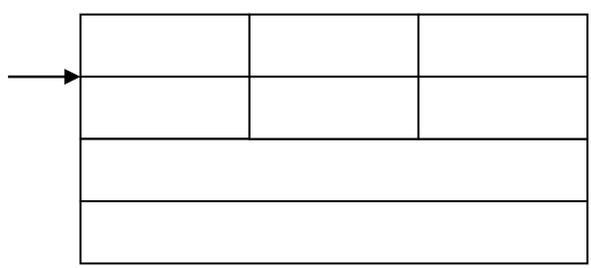
IBM CP/CMS VM/370
 IBM S/390 1 2
 VM/370



OS/360
 1-8(b)

CMS CMS VM/370 I/O VM/370 CMS

VM/370



1.4.5 /

1 /

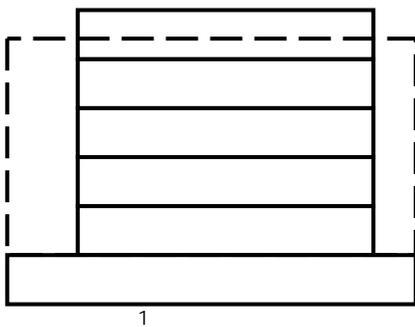
- / Client/Server Mach /

C/S / ;

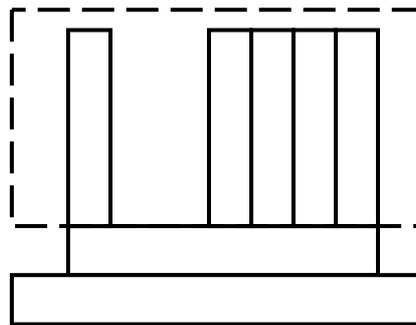
C/S) (

/ ;

(microkernel) /



1



2

1-9

1-9

I/O

(1

(2

(3

CPU

CPU

(4

API

(5

(6

Windows2000/XP

Mach

Chorus

/

300KB

140

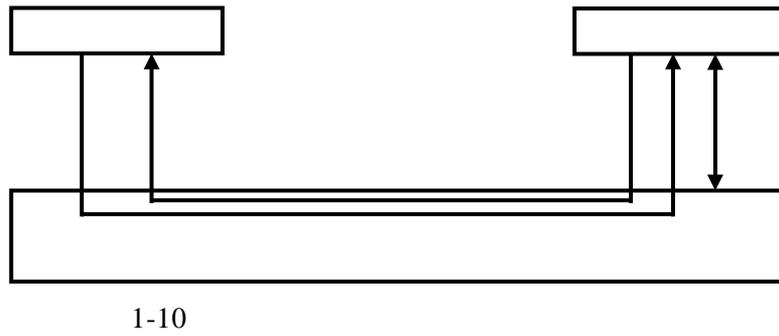
L4

12KB

7

2

(1)



1-10

- grant
- map
- flush

(2)

header

Ports

CPU

(3)I/O

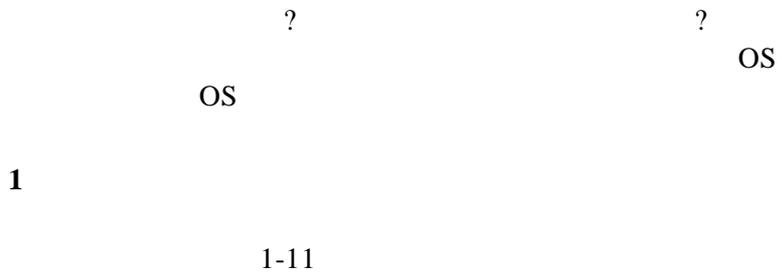
I/O

```

driver thread;
do
  wait for (mhg sender);
  if sender = my_hardware_interrupt
  {
    read/writer I/O ports;
    reset hardware interrupt
  }
  else ...
while (true);

```

1.4.6



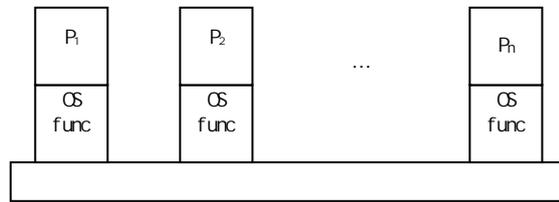
2 OS

UNIX

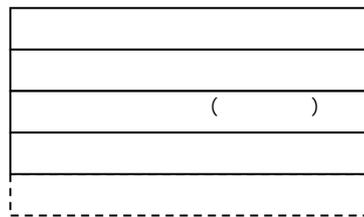
OS

1-12

1-13 OS



1-12 OS



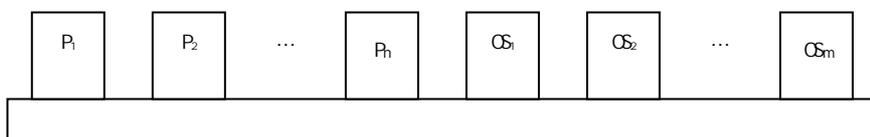
1-13 OS

3 OS
OS

Client/Server

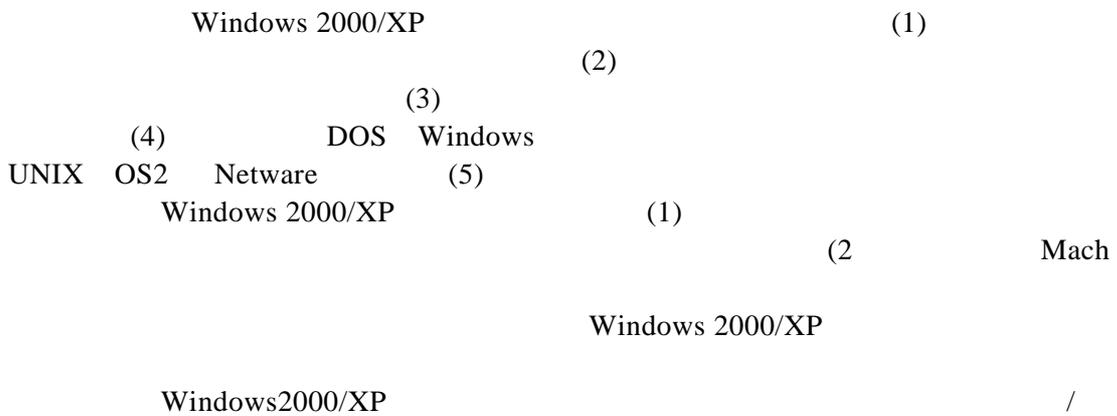
Window 2000/XP

1-14

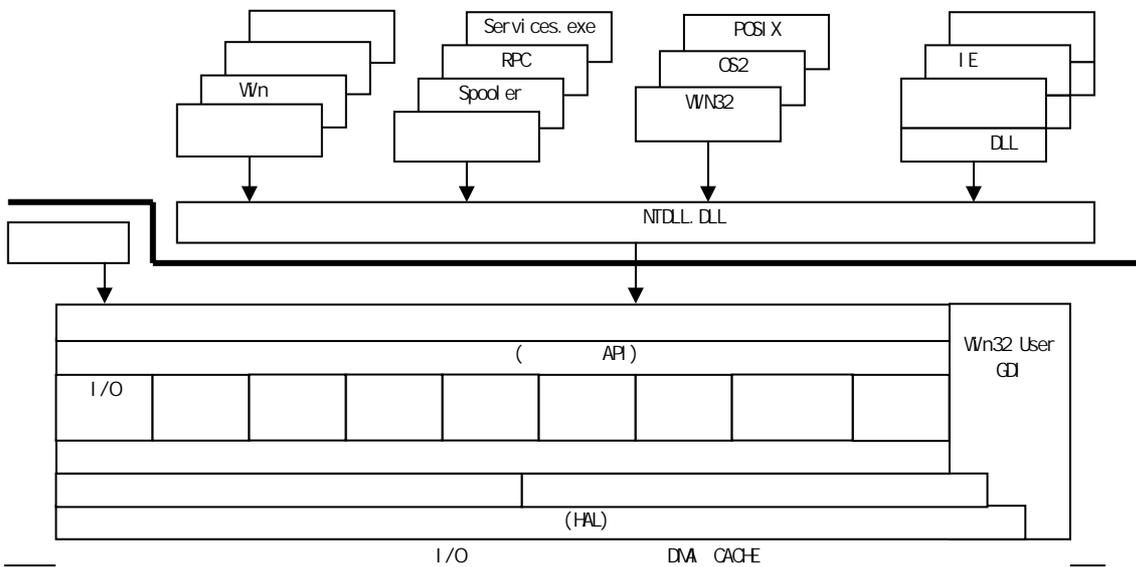


1-14 OS

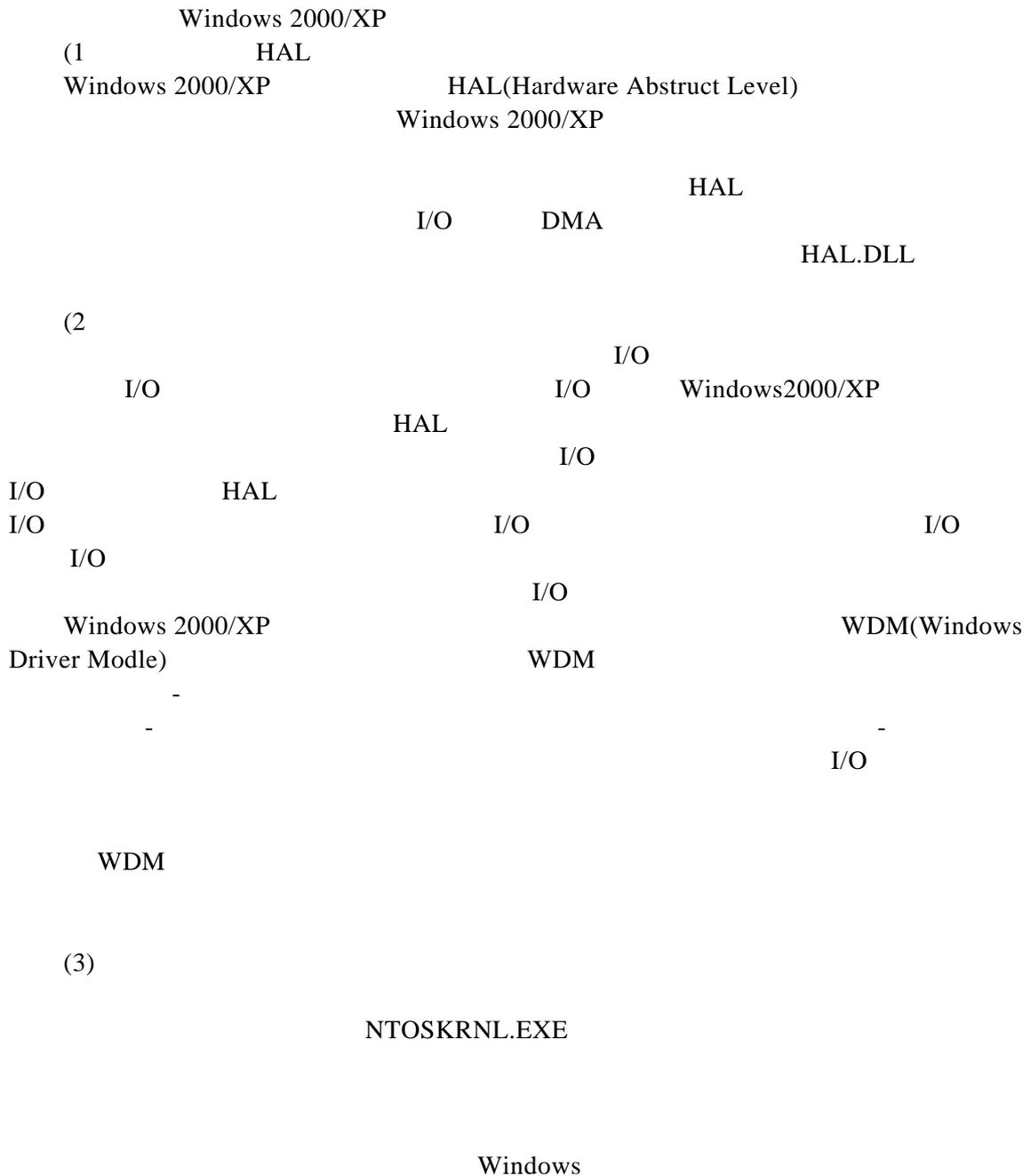
1.4.7 Windows 2000/XP /



1-15 Windows 2000/XP



I/O



" "

" "

DPC Deferred procedure Call APC Asynchronous Procedure Call
I/O

" "

Windows 2000/XP

HAL

(4)

Windows 2000/XP NTOSKRNL.EXE

NTDLL.DLL WIN32 API

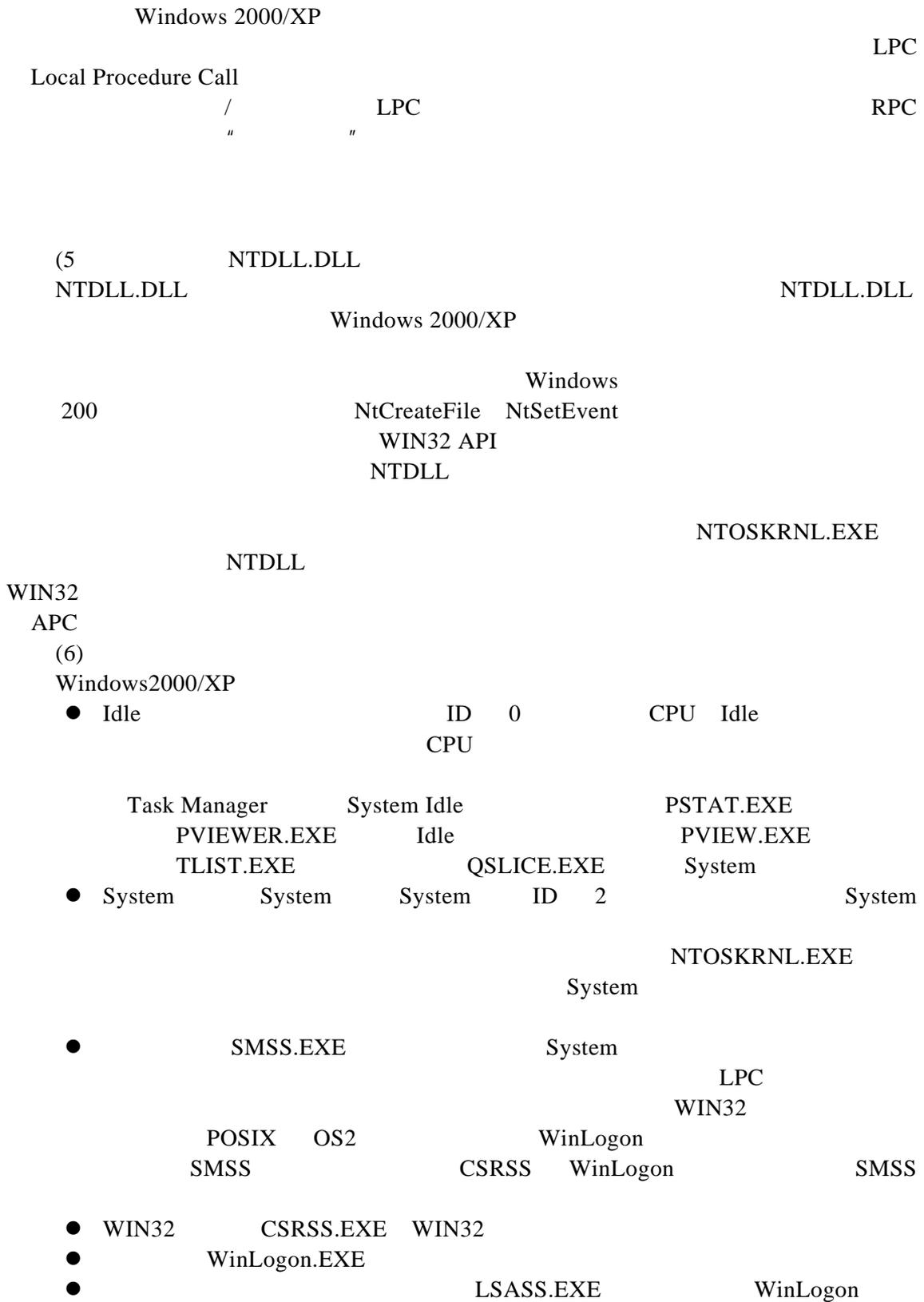
NtCreatePagingFile LPC NtQueryInformationxxx
Windows 2000/XP DDK

Windows 2000/XP

" "

I/O I/O I/O API
I/O

I/O



(7)

SERVICES.EXE

UNIX

Win32

Win32

Windows2000

RPC

(8

" "

Windows 2000/XP

WIN32 POSIX(

POSIX.1)

OS/2 1.2

DLL

Win32

Windows2000/XP

Win32

Win32

CSRSS

GDI

Windows

CPU OS

Windows2000/XP

Windows 2000/XP

2000/XP

DOS Windows

OS/2

LAN Manager

POSIX

(9

-

(9

WIN32

Windows 3.1

MS-DOS

POSIX

OS/2

1-13

"

"

"

"

Windows 2000/XP

(

)

Windows

2000/XP

"

"

(

)

DOS CPU DOS4.0
 DOS DOS PC
 DOS PC
 Windows DOS
 DOS

1.5.2 Windows

1 Windows

Microsoft 1975
 1983 11 Microsoft Windows WindowsXP Windows
 20
 Windows 90% PC
 GUI 1981 Apple Microsoft Xerox
 Lisa 1983 Macintosh 1984
 Microsoft 1983 11 Windows
 1985 11 Windows
 11
 1.01 1987 Windows 2.0
 1990 Windows 3.0
 1992 4 Windows 3.1 Windows DOS
 Windows 1.x Windows 3.x DOS
 1995 8 Microsoft Windows 95
 DOS Windows95
 Microsoft Windows 97 Windows 98 Windows 98 SE Windows Me
 Microsoft Windows Millennium Edition Windows3.x Windows 9x
 Windows
 Windows 2000 Windows NT
 PC

(7)

(8)

(9)

3 Windows 9x

Windows 9x Windows95 Windows 97 Windows 98 Windows 98 SE
Windows Me Windows

Windows 9x

(1) 32

16 (2)

32

(3) " "

Windows

(4) USB(Universal

Serial BUS) AGP(Accelerated Graphics Port) ACPI(Advanced Configuration and Power
Interface) DVD (5) MPEG WVA MPEG AVI

Apple Quiet Time (6)

Internet (7)

(8) FAT32

4 Windows NT

Windows RISC
CPU /

1993 Windows NT New Technology
Windows (1)

CPU SMP

(4) API (2) (3)32

1.5.3 UNIX

1 UNIX

UNIX
 Kenneth Lane Thompson Dennis MacAlistair Ritchie 1969 DEC
 PDP-7 1971 PDP-11 1973 Ritchie
 BCPL(Basic Combined Programming Language M.Richard 1969)
 C UNIX C 3 UNIX
 1974 7 "The
 UNIX Time-Sharing System" Communication of ACM
 UNIX 1975 UNIX 6 1978 UNIX
 7 UNIX UNIX 70
 UNIX
 UNIX
 BSD UNIX AT&T UNIX
 AT&T UNIX 1981 System 1983 System
 (PWB) 1984 System .2 1987 System .3
 SCO 1978 PC UNIX XENIX
 UNIX AT&T UNIX SVR3.2
 AT&T UNIX Berkeley UNIX BSD
 Berkely Software Distribution BSD PDP UNIX BSD
 1978 1BSD 2BSD 1979 3BSD 1980
 4.0BSD 4.1BSD 4.2BSD 4.3BSD UNIX BSD
 1993 4.4BSD
 TCP/IP UNIX
 Sun OS Solaris 4BSD
 4.3BSD Sun OS AT&T UNIX SVR3.2 SVR4.0(System V
 Release 4) 1989 UNIX BSD UNIX
 UNIX
 UNIX UNIX
 UNIX (1) C
 (2) (3)
 (4) I/O (5)
 CPU (6)
 Shell (7) (8)
 UNIX
 UNIX IBM AIX SUN Solaris
 Berkeley UNIX BSD DEC Digital UNIX
 Compaq Tru64 UNIX HP HP-UX SGI Irix
 SCO SCO UNIXWare Open Server AT&T SVR
 20 90 UNIX 100
 UNIX IEEE UNIX

POSIX
 UNIX
 UNIX
 Portability Guideline

UNIX
 (ACM Turing Award)

2 Solaris
 SunMicrosystem

Intel 80x86
 Solaris

SUN
 1M
 20 80
 SUN
 OS 2.0

RPC
 NFS SUN NFS Network File System
 NFS C/S
 Remote Procedure Call

SUN 1988
 SPARC
 SUN OS 4.1.3

SUN OS 4.0
 Intel 1990

MOTOROLA 680X0
 UNIX
 SUN OS 1982
 SUN OS 1.0
 1 MIPS
 SUN
 1985

NFS C/S
 NFS C/S
 (1)

(2)

(3)

MOTOROLA 680X0
 SUN OS 4.1 1992
 ASMP

SUN

1992 SUN

Solaris 2.0 Solaris 2.0
 SUN OS 5.X 1992 Sun
 2.6 1998 Sun
 2.8

SVR4
 SUN OS 4.X
 Sun
 Sun

Solaris 1.X Solaris 2.X
 Solaris 2.1 Solaris
 64 Solaris2.7

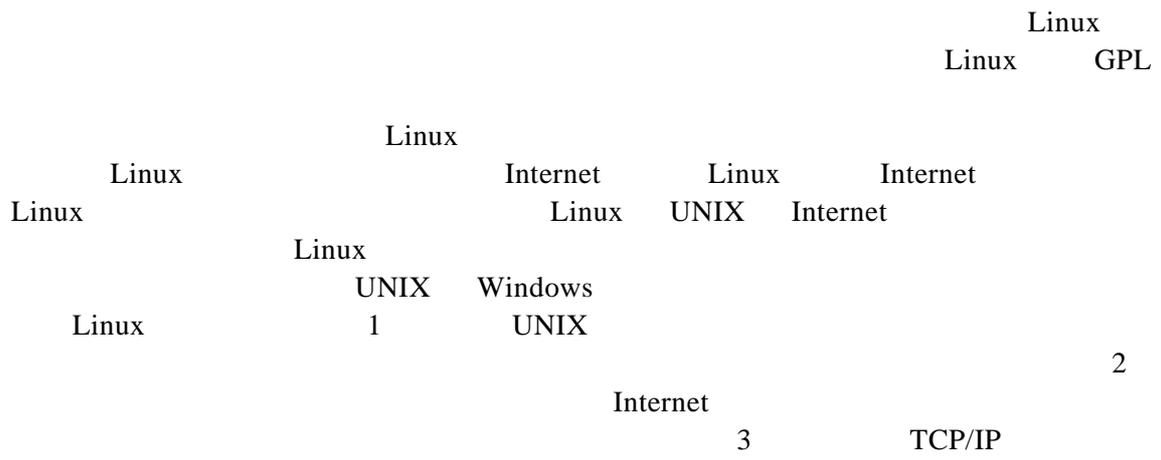
1.5.4 Linux

1

Free Software or Freeware
 License

POSIX
 XPG X-open
 UNIX
 SVID System V Interface Definition
 UNIX
 C
 1983 ACM
 (Software System Award)
 Ritche Thompson
 Solaris UNIX SVR4
 32 Solaris2.x
 SUN OEM SPARC
 SMP
 C/S
 BSD 4.1
 SUN OS 1.0
 SUN
 1985
 SUN
 NIS
 1984
 UDP RPC

Stallman free Richard
 1 2 0
 GPL Bill Gates
 " "
 GNU GNU is not UNIX Richard
 Stallman GPL
 1984 MIT Stallman GNU
 FSF Free Software Foundation GNU
 UNIX UNIX GNU
 gcc/gcc++ shell Linux GNU
 gcc/gcc++ Objective
 C Free BSD Open FORTRAN BSD C BSD email BIND Perl Apache
 TCP/IP IP accounting HTTPserver Lynx Web
 GPL
 GPL
2 Linux
 Linux Linus Torvalds 1991
 Linus
 Internet
 Linux Internet
 Linux
 Linux GNU
 Linux 1998 Linux Internet
 Linux Windows NT
 IBM Intel Oracle Sun Compaq
 Linux Linux
 Linux Linux
 Linux
 Linux UNIX
 UNIX



2⁵²

UNIX

Internet

RS/6000 AIX FTP email Web

CAD CAE

3 OS/390 VM DOS/VSE

IBM S/390(System/390) CMOS OS/390

VM Virtual Machine DOS/VSE Disk Operating System/Virtual

Storage Extended 2000 12 IBM Z900 2001

3 OS/390 ZOS

IBM S/360 S/370 S/390

IBM S/390

70% S/390

IT S/390

OS/390 S/390 G6

1600MIPS

OS/390 MVS Multiple Virtual Storages 1996 IBM

OS/390 1.1 1998 IBM OS/390 2.5 OS/390 2.7

OS/390 X/Open UNIX

OS/390 S/370 TCP/IP

LAN

OS/390

Client/Server

OS/390 S/370 S/370 MVS/ESA390

Enterprise System Architecture ESA 10 240MB 256

ESA/390 LPAR Logical Partitioning

CPU 20 LPAR CPU

OS/370 MVS/ESA370 MVS/ESA390

IBM VM DOS/VSE

4 OS/400

AS/400 64 RISC OS/400

SSP IBM System Support Program

AS/400 IBM

AS/400 OS/400

OS/400 OS/400

OS/400

C

C++ Cobol RPG Java

Licensed Internal Code LIC IBM

AS/400

LIC CPU

OS/400 IBM AS/400

IBM DB2 4.4 AS/400 OS DB

OS/400

5 PC

PC	DOS	OS2	Windows				
OS/2	Microsoft		IBM	1987		PS/2	
			IBM			OS/2 2.0	1994
						Workplace Shell	
32			OS/2 Warp			Windows 95	1996
		OS/2	OS/2 Warp Server4.0			OS/2 Warp Server SMP	
		TCP/IP	6			I/O Navigator	Java
OS/2		1				2	16
32	CPU	3		4GB	4		
		5					6
		API				7	MS-DOS
	MS-DOS		OS/2				

1.5.6

1	Mach						
	Mach		Carnegie-Mellon		Richard Rashid		Accent
1984		Accent			Mach		
	Mach	UNIX		UNIX			
		ARPA		Mach			
		Berkeley UNIX			ARPA		
	Mach		1986		CPU VAX 11/784		
		IBM PC/RT	Sun3		1987 Mach		Encore
	Sequent				IBM DEC HP		

2 Mac OS

Mac OS Apple Macintosh
Apple Mac OS
IBM PC Mac OS Apple 1
2 3 4
5 6 7
Macintosh

3 Netware

Novell PC 1983
IBM Apple UNIX Dec
Netware NOS Networking Operating System
Netware DOS OS/2 MAC UNIX

4 Netware Netware5
Minix
UNIX UNIX AT&T UNIX 7
UNIX
UNIX

Vrije A. S. Tanenbavm UNIX
Minix Minix AT&T

UNIX' Minix C Minix
Minix " Small is Beautiful" Minix
Minix2.0
TCP/IP 4GB 5 200
Minix Internet USENET
Minix
Minix
Minix
Linus Torvalds Minix
Linux

1.6

()

20

60

CPU

CPU

()

()

(SMP) / /

/

Windows 2000/XP

/

DOS Windows UNIX
Mach Mac OS

Solaris Linux IBM
Netware Minix

,

- 1.
2. ?
3. ? ?
4. ?
5. ?
- 6.
7. ? ?
8. ? ?
- 9.
- 10.
- 11.
12. I/O I/O
13. ?
- 14.
15. ? ?
- 16.
- 17.
18. ? ?
19. ? ?
- 20.
- 21.
- 22.
- 23.

- 24.
- 25.
- 26. ?
- 27.
- 28.
- 29. ? ?
- 30. ?
- 31. ?
- 32.
- 33.
- 34. ?
- 35. ?
- 36. / ?
- 37.
- 38. MS-DOS UNIX Windows
- 2000/XP VM/370 Mach
- 39. Windows 2000/XP
- 40. Windows 2000/XP
- 41. Windows 2000/XP ?
- 42. ? Windows 2000/XP
- ?
- 43. /
- 44.
- 45.
- 46.

1		1MB		200KB		200KB
	I/O	80%	1MB	CPU		
2						A
	B		A	50ms	100ms	
50ms	100ms		B	50ms	80ms	
100ms		(1)		CPU		
		(2)	A B	CPU		
3		A B C		I/O		
	A		B		C	
	C ₁₁	30ms	C ₂₁ =60ms		C ₃₁ =20ms	
	I ₁₂	40ms	I ₂₂ =30ms		I ₃₂ =40ms	
	C ₁₃	10ms	C ₂₃ =10ms		C ₃₃ =20ms	

4 CPU I/O(I1,I2)

Job1 I2(30ms) CPU(10ms) I1(30ms) CPU(10ms) I2(20ms)

Job2 I1(20ms) CPU(20ms) I2(40ms)

Job3 CPU(30ms) I1(20ms) CPU(10ms) I2(10ms)

CPU I1 I2 Job1 Job2 Job3

CPU I1 I2 (1)

(2) CPU (3)I/O

5 CPU I/O(I1,I2)

Job1 I2(30ms) CPU(10ms) I1(30ms) CPU(10ms)

Job2 I1(20ms) CPU(20ms) I2(40ms)

Job3 CPU(30ms) I1(20ms)

CPU I1 I2 Job1 Job2 Job3

CPU (1)

(2) CPU (3)I/O

6 3 A B C A B C

A (20) I/O(30) (10)

B (40) I/O(20) (10)

C (10) I/O(30) (20)

I/O()

CPU ?

7 3 A B C A B C

A (30) I/O(40) (10)

B (60) I/O(30) (10)

C (20) I/O(40) (20)

I/O()

CPU ?

8 3 A B C A B C CPU

I/O

A 60 20 30 10 40 20 20 (ms)

I/O2 CPU I/O1 CPU I/O1 CPU I/O1

B 30 40 70 30 30 (ms)

I/O1 CPU I/O2 CPU I/O2

C 40 60 30 70 (ms)

CPU I/O1 CPU I/O2

I/O

CPU ?

9	A	(CPU)10	()5	(CPU)5	()10
	B	()10	(CPU)10	()5	(CPU)5
		()10	A	B	CPU
10				2ms	?
	60HZ	CPU			?

SMP

64

SMP

I/O

Dec HP IBM SUN SGI

SMP

MIMD

Cluster

MIMD

RAID

SMP

2.1.2

●

●

●

● I/O

● I/O

●

I/O

I/O

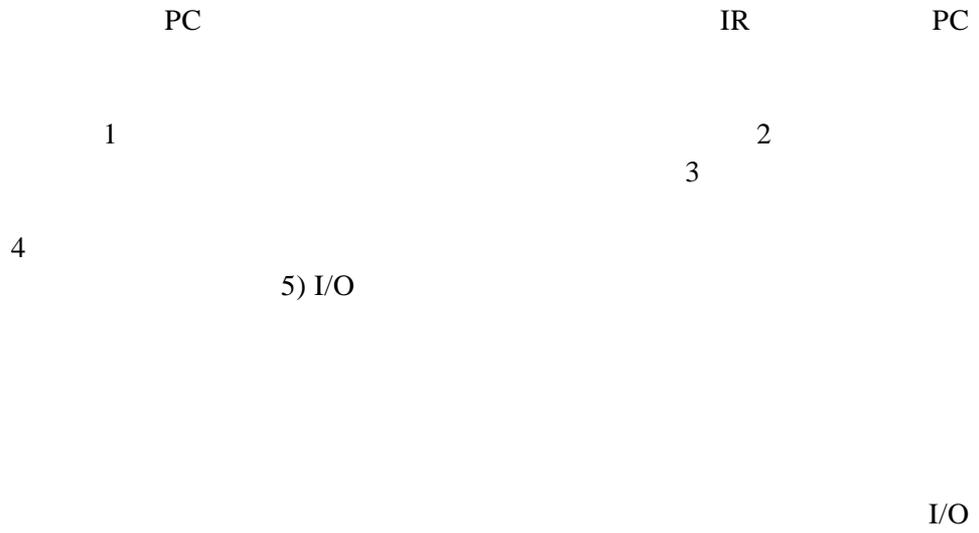
/

PC Program Counter

IR Instruction Register

I/O

2.1.3



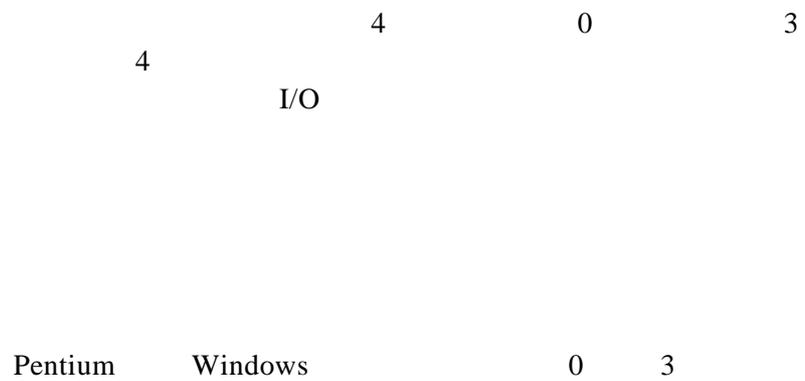
Privileged Instructions

PSW

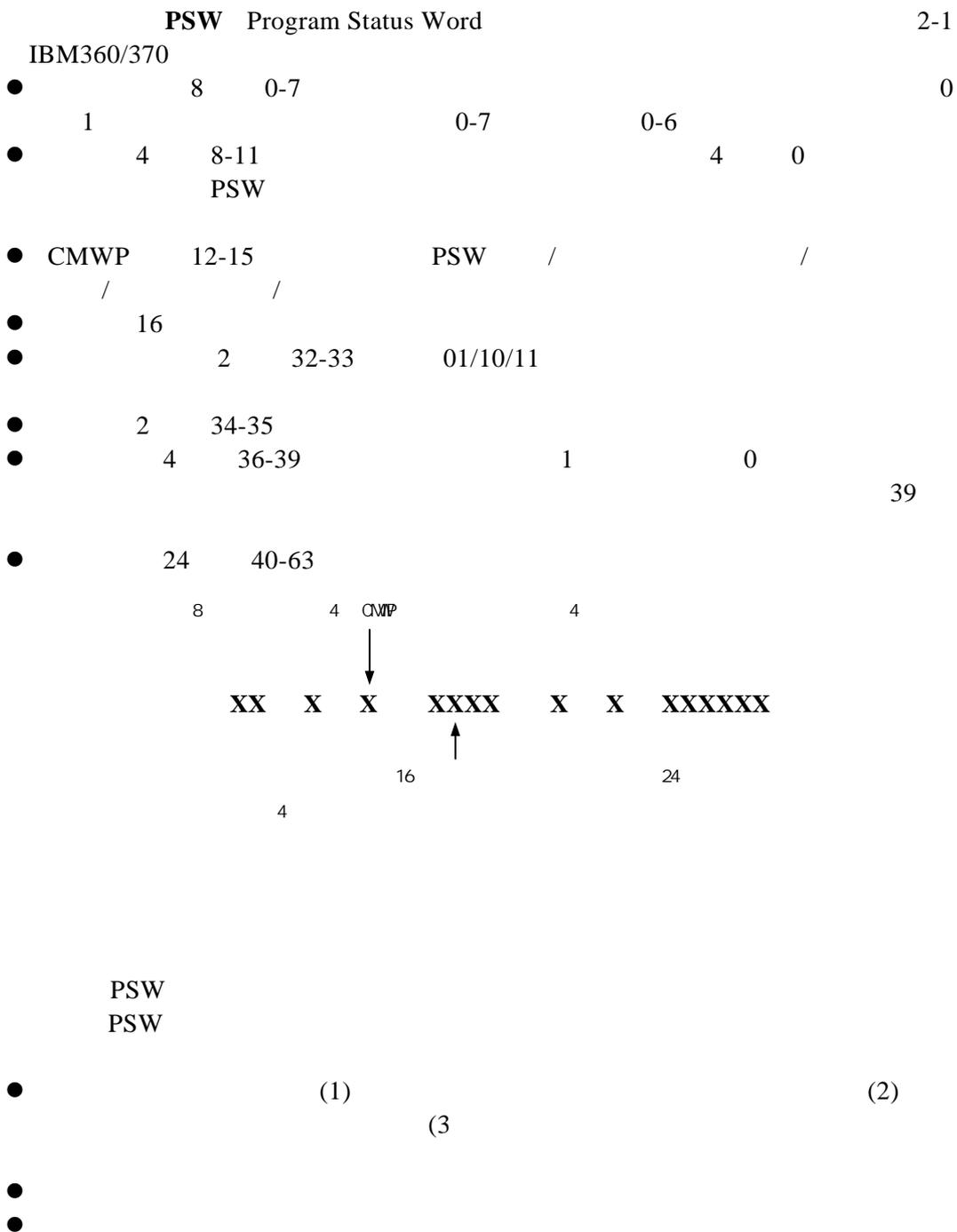
2.1.4

Intel Pentium

- 0
- 1
- 2
- 3



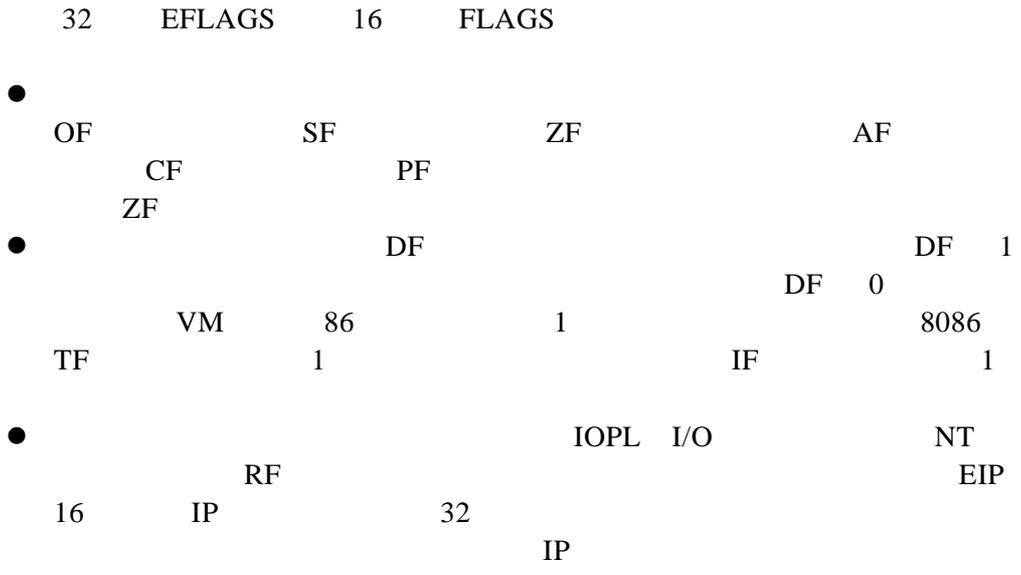
2.1.5



Intel Pentium

EFLAGS

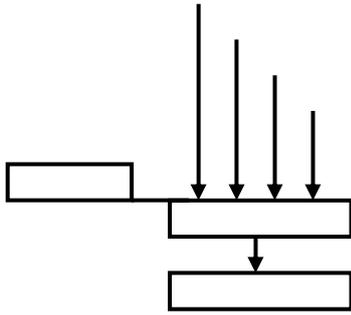
EIP



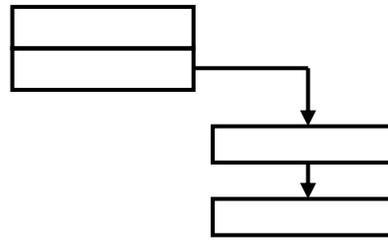
-
-
-
-

0

2-2



2-2



-
-

I/O

IBM

Windows2000/XP

()

CPU

CH3

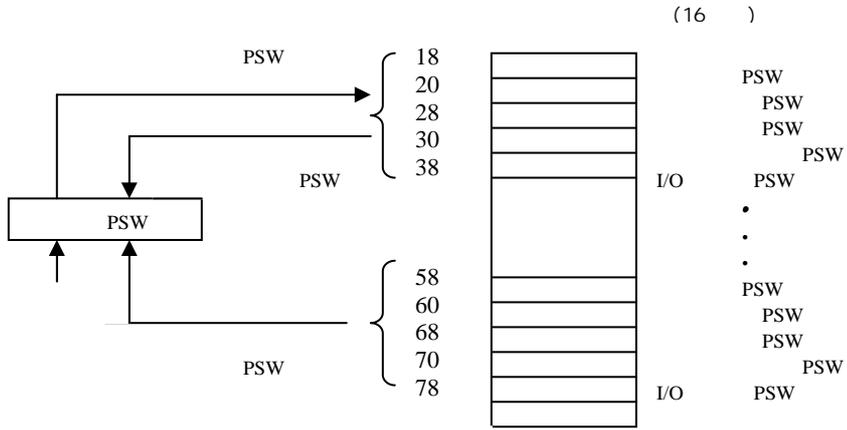
2.2.3

•

•

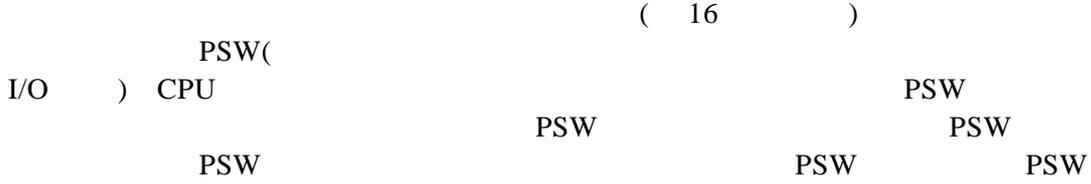
•

" 0"



2-3 IBM

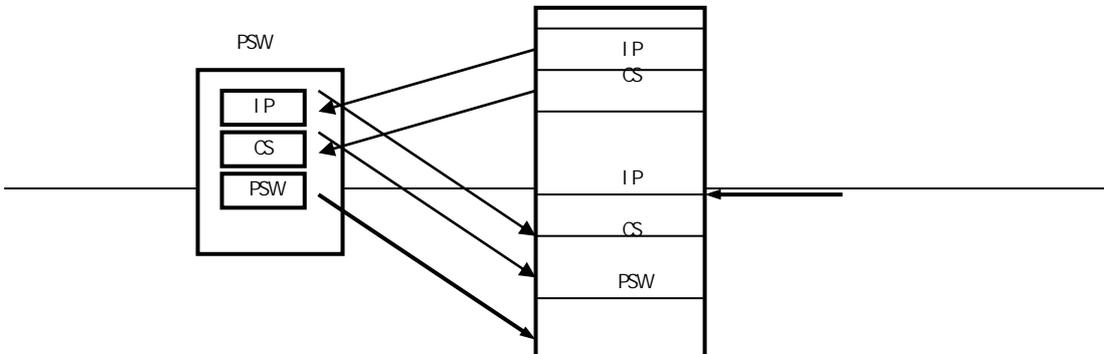
2-4 IBM



PSW() CPU
I/O)
PSW

PSW

PSW
PSW
PSW



•

•

•

•

2

2

?

on

on < > < >

on fixed overflow go to LA
LA

on fixed overflow go to LB
LB LA

2-5

0	1	0	0
0			
1			
...			
N			

2-5

0

?

" 0"

" 0"

2

0

1

•

•

on

•

on

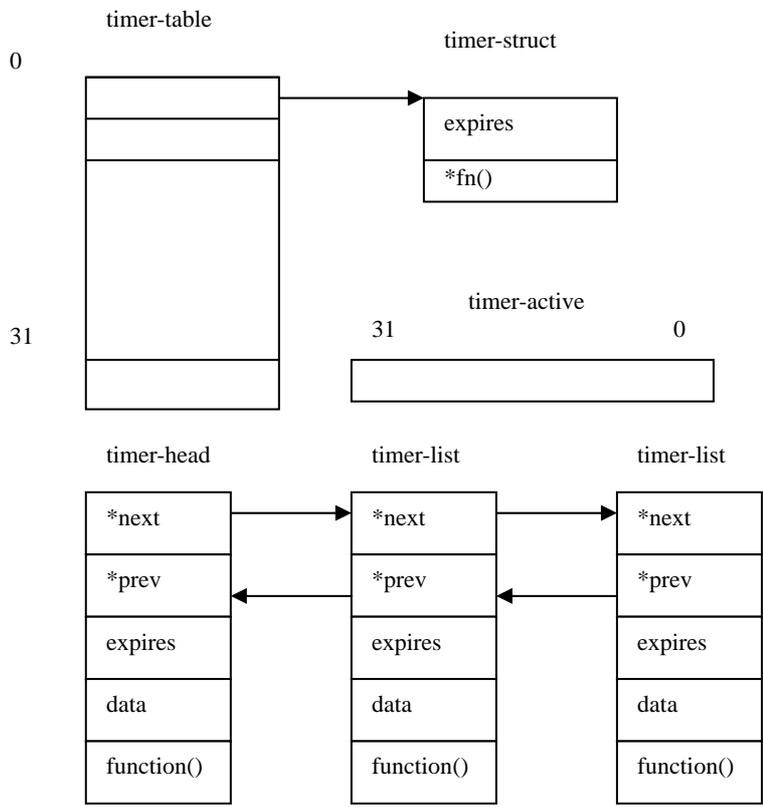
•

•

3

1

Linux tick) 2-6 jiffies timer-active 32 (clock Linux timer-struct timer-list



2-6

jiffies bottom half bottom half timer-active timer-list

Linux CPU Linux (accounting time) SIGALRM tick SIGVTALRM

- real
- virtual

- profile
 - SIGROF
 - Linux
 - task-struct
 - profile
 - timer-list
 - Real
 - real
 - Virtual
 - tick

2

4 I/O

1 I/O

2 I/O



n

2.2.6

1

?

I/O

I/O

I/O

2

1

0

3

CPU

?

•

—

I/O

I/O

•

•

2.2.7 Windows 2000/XP

1 Windows 2000/XP

Pentium Windows

" "

I/O

I/O

2-7

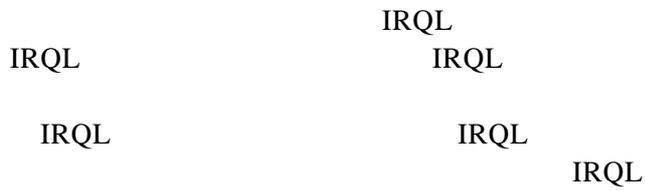
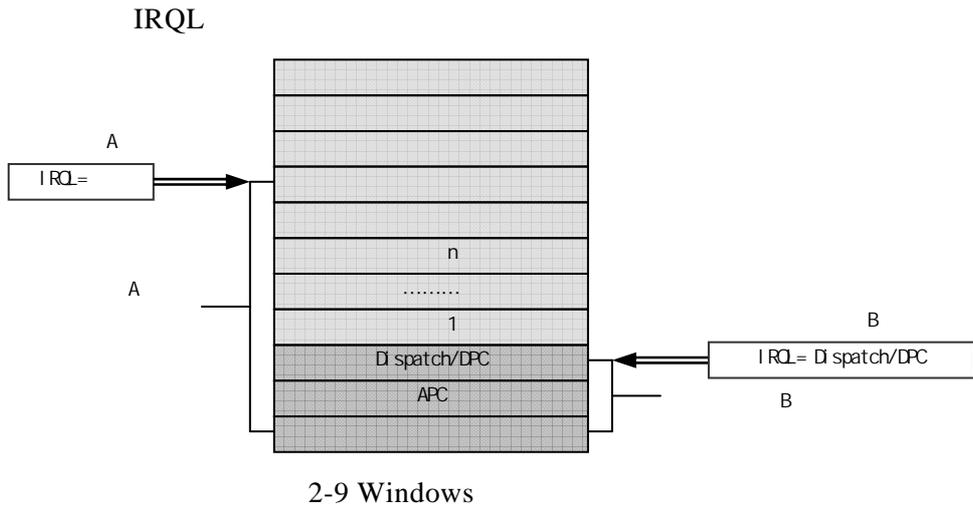
" Trap Frame "

Interrupt Service Routine

2 Windows 2000/XP

1

Level	IRQL	Windows2000/XP IRQL Interrupt Request
2-8	x86	IRQL IRQL
	IRQL	IRQL



3 Windows2000/XP

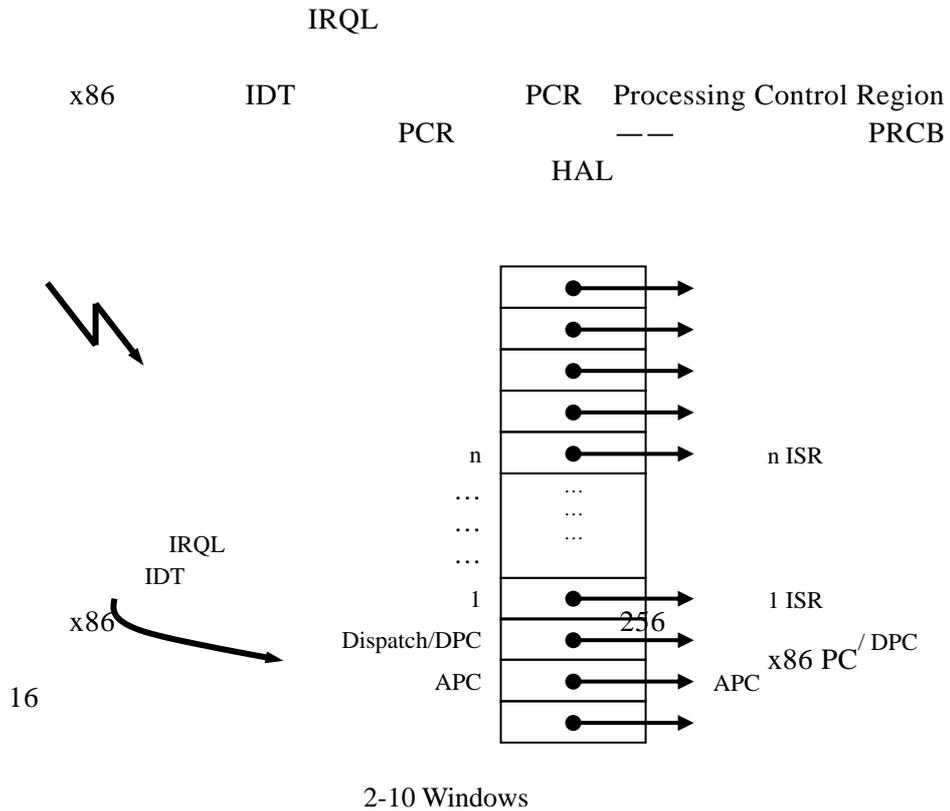
1)

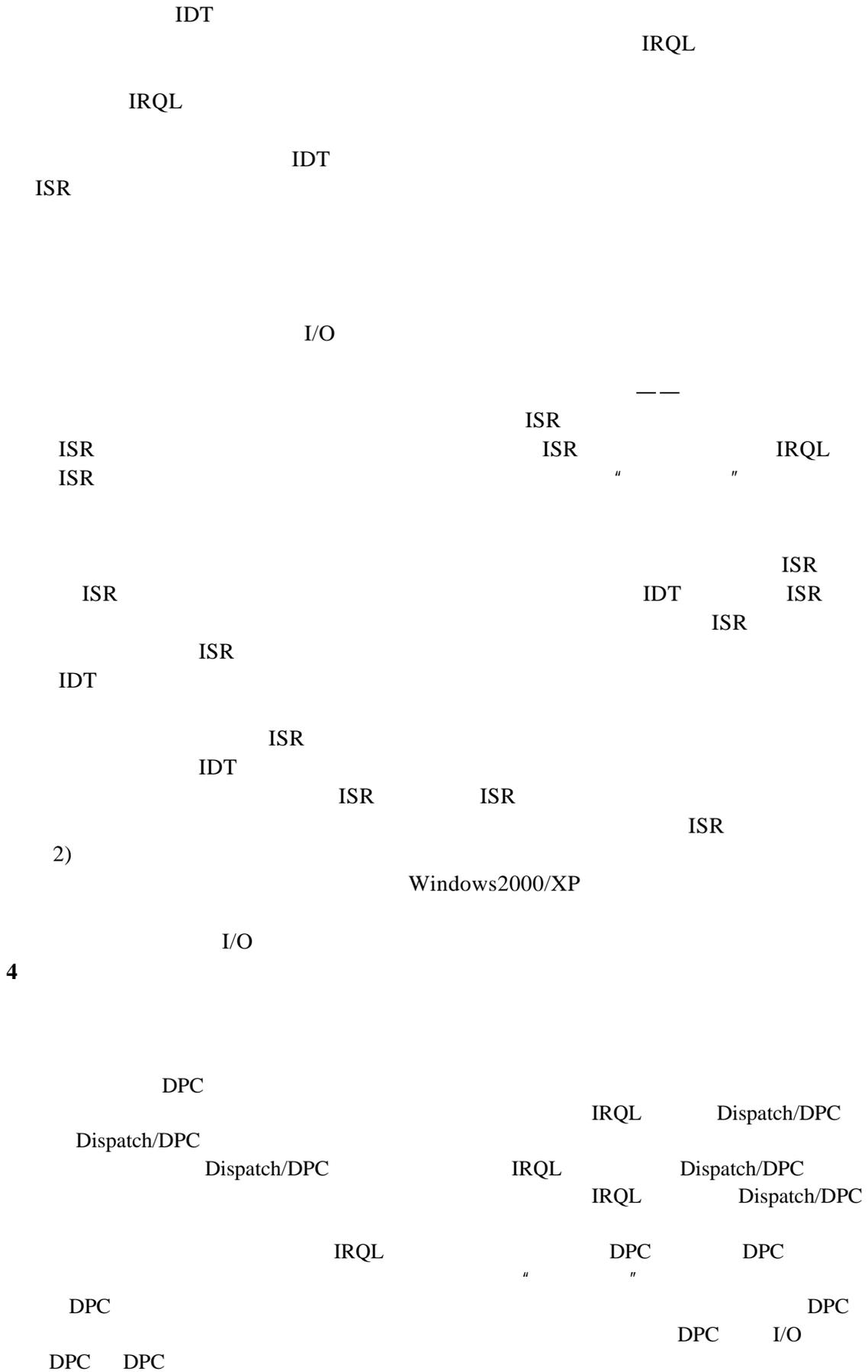
IRQL

Windows2000/XP

IDT Interrupt Dispatch Table

2-10





" DPC DPC DPC " DPC
 " DPC DPC DPC DPC
 DPC DPC DPC Dispatch/DPC
 APC IRQL IRQL IRQL IRQL IRQL
 APC APC IRQL DPC DPC
 DPC DPC
 IRQL
 Dispatch/DPC IRQL 0
 DPC IRQL DPC
 DPC
 APC ()
 2 IRQL APC APC ()
 APC APC APC APC APC APC APC APC
 DPC DPC APC APC APC APC APC APC APC
 APC APC APC APC APC APC APC APC APC APC
 " " " " APC
 APC ()
 APC APC
 I/O APC APC
 POSIX APC APC POSIX
 POSIX APC I/O
 I/O I/O
 Win32 API ReadiEX WriteFileEX QueueUserAPC APC
 ReadiEX WriteFileEX I/O
 APC I/O

5 Windows 2000/XP

WIN32

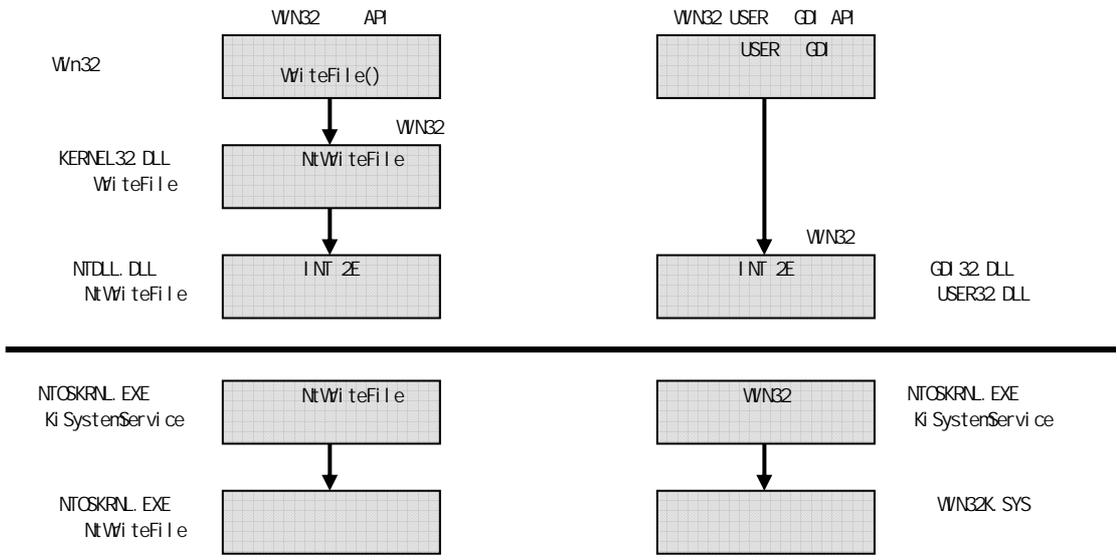
Structure Exception Handling

" " " "

6 Windows 2000/XP

x86 INT 2E System
Service Dispatcher 2-11
(System Service Dispatch Table)

Windows2000/XP
NTOSKRNL.EXE
WIN32 WIN32K.SYS WIN32 USER GDI
WIN32 WIN32 USER GDI
WIN32 USER GDI
Windows2000/XP NTDLL.DLL
DLL NTDLL
WIN32 USER GDI USER32.DLL
GDI32.DLL
2-12 KERNEL32.DLL WIN32 WriteFile NTDLL.DLL
NtWriteFile NtWriteFile
NTOSKRNL.EXE KiSystemServic2afc>Tj/TT2 1 Tf2XZ



2-12 Windows 2000

2.2.8 Solaris

1 Solaris

SPARC ULTRA SPARC Solaris

system call

interrupt

trap

2 ULTRA SPARC

SPARC ULTRA SPARC
SPARC

SPARC ULTRA SPARC

3

TBA
4M

8

2

1

Solaris

SPARC ULTRA SPARC
SPARC ULTRA SPARC

1

2

		3		4
5		6		
SPARC	ULTRA SPARC			
Interrupt Level)			PIL	PIL(Processor
ULTRA SPARC				5 0
	4		3	
3 Solaris				
Solaris			Solaris	
	1	15	15	
	9	9		9
			Solaris	

1)

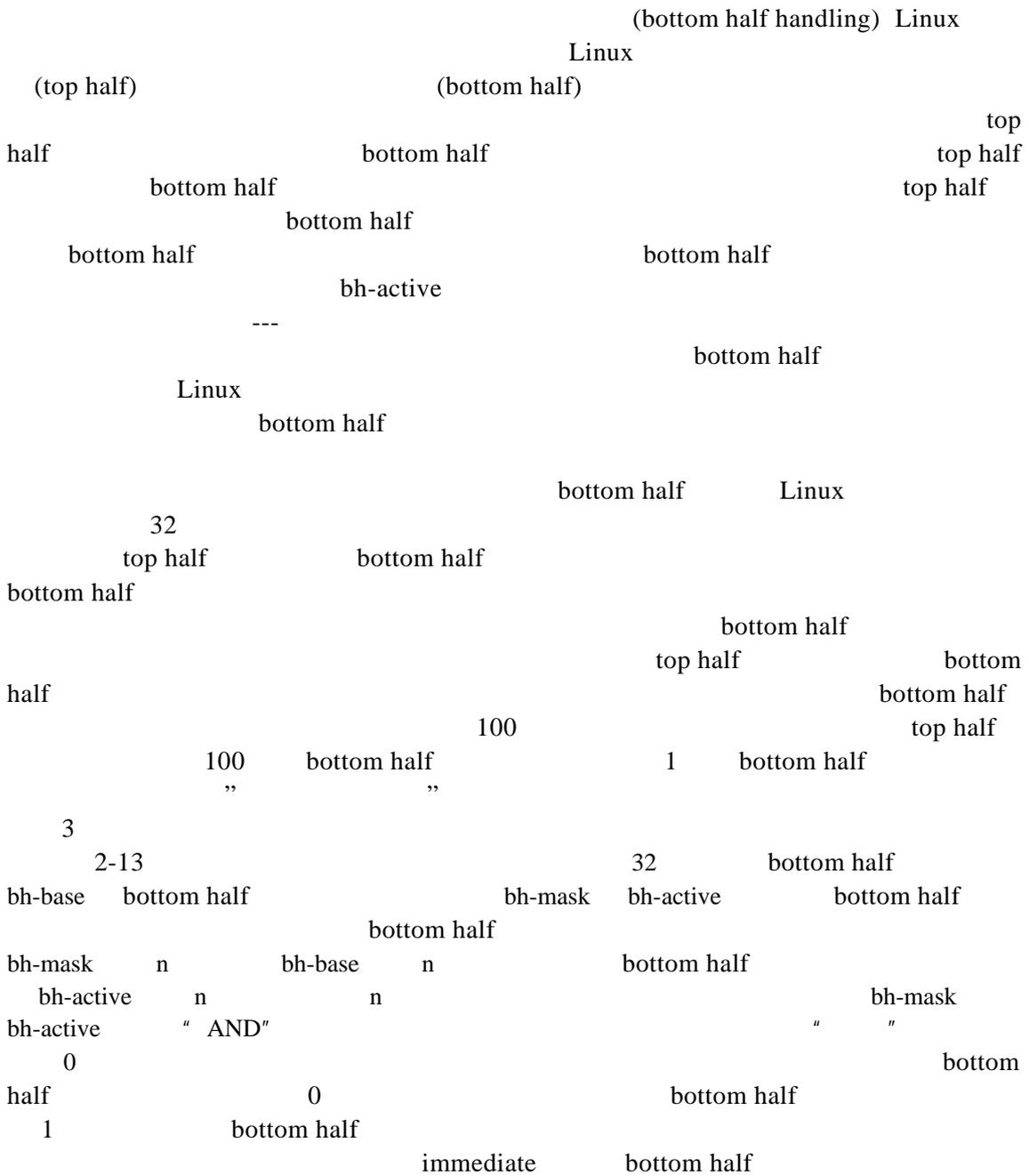
C

2)

3)

)
)
2

(
(



(tq-immediate)

4

(1) (TIMER-BH)
(2) (CONSOLE-BH) (3)
(NET-BH) (4) (IMMEDIATE-BH)

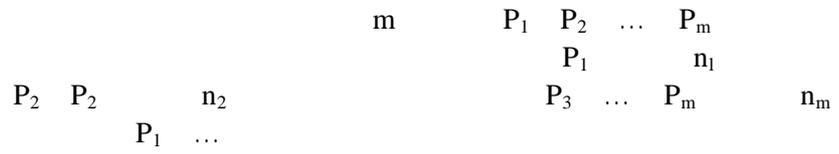
bh-active immediate
bh-active immediate bottom half
1
bh-active
bottom half 0 1 31
bh-active bottom half
top half bottom half
() timer-interrupt top half do-timer
bottom half timer-bh timer-interrupt CPU
top half do-timer jiffies(1)
() W,X

active

1978

-
-
-
-
-
-

CPU

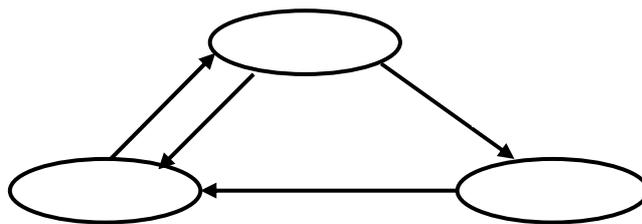


2.3.2

1

- running
 - ready
 - wait
- blocked
- sleep

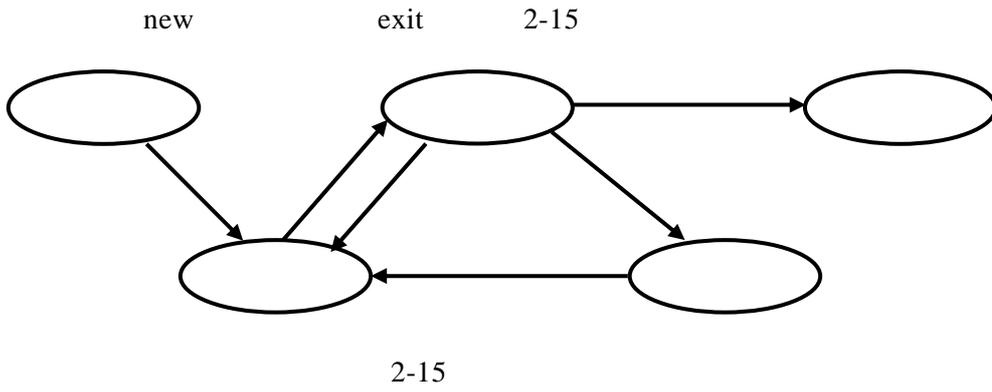
2-14



2-14

- —
- —
- —
- — CPU

2

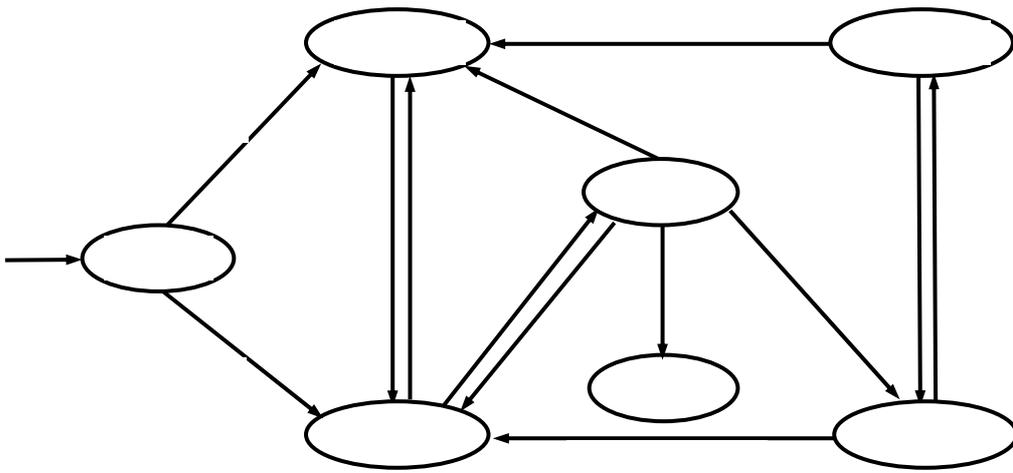


- NULL—
- —
- —
- — NULL
- —
- —

3

suspend

-
-
-
-
-
-



2-16

2-16

ready suspend

blocked suspend

- —
- —
- —
- —
- —

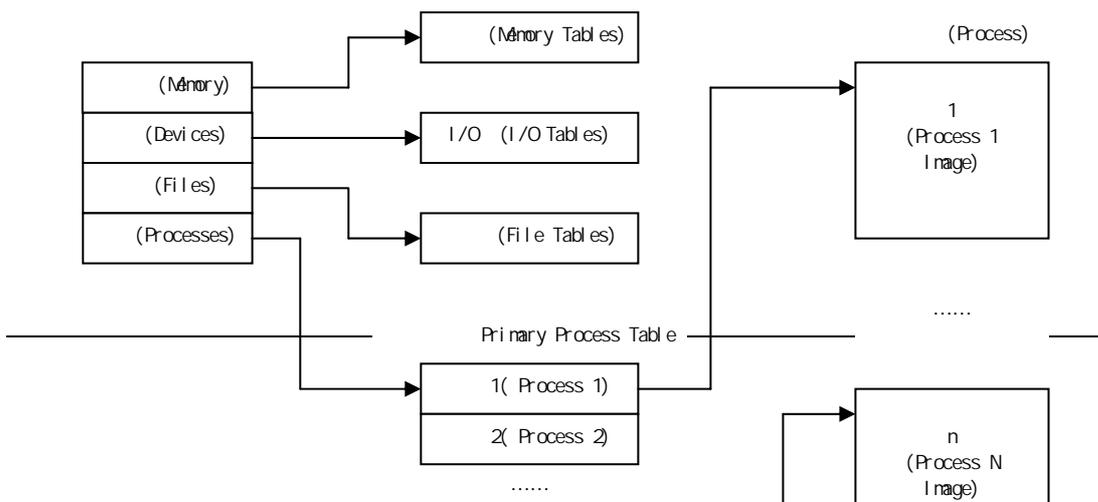
- — CPU
- —
-
-
-
-

2.3.3

1

- I/O
-
-
- I/O
- I/O
- I/O
- I/O
- I/O
- I/O
-

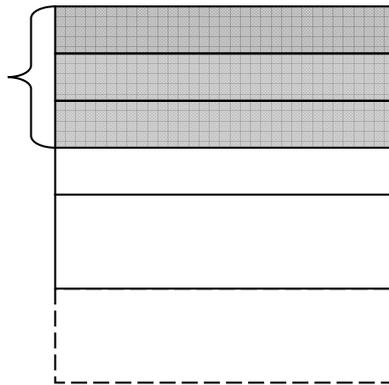
2-17



2

process context

- user -level context
- system -level context
- register context



2-18

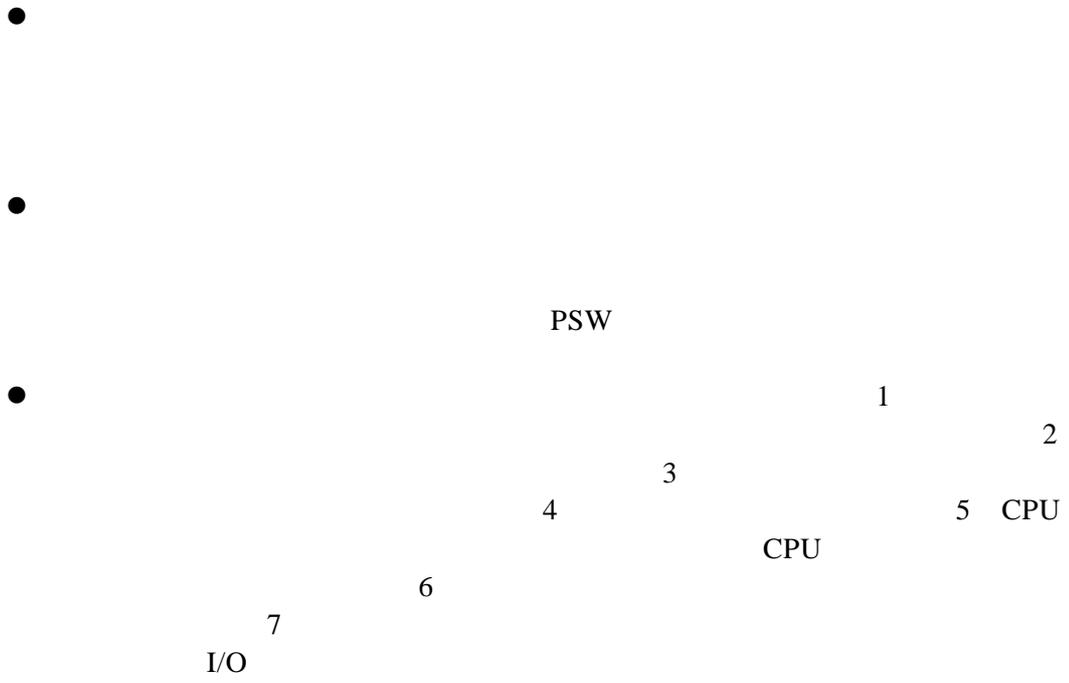
Image

-
-
- /
-

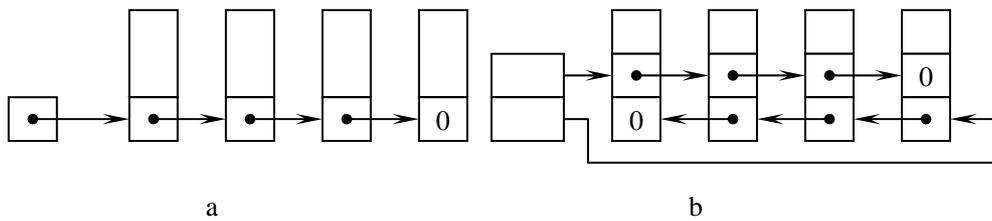
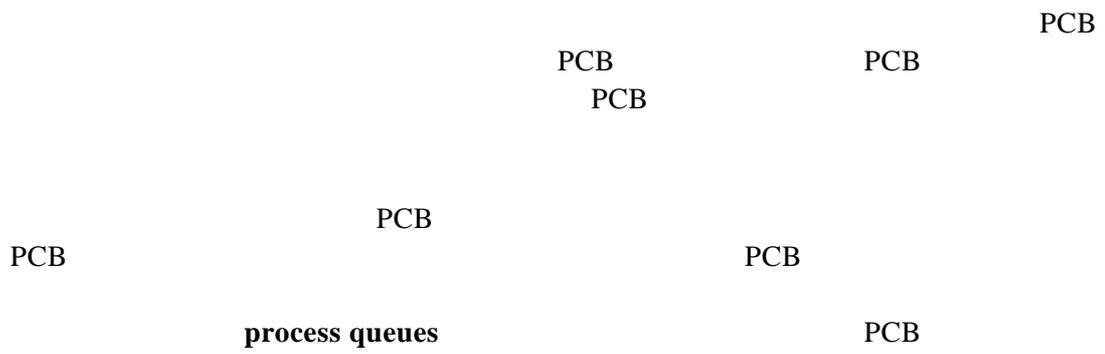
Process

3

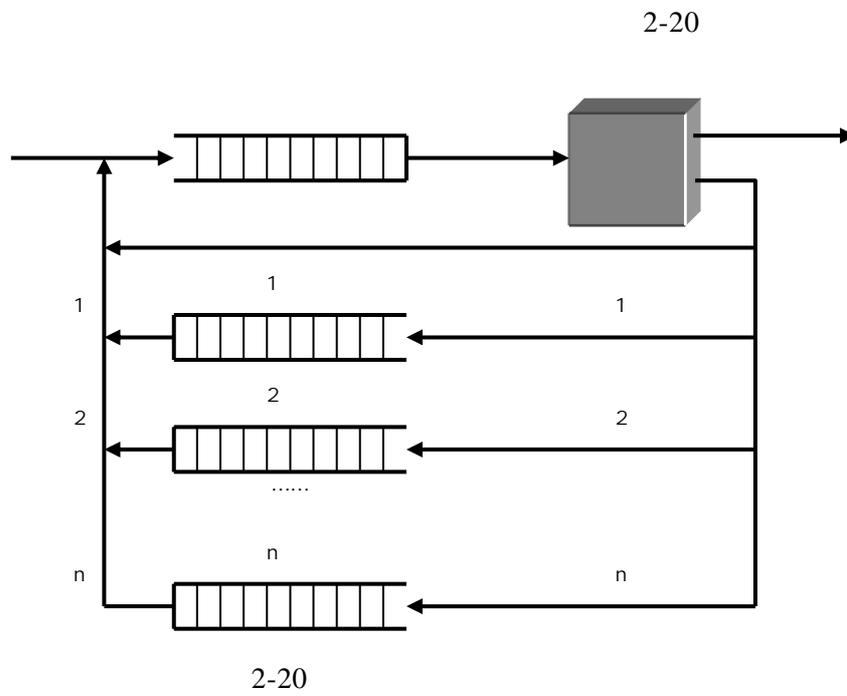
PCB Process Control Block



4

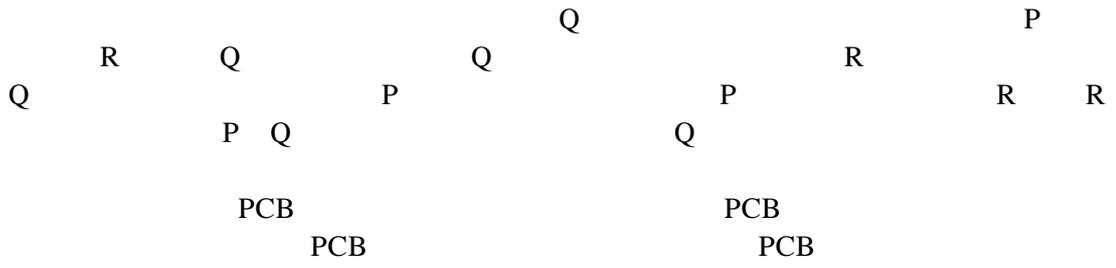


2-19 a b



2-15 b

- 1 0
- 2 0
- 3



2.3.4

context layer

- (1)
- (2)
- (3)
- (4)



CPU
CPU

CPU

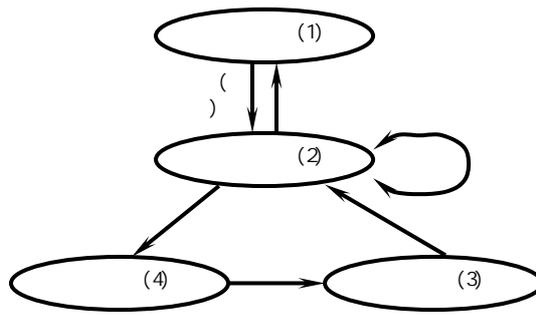


PCB

2-21

-
-
-
-

I/O



2-21

1

2

4

3

2.3.5

(Primitive)

() ()

1

-
-
-
-

" "

parent process

child process

UNIX/Linux

-
-
-
-
-
-
-

task_struct PCB PCB UNIX/Linux

PSW

Linux clone

fork()

fork()

task_struct Linux fs clone fork() clone() files sig

mm
count 1 count 0
copy on

write
mm
mm

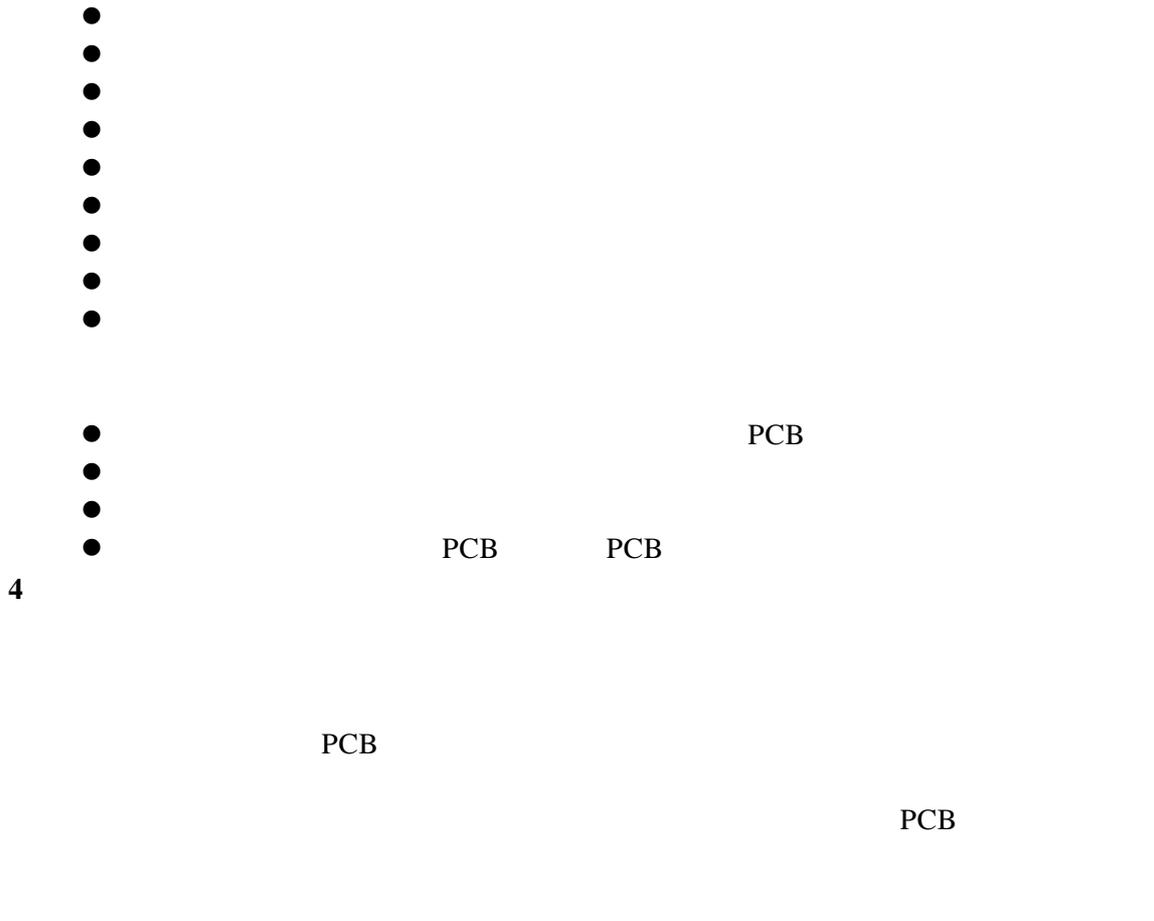
UNIX Solaris
Linux2.4
size-of-memory-in-the-system/kernel-stack-size/2
512MB 512 × 1024 × 1024 / 8192 / 2 = 32768
Linux
2 I/O

-
-
-
-
-
-
-

PCB

UNIX/Linux
) wait() kill() sleep() pause()
SuspendThread ResumeThread
Windows 2000/XP
(suspend count) 1 1 0 1
0 1 1 0

- 3
- -
 -
 -
 -

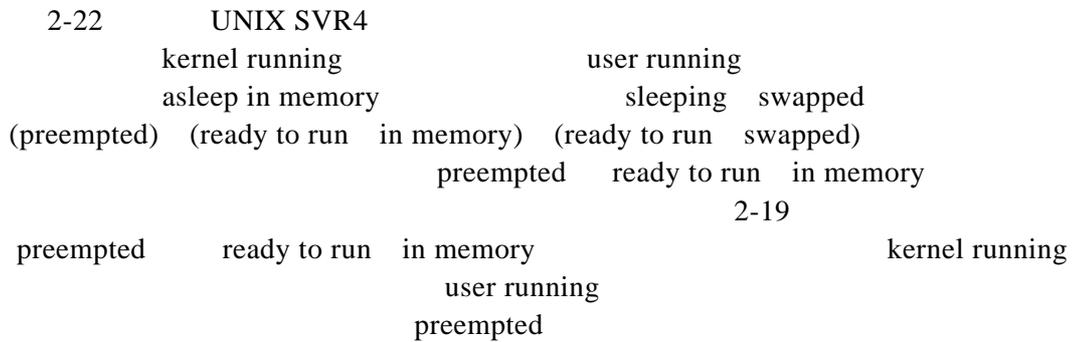


2.3.6 UNIX SVR4

UNIX SVR4

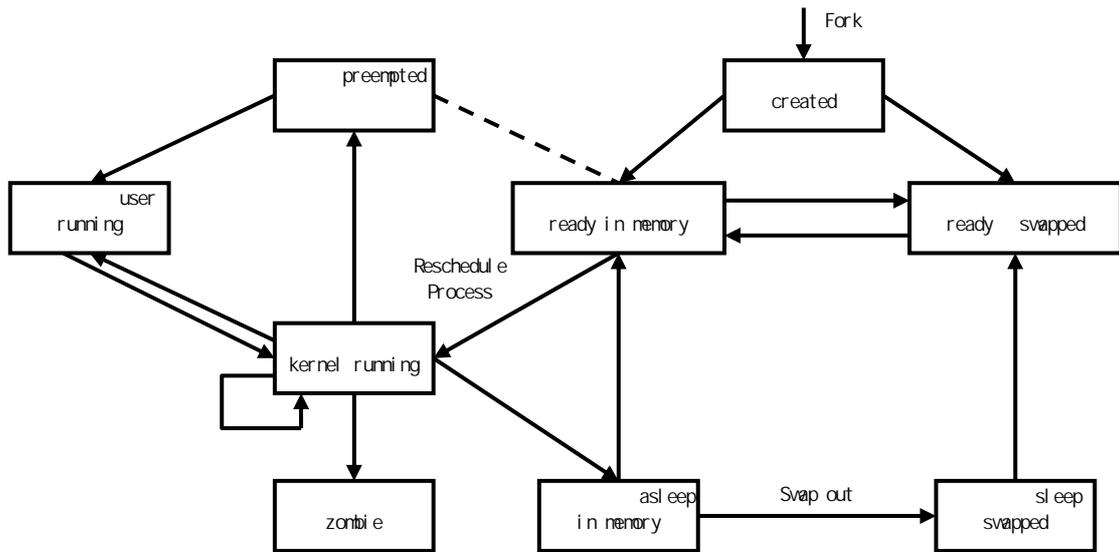
UNIX SVR4

1 UNIX SVR4



- user running
- kernel running
- preempted

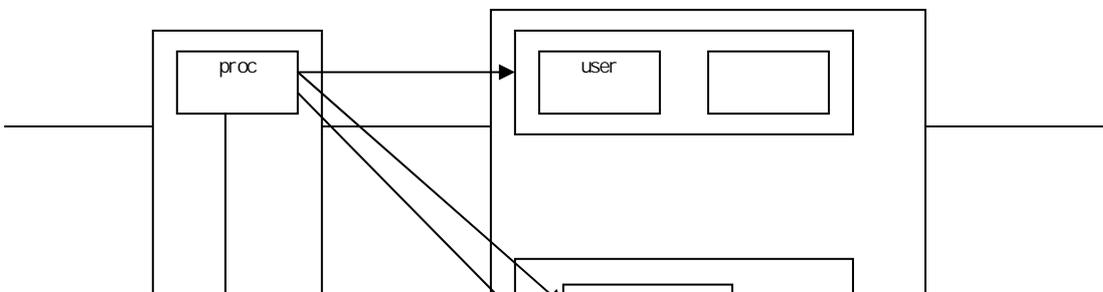
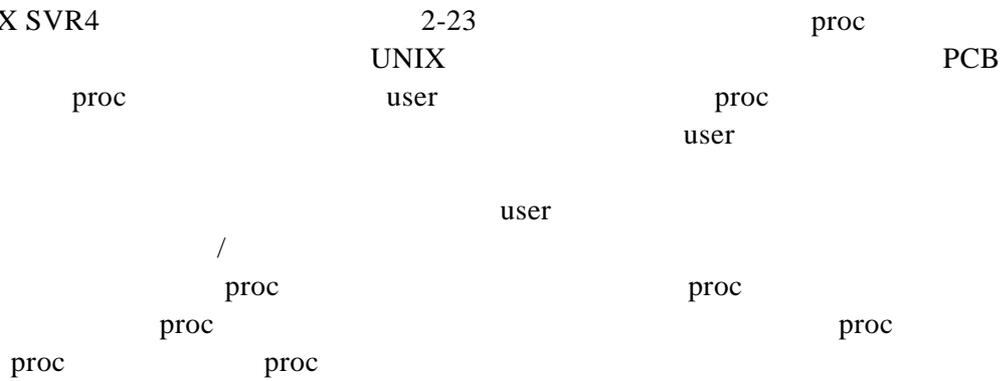
- ready to run in memory
- asleep in memory
- ready to run swapped
- sleeping swapped
- created
- zombie



2-22 Unix SVR4

UNIX 1 init 0 0 swap 1

2 UNIX SVR4
UNIX SVR4



proc / user / text/data/stack / proc

user / / I/O / proc

user / / user / PCB / proc

user / 1024 /

text / UNIX / text

proc user text text text UNIX proc

PPRT Per Process Region Table

/

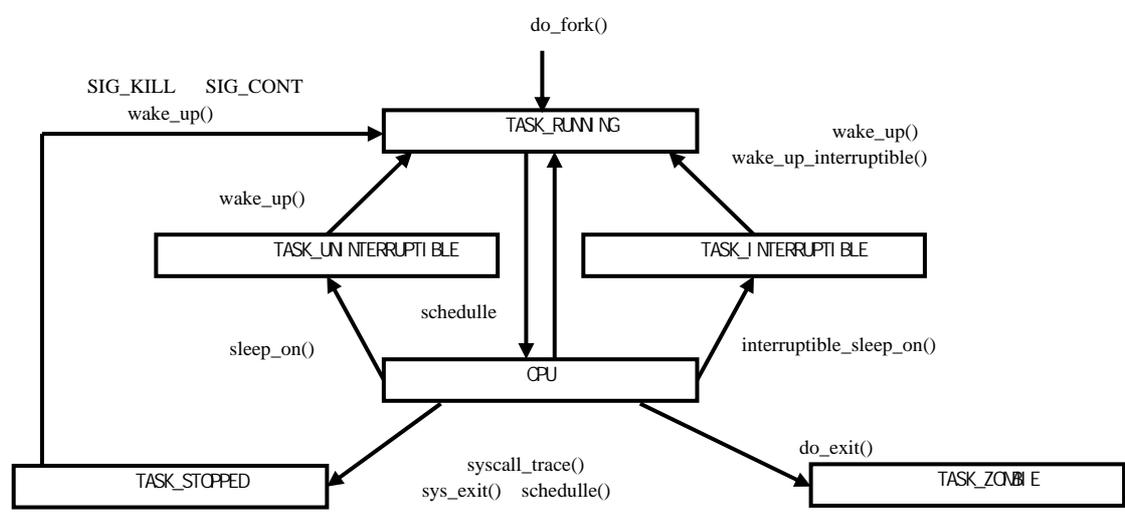
3 UNIX SVR4

UNIX SVR4 fork() pid=fork()
 ●
 ●
 ●
 ●
 ●
 ● 0
 ● fork()
 ● Ready to Run fork()
 ● Ready to Run ready to run
 ● fork() pid 0
 pid 0 pid 0

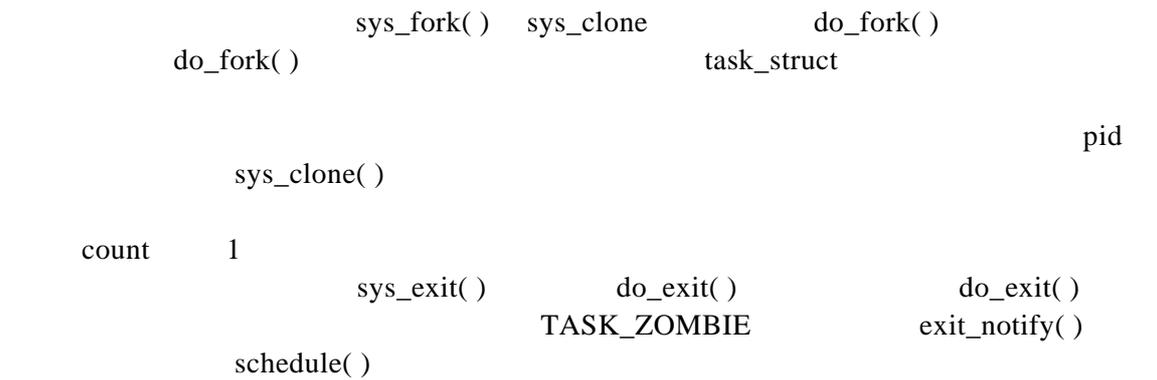
2.3.7 Linux

Linux
 1 Linux
 Linux 2.2.x SMP
 2-24 Linux 6 CPU CPU
 ● TASK_RUNNING current
 ● TASK_INTERRUPTIBLE SIGNAL
 ● TASK_UNINTERRUPTIBLE SIGNAL

- TASK_ZOMBIE
- TASK_STOPPED
- TASK_SWAPPING



2-24 Linux



2 Linux

Linux
 process table NR-TASKS task_struct
 /include/linux/sched.h Linux 2.2.10 512 x86 NR-TASKS
 4096 Linux task-struct
 PCB task-struct

task-struct	
state	(6)
flags	(10)
priority	
rt_priority	
counter	
policy	0 2 1

32

CB

Li nux

sui d, sgi d

ui d/gi d

ui d/gi d

	nswap		
	min_flt, maj_flt		
	cnswap		
	swap_cnt		
SMP	processor	SMP	CPU
	last_processor		CPU
	lock_depth		
	used_math		FPU
	Comm[16]		
	rlim		
	errno		0
	Debugreg[8]		
	*exec_domain personality		i386
	*binfmt		elf, java 4
	exit_code, exit_signal		
	dumpable		memory dump
	di_d_exec		POSIX
	tty_old_group		
	*tty		
	*wait_chldexit		

Linux

(1)current SMP CPU

CPU

(2)init-task 0 PCB

(3)*task[NR-TASKS] PCB

(pid) task[0] 0 init-task

tasks[] PCB for-each-task()

next-task PCB

(4)jiffies Linux 0 10ms

do-timer() 1

(5)need-resched 1 schedule() CPU

(6)intr-count

2.4

2.4.1

process

-
-
-
-
- / C/S

multiple threaded process

" " " "

MS-DOS

UNIX

Solaris Mach SVR4 OS/390 OS/2

WindowNT Chorus JAVA

Solaris thread OS/2 thread

Windows NT thread IEEE UNIX

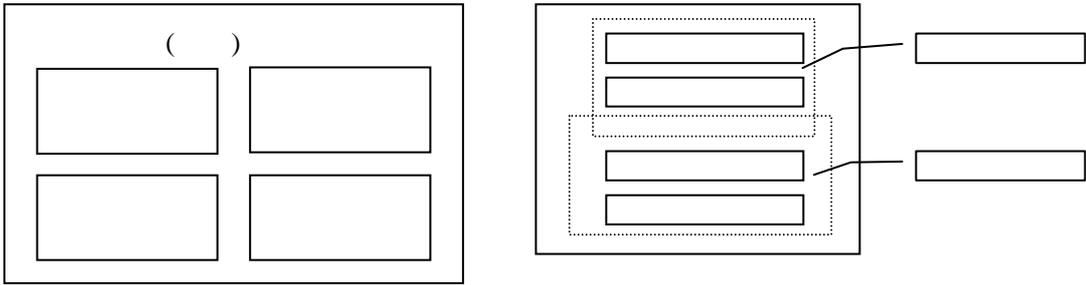
POSIX 1003.4a

2.4.2

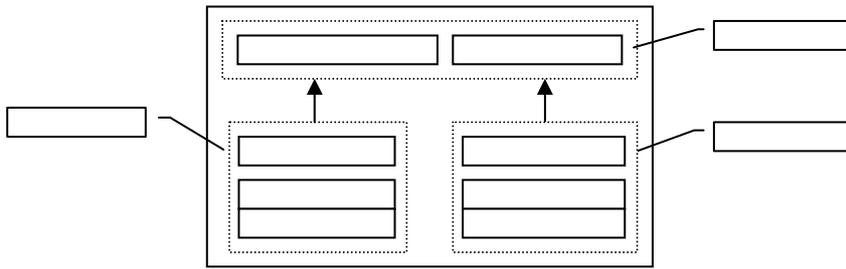
1

2-25

/



2-25



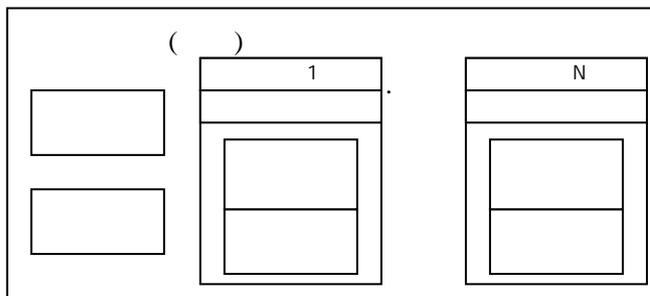
2-26

2-26

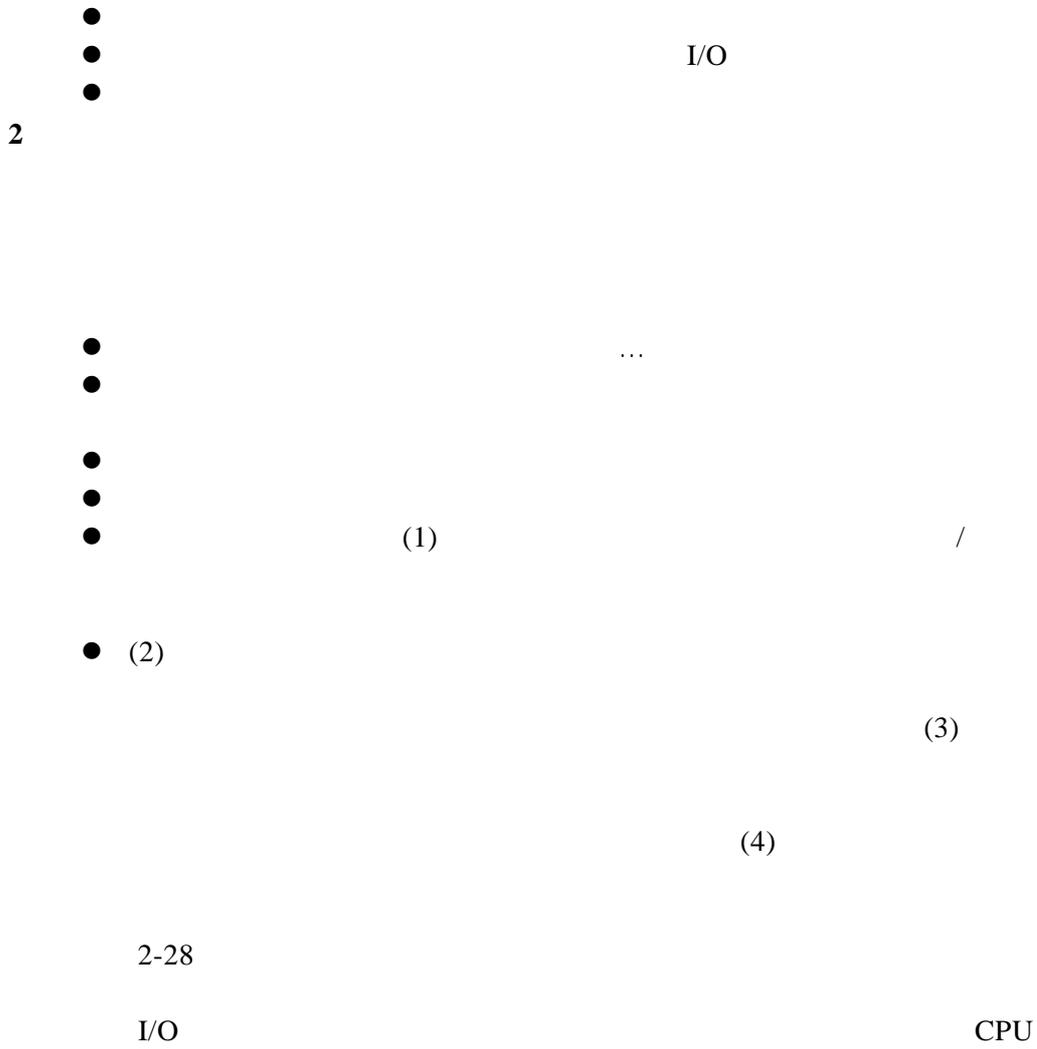
PCB

2-27

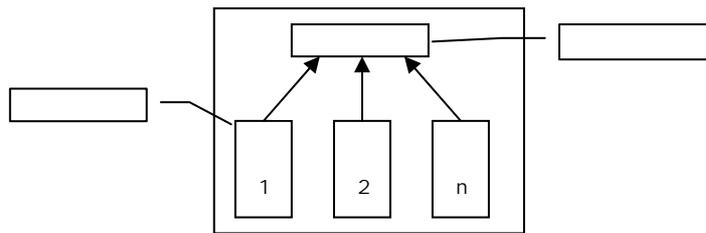
TCB Thread Control Block



2-27



LWP(Light-Weight Process)



" " " " ()
()

() () ()

3

I/O

?

?

Windows

4

()

()

()

()

()

()

(

)

(

)

API

● spawn

● block

● unblock

● finish

TCB

()

()

U

)

()

()

()

API

- -

()
Thread) POSIX P-threads Java
Windows2000/XP OS/2 Mach C-thread
ULT KLT
ULT(User Level
KLT(Kernel Level Thread)
Solaris ()

Intel

4

cache

5

●

●

●

10

●

●

●

●

●

UNIX

Mach

UNIX

50

/

•

• C/S

CPU

CPU

•

•

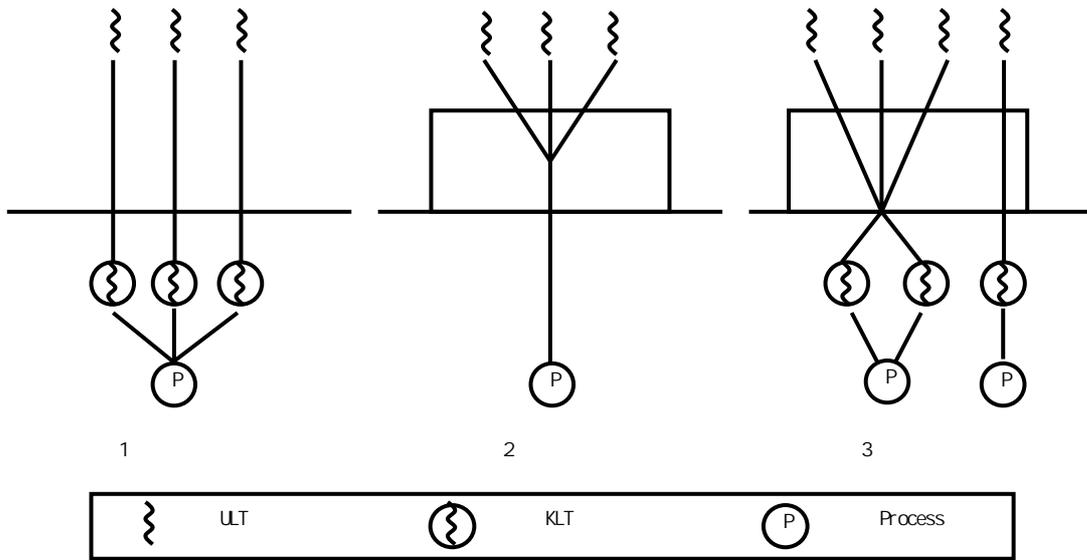
•

GUI

GUI

2.4.3

	ULT	POSIX	P-threads	Java
KLT	Windows2000/XP	OS/2	Mach	C-thread
		Solaris		
2-29				



2-29

1

KLT Kernel Level Threads

KLT

API

Windows 2000/XP

OS/2

TCB

PCB

KLT

KLT

- -

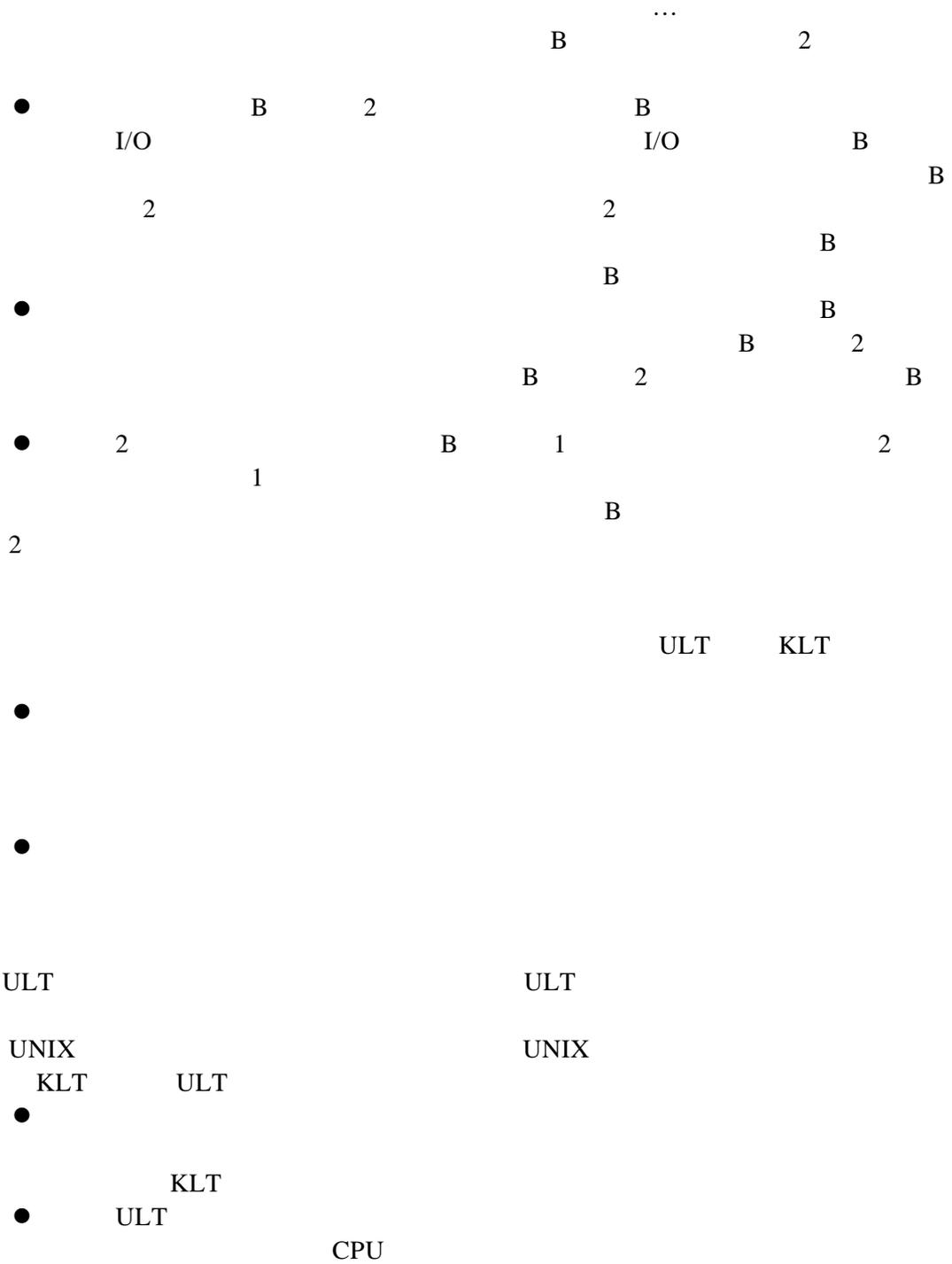
2

ULT User Level Threads

ULT

" " " "

TCB



jacketing

jacketing I/O
 3 jaketing I/O
 ULT/KLT Solaris
 KLT
 ULT KLT
 KLT

2.4.4 Solaris

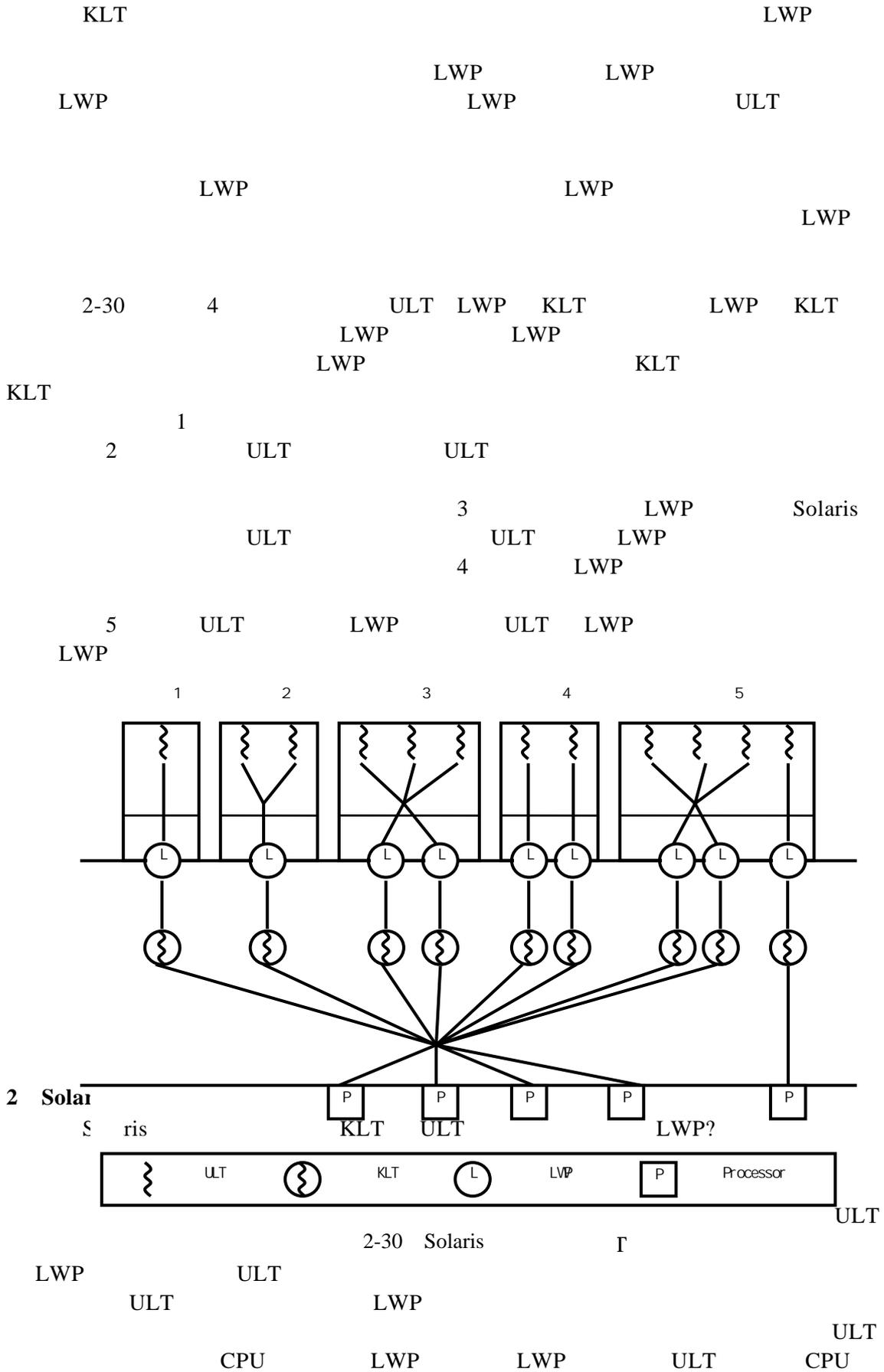
1 Solaris

- Solaris 4 SMP
- Process UNIX PCB
 - User-Level Threads
 - Light Weight Process LWP ULT KLT
LWP ULT KLT LWP
 - Kernel-Level Threads KLT

Solaris

LWP KLT
 KLT LWP ULT LWP ULT KLT
 ULT LWP KLT
 ULT

LWP(LWP ULT ULT
 LWP ULT LWP LWP
 KLT LWP ULT KLT LWP ULT LWP ULT
 LWP ULT ULT
 ULT LWP ULT LWP ULT
 LWP ULT LWP ULT
 ULT KLT LWP ULT LWP
 ULT ULT LWP LWP 5
 ULT 5 LWP LWP
 KLT KLT
 4 LWP 4 ULT KLT
 ULT



LWP

ULT

ULT

LWP

LWP

LWP

Solaris

2-31

(

)

(

) Solaris
UNIX

LWP

LWP
LWP

LWP

(

)

/

/

LWP

KLT

LWP

LWP

3 Solaris

2-32

ULT

LWP

ULT

LWP

ULT

LWP

KLT

ULT

● KLT

ULT
ULT

thread_exit()	
thread_wait()	
thread_get_id()	
thread_sigsetmask()	
thread_sigprocmask()	
Thread_kill()	
Thread_stop()	
Thread_priority()	
mutex_enter()	
mutex_exit()	
mutex_tryenter()	
Sem_p()	P
sem_v()	V
sem_try()	P
rw_enter()	/ /
rw_exit()	/ /
rw_tryenter()	/
rw_downgrade()	write lock read lock
rw_tryupgrade()	read lock write lock
cv_wait()	
cv_signal()	cv_wait()
cv_broadcast()	cv_wait()

2.4.5 Windows 2000/XP

1 Windows 2000/XP

Windows2000/XP
Windows2000

Windows NT4

CPU

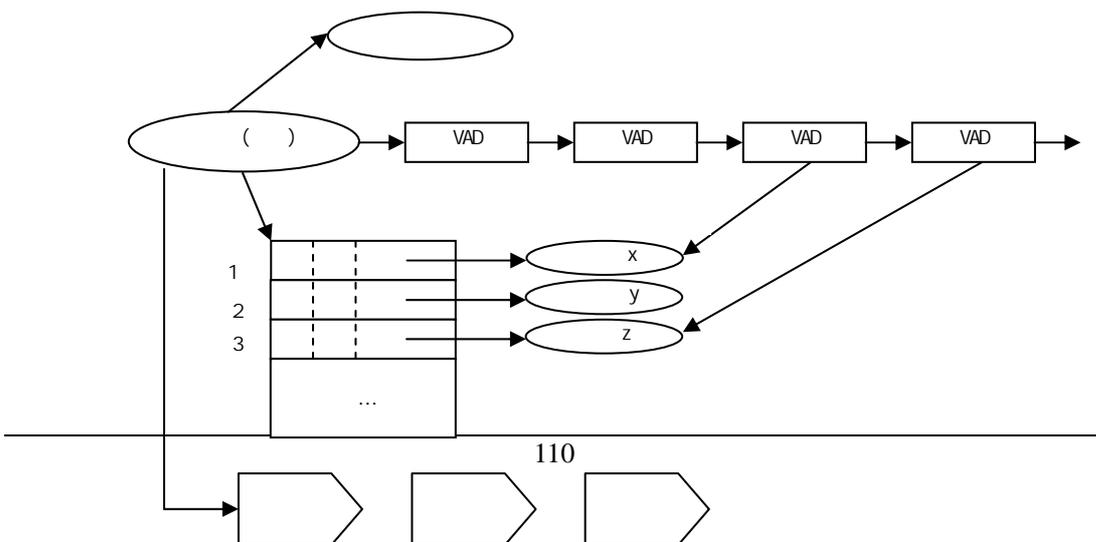
Windows 2000/XP

SMP

C/S

CPU

-
-
-



Windows2000/XP

2-33
access token

2-33

2

(object oriented)

()

()

()

(object class)

(instance)

(encapsulati on)

(subclass)

(superclass)

Windows 2000/XP

(object-based)

Windows 2000/XP

?

Windows 2000/XP

(1)

(2)

Windows2000/XP

/

/

Window2000/XP

EPROCESS

(2-33)

3

EPROCESS

EPROCESS

ETHREAD

PEB

Process Environment Block

EPROCESS

CSRSS

WIN32

WIN32K.SYS

WIN32

WIN32

WIN32 USER

GDI

EPROCESS

● KRPROCESS

●

● CPU

● VAD

●

●

● /

●

●

● PEB

GDI

● WIN32 WIN32
WIN32

● CreateProcess
● CreateProcessAsUser
EXE

● OpenProcess
● ExitProcess
● TerminateProcess
● FlushInstructionCache
● GetProcessTimes

● GetExitCodeProcess

● GetCommandLine
● GetCurrentProcessID ID
● GetProcessVersion

Windows

● GetStartupInfo CreateProcess STARTUPINFO

● GetEnvironmentStrings
● GetEnvironmentVariable
● GetProcessShutdownParameters
● SetProcessShutdownParameters

CreateProcess WIN32 WIN32

3
KERNEL32.DLL Windows2000/XP
●
● Windows2000/XP
●

WIN32
.EXE

WIN32
CSRSS

- GetExitCodeThread
- GetThreadTimes
- GetThreadSelectorEntry
- GetThreadContext CPU
- SetThreadContext CPU
- CreateThread WIN32
-
-
- NtCreateThread ID

WIN32 TEB KeInitializeThread KTHREAD

- WIN32
- ID
- CREATE_SUSPEND

Windows2000/XP 2-35

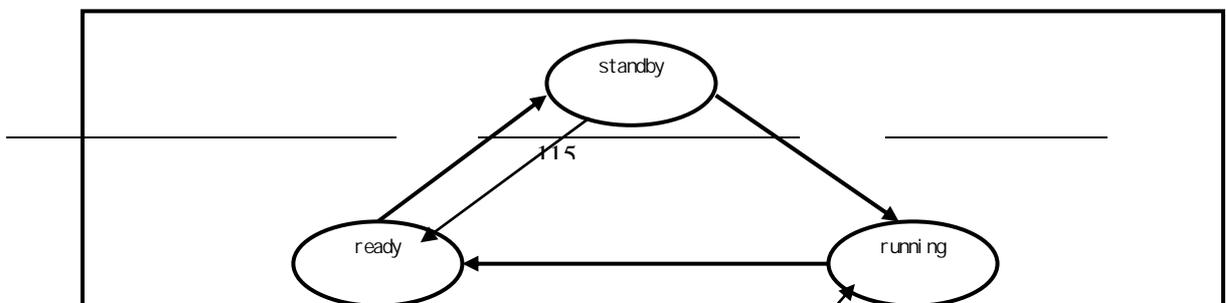
7

- --- CPU
- ---
- ---
- --- 1
- I/O 2 3
- --- ()
- ---
- ---

5

Windows2000/XP " "

Windows NT4



-
- CPU
- CPU
-
-
-

ID SID

SystemParameterInfo

WIN32

- CreateJobObject
- Open JobObject
- AssignProcessToJobObject
- TerminateJobObject
- SetInformationToJobObject
- QueryInformationToJobObject

CPU

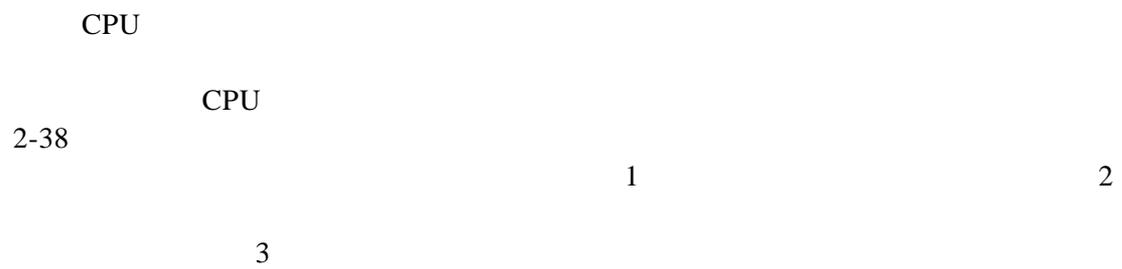
ID

2.5

2.5.1

2-36

2-36



2.5.3

Medium Level Scheduling
Medium-term Scheduling

" "

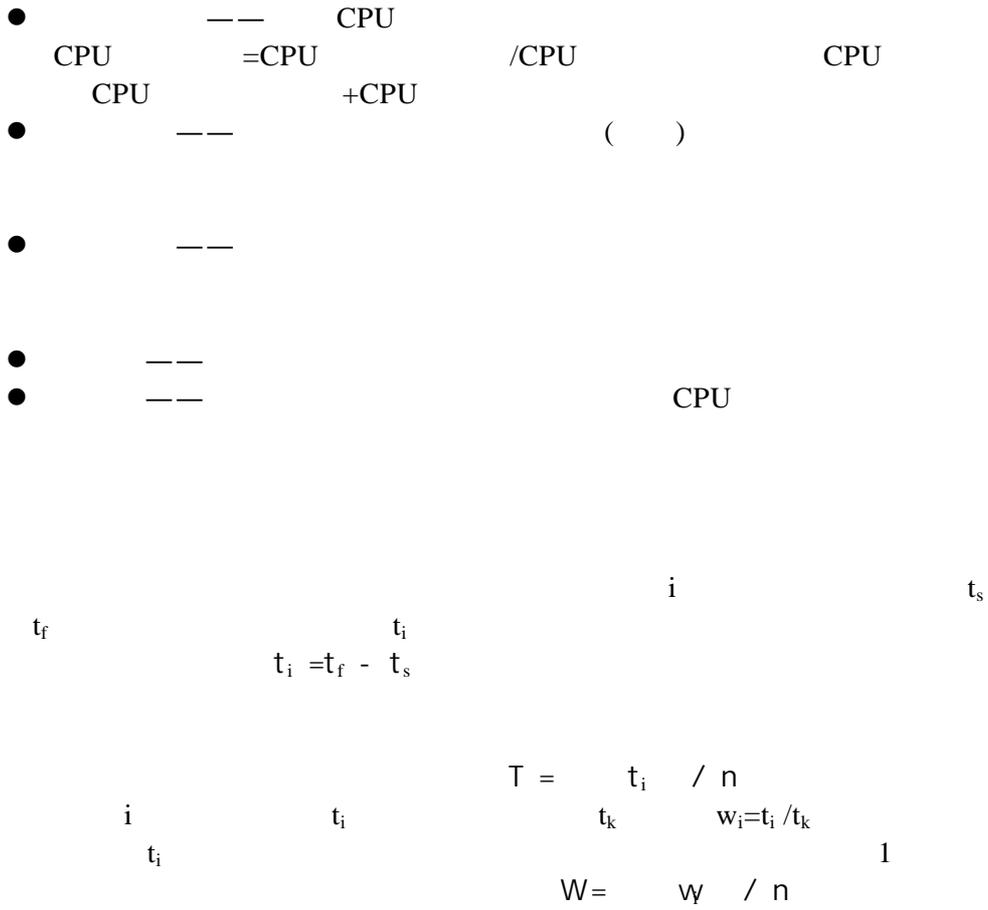
2.5.4

Low Level Scheduling
Short_term Scheduling

CPU (dispatcher)

2.5.5

scheduling algorithm scheduler



2.6

2.6.1

JOB

Job Step

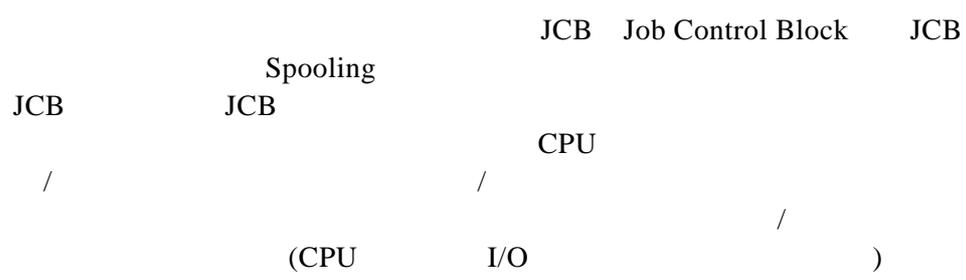
" " " " " "

" " "

CPU " " "

UNIX

2.6.2



-
-
-
-

P

9

2.6.3

$$T = \frac{T_1 + T_2 + \dots + T_n}{n} \quad \text{CPU}$$

2.6.4

1

FCFS First Come First Served

FCFS

CPU

	CPU
1	9
2	4
3	10
4	8

SJF

2 4 1 3

$$T = 4+12+21+31 / 4 = 17$$

$$W = 4/4+12/8+21/9+31/10 / 4 = 1.98$$

FCFS

$$T = 9+13+23+31 / 4 = 19$$

$$W = 9/9+13/4+23/10+31/8 / 4 = 2.51$$

SJF

FCFS

FCFS

SJF

SJF

CPU

CPU

SRTF Shortest Remaining Time First

CPU

		CPU
1	0	8
2	1	4
3	2	9
4	3	5

J1	J2	J4	J1	J3
----	----	----	----	----

0 1 5 10 17 26

Job1 0 Job2 1 Job1
7 JOB2 4 Job1 Job2

SRTF

10-1 + 1-1 + 17-2 + 5-3

$$/4=26/4=6.5$$

SJF

$$7.75$$

Highest Response Ratio First

SJF

/ =1+ /

HRRF

HRRF

		CPU
1	0	20
2	5	15
3	10	5
4	15	10

SJF

1 3 4 2

T= 20+25+35+50 /4=32.5

W= 20/20+25/5+35/10+50/15 /4=3.2

FCFS

1 2 3 4

T= 20+35+40+50 /4=36.25

W= 20/20+35/15+40/5+50/10 /4=4.1

HRRF

- 1 1 20
- 1 1+15/15 1+10/5 1+5/10 3
- 3 5 1+20/15 1+10/10 2
- 15
- 2 4 10

T= 20+25+40+50 /4=33.75

W= 20/20+25/5+40/15+50/10 /4=3.4

HRRF

SJF FCFS

I/O

5

6

2.7

2.7.1

- dispatcher
-
-
-

2.7.2

1



$$p\text{-pri} = \min \{ 127, (p\text{-cpu}/16 + PUSER + p\text{-nice}) \}$$

$$p\text{-pri} = -100 \sim +127$$

$$p\text{-nice} = PUSER - 100$$

$$p\text{-pri} = 127 - (p\text{-cpu}/16 + p\text{-nice})$$

p-cpu 1 1
 p-cpu < 10
 p-cpu 0
 p-cpu > 10
 200ms 200ms 10

(1) p-pri p-cpu

(2) 200ms p-cpu 0 p-pri

UNIX 100

UNIX

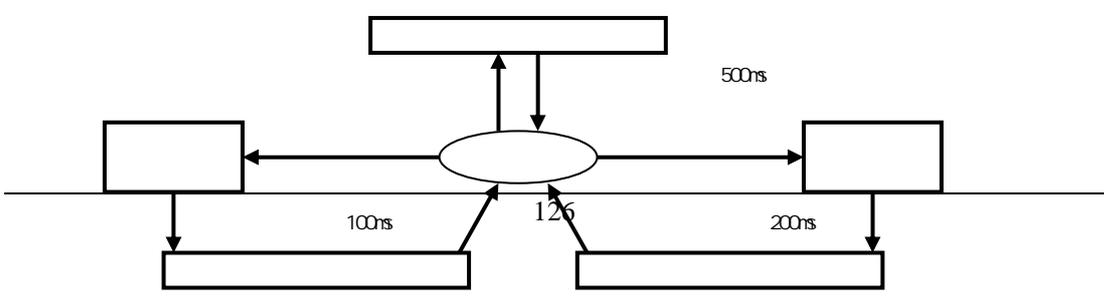
$$p\text{-pri} = (p\text{-cpu}/2) +$$

4

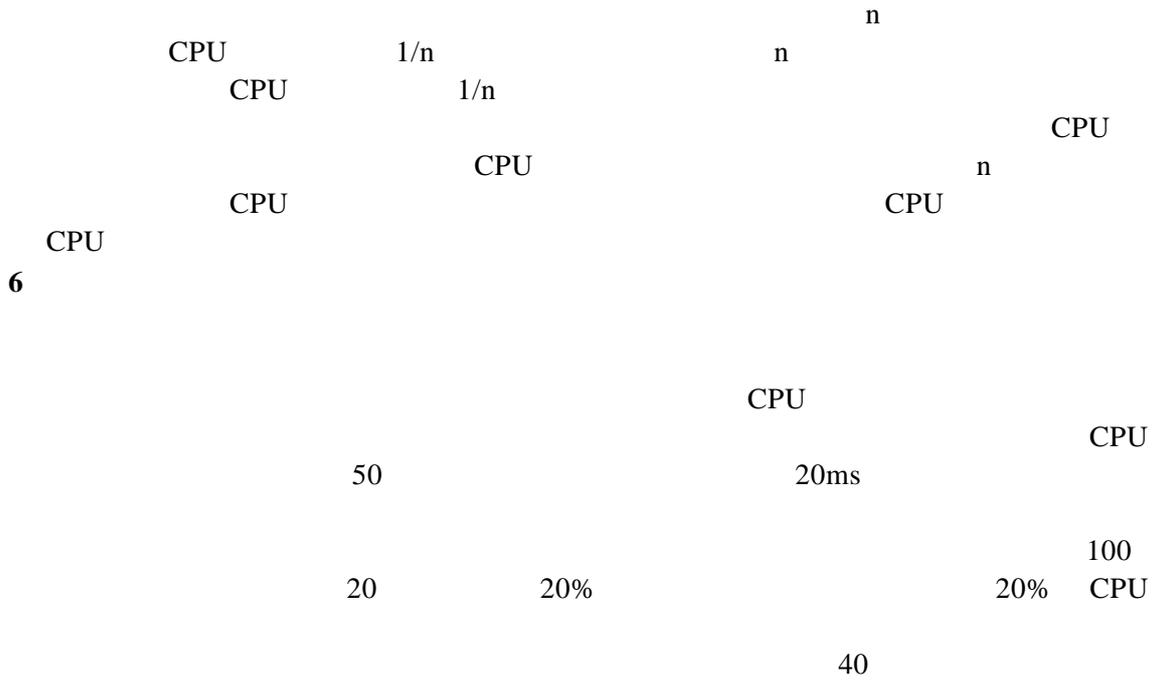
2-40 ()

XDS940
 I/O ()

I/O

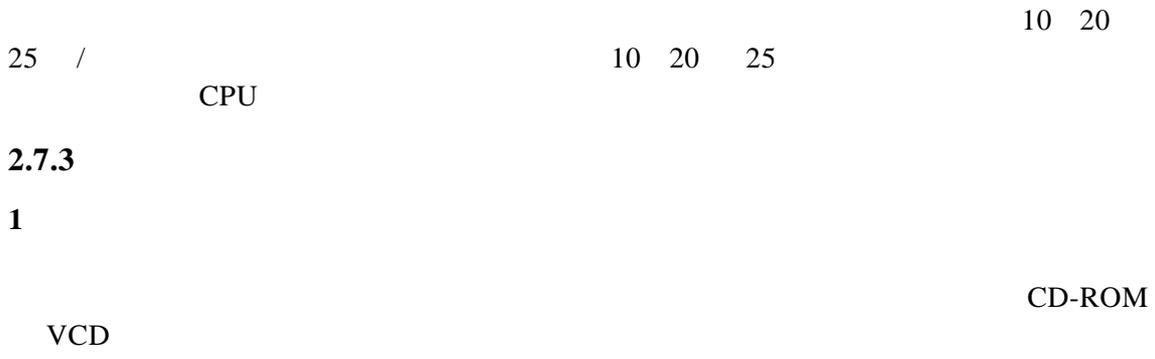


5



CPU

6



25 /

2.7.3

1

VCD

CD-ROM

hard real time

soft real time

m

i

P_i

C_i

CPU

$$C_1/P_1 + C_2/P_2 + \dots + C_m/P_m \leq 1$$

schedulable

500ms

50ms

30ms

100ms

100ms

200ms

$$0.5 + 0.15 + 0.2 = 1$$

150ms

2

1

10

20ms

50

100ms

2

3

$$= -(\text{laxity} + \dots)$$

2.7.4



cluster

I/O



1

master/slave

2 1 peer-to-peer

3

1 load sharing

-
-
-

Leutenegger S Vernon M 1

2

3

-

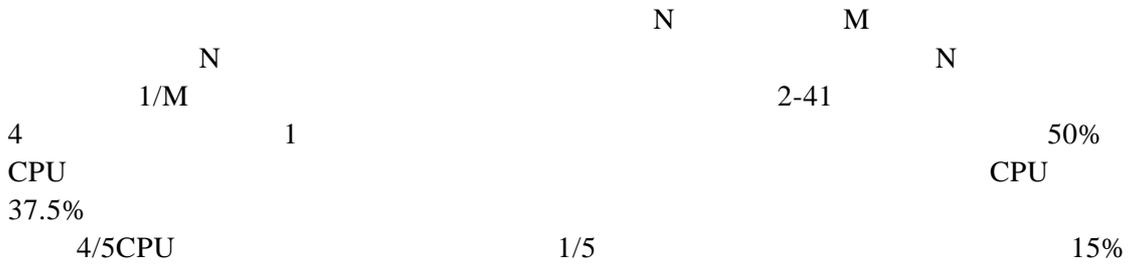
-
-

Mach

2

gang scheduling

-
-



1	2

1	2

37.5% 2-41

15%

3

dedicated processor assignment

dynamic scheduling

-
-
-

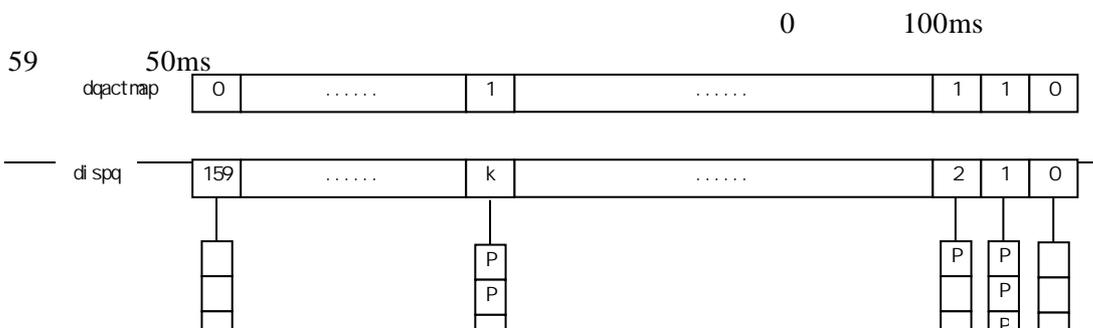
2.7.5 UNIX SVR4

UNIX SVR4 UNIX () UNIX
 SVR4 160 UNIX 3
 safe place
 safe place

UNIX SVR4

- 159-100
- 99-60
- 59-0

UNIX SVR4 2-42
 dispq
 dqactmap (1)
 dqactmap
 kprunrun ()



2.7.6 Windows 2000/XP

1 Windows 2000/XP

Windows 2000/XP
2000/XP
server

Windows

Window4 2000/XP

Windows 2000/XP

A 10

B 2

12

1/12

2 Windows2000/XP

Win32 API

Win32 API

Win32 API

API	
Suspend/ResumeThread	
Get/SetPriorityClass	
Get/SetThreadPriority	
Get/SetProcessAffinityMask	
SetThreadAffinityMask	
Get/SetThreadPriorityBoost	
SetThreadIdealProcessor	
Get/SetProcessPriorityBoost	
SwitchToThread	
Sleep	0
SleepEx	APC I/O

3

Windows 2000/XP

Windows

2000/XP

2-43 Windows 2000/XP

32

0 31

● 31-16

● 15-1

15

● 0

Windows 2000/XP

Win32

Win32

Win32

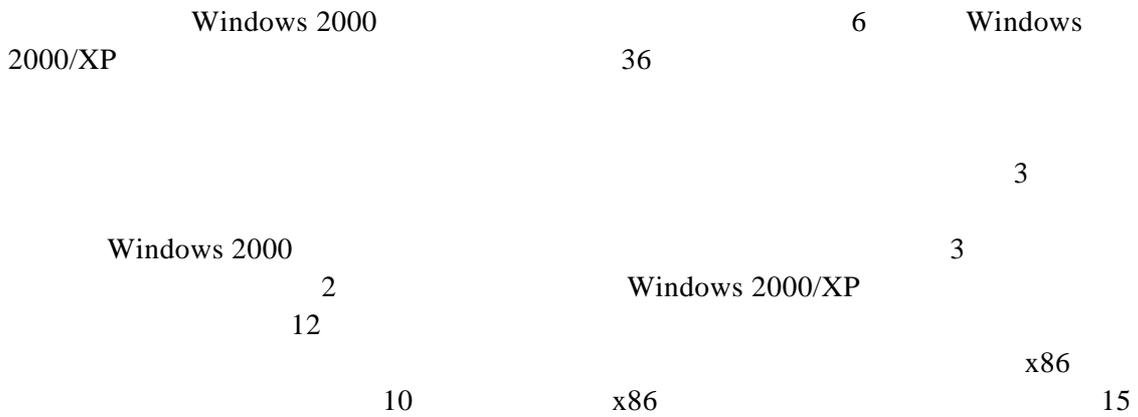
Task Manager

Win32

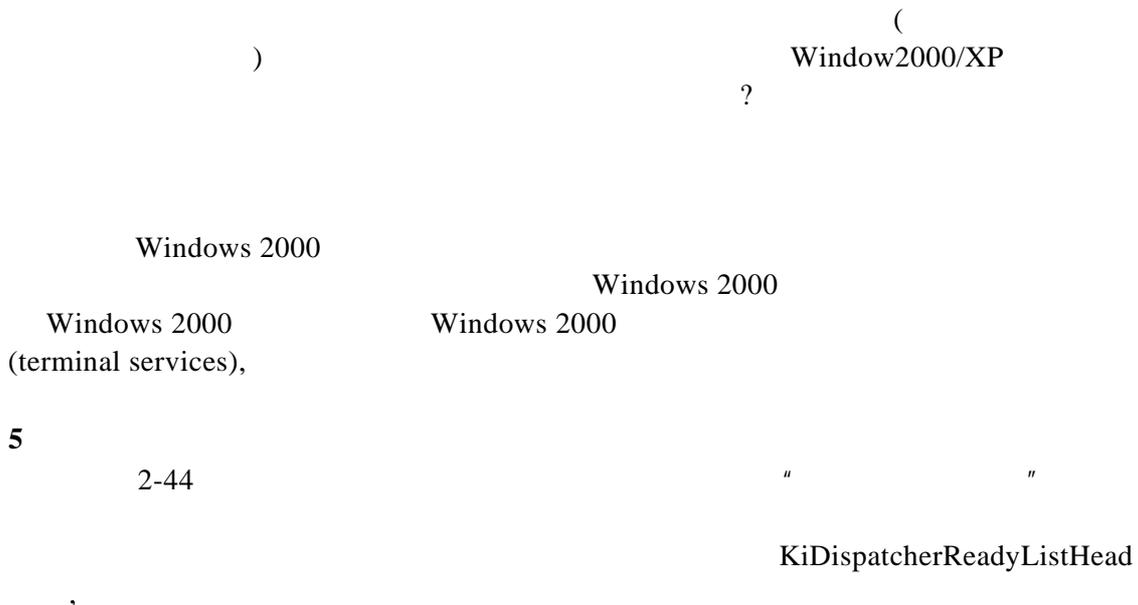
SetPriorityClass

quantum unit

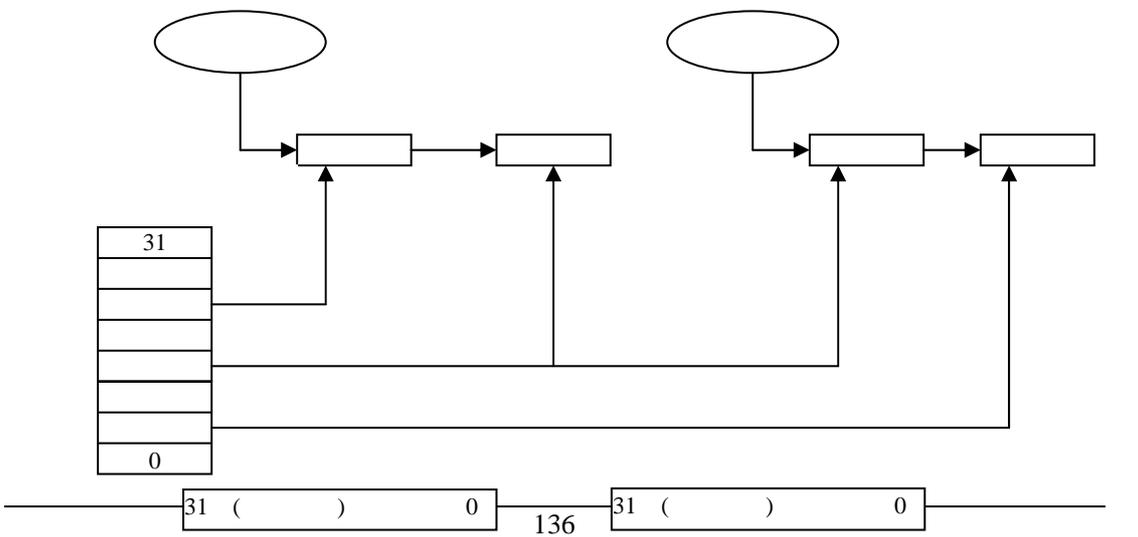
1



2)



5



	Windows 2000/XP				
KiReadySummary	32				
	0	0	1	1	
		KiIdleSummary	32		
		DPC/			
				KiDispatcherLock	,

6

Windows 2000/XP

Windows 2000/XP

1.4.5-1.4.5

4)

ExitThread

TerminateThread

7

5

Windows 2000/XP

1 I/O

2

3

4

5

16 31

I/O

I/O

I/O

I/O

IoCompleteRequest

I/O

1

2

6

8

I/O

CPU

KiUnwaitThread

PsPrioritySeparation

2

KeSetEvent

Windows2000/XP

(balance set

manager)

300

(300

3 4

)

15

2

300

16

10

10

CPU

8

Windows 2000/XP

1)

SetProcessAffinityMask SetThreadAffinityMask

2)



SetThreadIdealProcessor

3)

Windows 2000/XP

2000/XP

Windows

Windows 2000/XP

4)

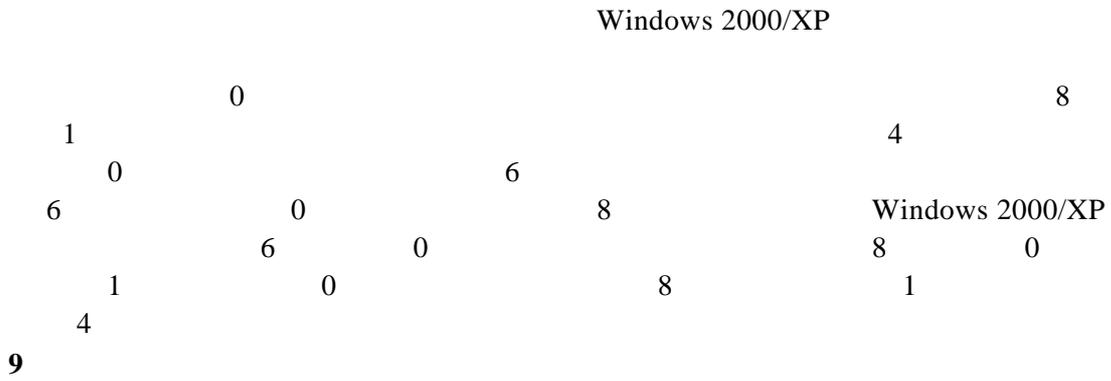
Windows 2000/XP

-
-
-
-

24 2

5)

Windows 2000/XP



0 Windows 2000/XP

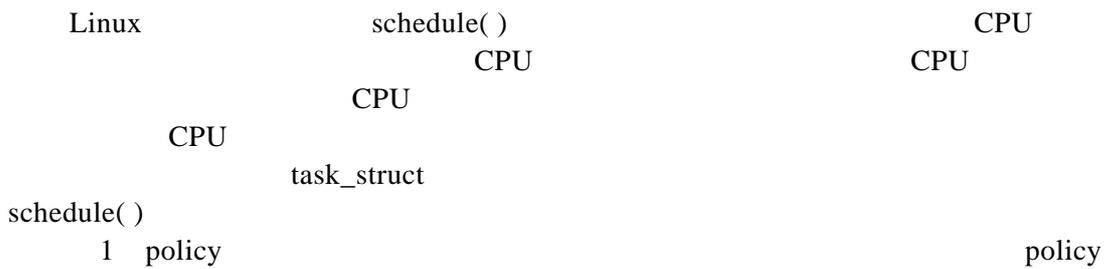
-
-
-
-

DPC

DPC

2.7.7

Linux



```

● #define SCHED_OTHER 0

● #define SCHED_FIFO 1          POSIX 1b FIFO

          rt_priority          CPU
● #define SCHED_RR 2          POSIX 1b RR

          SCHED_FIFO
          rt_priority
2 priority          -20 20          1 40
          CPU

3 rt_priority          CPU          0 99
          rt_priority+1000
          (
999          0 56)
4 counter          CPU          counter          0

          goodness()

Linux

          sched_setscheduler()
          sys_setpriority()

Linux          CPU          0
          Linux

task_struct          counter          CPU

          counter          CPU
10ms          60

          counter          priority
600ms          priority          rt_priority          priority
          policy

          conuter          priority
          counter          1

● counter          0          0
          0
          0          0 (
          0
          )          conuter          0

```

- 0
- 0 (priority 0) (priority+counter/2)
- I/O CPU
- 0 Linux Linux CPU
- counter 0
- Schedule()
- 1 Schedule() bottom half
- TASK_RUNNING
- TASK_INTERRUPTIBLE
- counter=0
- TASK_RUNNING
- TASK_RUNNING
- TASK_UNINTERRUPTIBLE
- 2 CPU
- goodness
- counter 0 999 0 56
- counter+1000 1001 1099
- counter
- counter
- 3 CPU task_struct
- 4 SMP Linux2.0.0 SMP CPU SMP

)

CPU

(ULT KLT)

/

Solaris Linux

Windows2000/XP

FCFS SJF SRTF

HRRF

FCFS

() ()

1. PSW? ?
2. ?
3. CPU
4. CPU CPU
5. ?
6. ? ?
7. ?
8. ?
9. ?
- 10.
11. ? ?
- 12.
13. ?
- 14.
- 15.
16. ?

- 17.
- 18.
19. Windows 2000/XP
20. Windows 2000/XP ?
21. Linux
22. Linux
23. Linux
24. ? ?
25. ?
26. ? ?
27. ?
28. ?
29. 1) 2)
30. ?
- 31.
- 32.
33. PCB ? ?
34. ?
- 35.
36. ?
37. ?
38. ?
39. ?
- 40.
- 41.
- 42.
43. UNIX SVR4
44. UNIX PCB proc user ?
- 45.
46. ?
- 47.
- 48.
49. (TCB)? ?
- 50.
51. ULT KLT
- 52.
- a)
- b)
- c)
- 53.
54. ? ?

55.

56. ?

57. Solaris

58. Solaris ? ?

59. Solaris

60. Solaris

61. Solaris

62. Solaris

63. Windows 2000/XP

64. Windows 2000/XP

65. Windows 2000/XP

66. Windows 2000/XP

67. Windows 2000/XP

68.

69.

70.

71.

72.

73.

74. 1 2 (3) (4)

75.

76. JCB

77.

78.

79.

80.

81. ?

82.

83.

84.

85.

86. " "

87.

88. CPU

89. UNIX

90. UNIX SVR4

91. Windows 2000/XP

92.

93. ?

94.

95.

96. (scheduling mechanism)

(scheduling policy)

97. ()

98.

0 /

99.

" 0"

100.

101.

1 ?

(1) 2 3 4 PSW 5

(6) (7) I/O

2 " " " "

" I/O "

3 ?

(1) (2) (3) (4)

4 ?

5 J1 J2 J3

a b c a<b<c

6 J₁ ... J_n S₁ ... S_n

CPU

7 0 1 2 3 4

5

1

()

2



1	10	3
2	1	1
3	2	3
4	1	4
5	5	2

8 S I/O T

CPU Q

1 Q 2 Q T 3 S Q T 4 Q S 5 Q 0

9 5 9 6 3 5 x

10 5 A E 2 4 6 8 10

1 2 3 4 5 5 1

2 3 4 (

C D B E A)

1 2 2 4

11 5 A E 10 6 2 4 8

3 5 2 1 4 5

1 FCFS(A B C D E) (2)

(3)

12 ? ?

13 CPU ?

14

15

1	10 00	2 00				
2	10 10	1 00				
3	10 25	0 25				
W=						

16 FCFS SJF HRRF

(1)
 (2)
 20 100K 2 1
 I/O

1	8:00	25	15K	1	1
2	8:20	10	30K	0	1
3	8:20	20	60K	1	0
4	8:30	20	20K	1	0
5	8:35	15	10K	1	1

FCFS
 CPU (1) ?(2)
 ?(3) ?(4) ?
 21 200K 5
 I/O

A	8:30	40	30K	3
B	8:50	25	120K	1
C	9:00	35	100K	2
D	9:05	20	20K	3
E	9:10	10	60K	1

(1)FIFO ?(4) SJF
 ?
 22 (1)FIFO
 ?(4) SJF ?
 23 (1)
 (2) 40 P2 25 P3 35 P1 20
 35 20 10 P4 60 35
 24 4 50 100 300 250ms
 xms x

CH3

3.1

3.1.1

()

Sequential

Programming

-
-
-
-

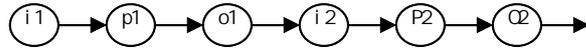
()

3.1.2

Concurrency

	A	B				a1	a2	a3	b1	b2	b3					
A	B								a1	a2	a3	b1	b2	b3		
			a1	b1	a2	b2	a3	b3	a1	b1	a2	b2	b3	a3	A	B

3-1 ?
while(TRUE){input process output}



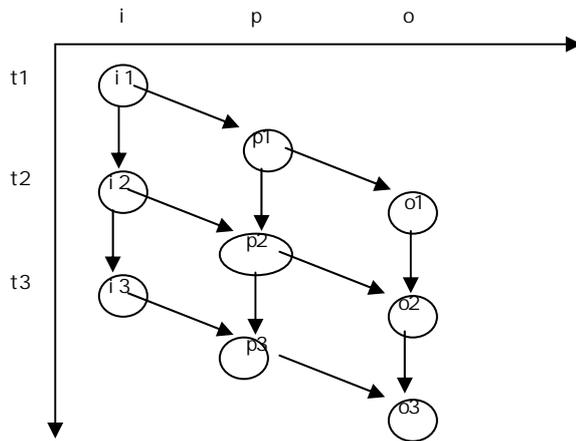
3-1

```

while(TRUE) {input send}
while(TRUE) {receive process send}
while(TRUE) {receive output}
( ) send receive

```

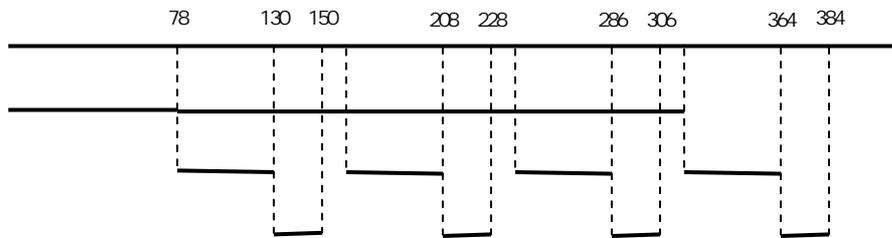
1 1
2 1 2
3 2



3-2

o1 3-2 t4 t5 n 3-3 t3 i3 p2

$$52 * n / 78 * n + 52 + 20$$



3-3

programming

concurrent

1966

Bernstein

Bernstein

$$R(p_i) = \{a_1 \ a_2 \ \dots \ a_n\} \quad p_i$$

$$W(p_i) = \{b_1 \ b_2 \ \dots \ b_m\} \quad p_i$$

Bernstein

$$R(p_1) \ \&W(p_2) \ \cup R(p_2) \ \&W(p_1) \ \cup W(p_1) \ \&W(p_2) = \{ \}$$

- S₁: a := x + y
- S₂: b := z + 1
- S₃: c := a - b
- S₄: w := c + 1

$$R(S_1) = \{x \ y\} \quad R(S_2) = \{z\} \quad R(S_3) = \{a \ b\} \quad R(S_4) = \{c\} \quad W(S_1) = \{a\} \quad W(S_2) = \{b\}$$

$$W(S_3) = \{c\} \quad W(S_4) = \{w\} \quad S_1 \quad S_2 \quad \text{Bernstein}$$

-
-
-
-

()

$$I/O \quad I/O \quad I/O \quad (1) \quad (2) \quad (3)$$

$$() \quad I/O$$

3.1.3

```

1
T1 T2
Aj(j=1 2 ...)
x1 x2
T1 T2
process Ti (i=1, 2)
var Xi : integer;
begin
  { Aj };
  Xi := Aj;
  if Xi >= 1 then begin Xi := Xi - 1; Aj := Xi; { }; end
  else { " " };
end;
T1 T2

T1: X1 := Aj;
T2: X2 := Aj;
X1 = m (m>0)
X2 (i = i)

```

```

coend
    borrow  return
B>x      {
{          }
          borrow
          x
          B  x
          return
          return

```

3.1.4 Interaction Among Processes

deadlock

starvation

()

FCFS

" "

mutual exclusion

()

input process output

Synchronization ()

3.2

3.2.1

```

1
Aj " " critical
section " " critical resource
T1
  X1 := Aj ;
  if X1 >= 1 then begin X1 := X1 - 1; Aj := X1;
T2
  X2 := Aj ;
  if X2 >= 1 then begin X2 := X2 - 1; Aj := X2;

```

Dijkstra 1965

shared

shared variable
region variable do statement

```

shared Aj
process Ti ( i = 1, 2 )
var Xi : integer;
begin
  Aj ;
  region Aj do begin
    Xi := Aj ;
    if Xi >= 1 then begin Xi := Xi - 1; Aj := Xi ; { } ; end
    else { } ;
  end;
end;

```

end;

-
-
-
-
-
-

region x do begin ... region y do ... end;

region y do begin ... region x do ... end;

3.2.2

	P1	P2	insidel	inside2
	true			false
P1(P2)		inside2(insidel)		
insidel(inside2)	true	P2(P1)	P1(P2)	
	false			
i n s i d e 1, i n s i d e 2 : B o o l e a n ;				
i n s i d e 1 := f a l s e; /* P1 */				
i n s i d e 2 := f a l s e; /* P2 */				
c o b e g i n				
p r o c e s s P 1				
b e g i n				
w h i l e i n s i d e 2 d o b e g i n e n d ;				
i n s i d e 1 := t r u e ;				
;				
i n s i d e 1 := f a l s e ;				
e n d ;				
p r o c e s s P 2				
b e g i n				
w h i l e i n s i d e 1 d o b e g i n e n d ;				
i n s i d e 2 = t r u e ;				
;				
i n s i d e 2 := f a l s e ;				
e n d ;				
c o e n d .				
insidel(inside2)	P2(P1)	P1(P2)	inside2(insidel)	insidel(inside2) false

```

inside2(inside1)  true          P1(P2)
inside1(inside2)  true          P2(P1)
                                     true

```

```

i n s i d e 1, i n s i d e 2: b o o l e a n;
i n s i d e 1 := f a l s e; /* P1 */
i n s i d e 2 := f a l s e; /* P2 */

c o b e g i n
p r o c e s s P1
b e g i n
    i n s i d e 1 := t r u e;
    w h i l e i n s i d e 2 d o b e g i n e n d;
    ;
    i n s i d e 1 := f a l s e;
e n d;
p r o c e s s P2
b e g i n
    i n s i d e 2 := t r u e;
    w h i l e i n s i d e 1 d o b e g i n e n d;
    ;
    i n s i d e 2 := f a l s e;
e n d;
c o e n d.

```

3.2.3

1 Dekker

```

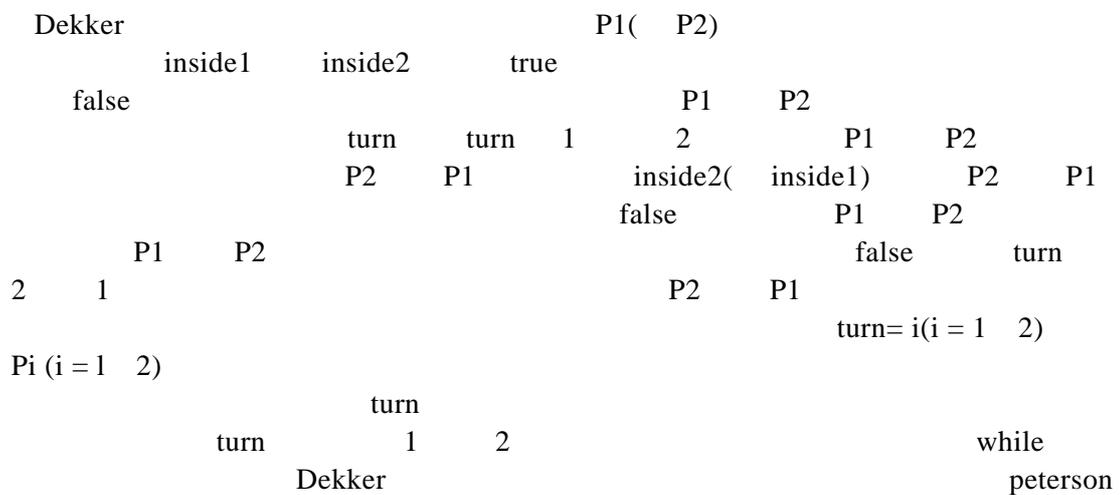
T.Dekker
turn
P1          turn=2          P2          Dekker          turn=1

```

```

var inside : array[1..2] of Boolean;
  Turn : integer;
  turn := 1 or 2;
  inside[1] := false;
  inside[2] := false;
cobegin
process P1
begin
  inside[1] := true;
  while inside[2] do
    if turn=2 then begin
      inside[1] := false;
      while turn=2 do begin end;
      inside[1] := true;
    end
  ;
  turn = 2;
  inside[1] := false;
end;
process P2
begin
  inside[2] := true;
  while inside[1] do
    if turn=1 then begin
      inside[2] := false;
      while turn=1 do begin end;
      inside[2] := true;
    end
  ;
  turn = 1;
  inside[2] := false;
end;
coend.

```



2 Peterson

1981 G.L.Perterson

false

turn=i

Pi turn
 Perterson

```
var inside: array[1..2] of boolean;
    turn: integer;
    turn := 1 or 2;
    inside[1] := false; /* P1 */
    inside[2] := false; /* P2 */
cobegin
process P1
begin
    inside[1] := true;
    turn := 2;
    while (inside[2] and turn=2)
        do begin end;
        ;
    inside[1] := false;
end;
process P2
begin
    inside[2] := true;
    turn := 1;
    while (inside[1] and turn=1)
        do begin end;
        ;
    inside[2] := false;
end;
coend.
```

turn while

```
                          inside[i]
                          Peterson                      while
" inside[i]    turn=1( 2)"
                          n
```

3.2.4

1

CPU

2

```

Set x true x false x x TS(Test and TS(x) TS
TS(x): x=true x:=false return true return false
TS s s
s true TS s true s false
s true TS
x
s : boolean;
s := true;
process Pi /* i = 1, 2, ..., n*/
pi : boolean;
begin
repeat pi := TS(s) until pi; /* */
;
s := true; /* */
end;

```

3

Swap

```

Swap (a, b): temp:=a; a:=b; b:=temp;
Intel 80x86 XCHG lock false

```

```

lock : boolean;
lock := false;
process Pi /* i = 1, 2, ..., n*/
pi : boolean;
begin
pi := true;

```

```

repeat Swap(lock, pi) until pi = false; /* */
;
lock := false; /* */
end;

```

3.3 PV

3.3.1

producer-consumer problem

```

-- n m pi cj k
cj pi
-
var k:integer;
type item:any;
buffer:array[0..k-1] of item;
in,out:integer:=0;
counter:integer:=0;
process producer
while (TRUE) /* */
produce an item in nextp; /* */
if (counter==k) sleep ( ); /* */
buffer[in]:=nextp; /* */
in:=(in+1) mod k; /* */
counter:=counter+1; /* 1*/
if (counter==1) wakeup( consumer); /* */
process consumer
while (TRUE) /* */
if (counter==0) sleep ( ); /* */
nextc:=buffer[out]; /* nextc*/
out:=(out+1) mod k; /* */
counter:=counter-1; /* 1*/
if (counter==k-1) wakeup( producer); /* */
consume thr item in nextc; /* */
sleep() wakeup()
counter counter counter counter

```

8
 counter 1
 counter 9 7 8 1
 counter 0
 counter 1 counter 1 wakeup counter
 0 counter
 0
 pi cj

PV
3.3.2 PV

busy waiting CPU
 1965 E.W.Dijkstra -- P V
 semaphore P
 V

- 0 1
-

1

```

s
V
• P(s) s 0
s 0
• V(s) s 1
P(s) V(s)
P(s) while s > 0 do null operation
      s:=s-1
V(s) s:=s+1
P s<=0
" s "
```

2

```

s value
queue value
P V
• P(s) s 1 0 P(s)
s
• V(s) s 1 0 s
P V
```

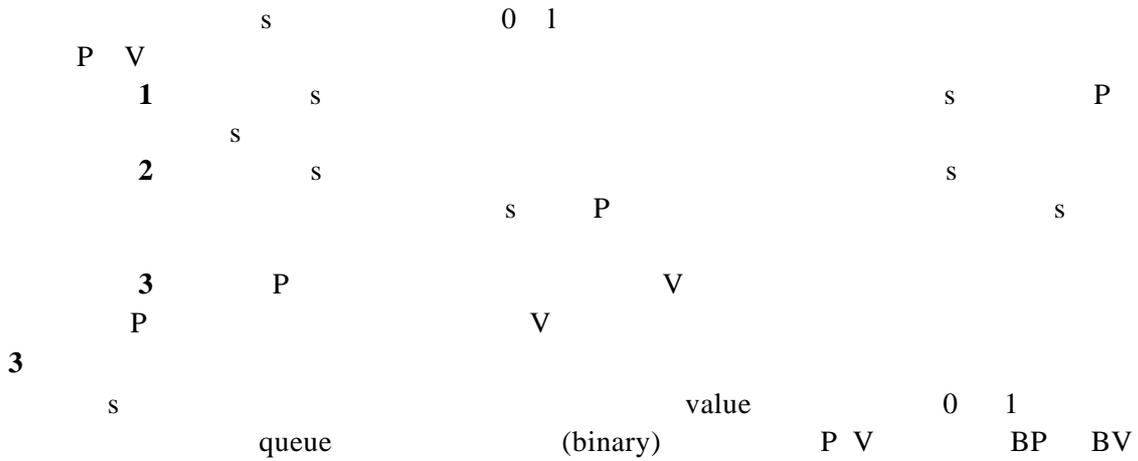
```

type semaphore=record
  value: integer;
  queue: list of process;
end
procedure P(var s: semaphore);
begin
  s.value := s.value - 1; /* 1 */
  if s.value < 0 then W(s.queue); /* 0 P(s)
                                W(s.queue)
                                s queue*/
end;

procedure V(var s: semaphore);
begin
  s.value := s.value + 1; /* 1 */
  if s.value > 0 then R(s.queue); /* 0 R(s.queue)
                                s queue
                                */
end;

W(s.queue) s s
CPU R(s.queue) s s
```

?



```

type binary semaphore=record
    value(0,1);
    queue: list of process
end;
procedure BP(var s:semaphore);
    if s.value=1;
    then
        s.value=0;
    else begin
        w(s.queue);
    end;
procedure BV(var s:semaphore);
    if s.queue is empty;
    then
        s.value:=1;
    else begin
        R(s.queue);
    end;
    0 1

```

```

var
    s1: binary-semaphore;
    s2: binary-semaphore;
    c:integer;
    s1=1 s2=0 c
s P V s
    P(s) V(s)
    begin begin
        BP(s1); BP(s1);
        c:=c-1; c:=c+1;
        if c<0 then BP(s2); if c 0 then BV(s2);
        BV(s1); BV(s1);
    end end
    s c P V s
s1 c 1 c s2 s

```

3.3.3

P V
TS

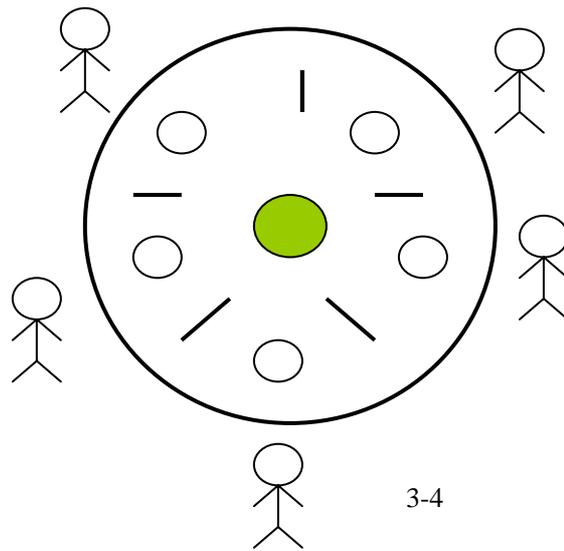
TS P V P V
P V

```

P      P      mutex      l      T1      T1
P      P      mutex      l      T1      T2
P      V      V      P      T2      mutex      T1
V      T1      T2      mutex      T2      T2
A[j]
s[j]
P V
else V
?      P V
      P V
P V

```

Dijkstra 1965



S_i ($i=0,1,2,3,4$) 1
P

V

```

var forki : array[0..4] of semaphore;
  forki := 1;
cobegin
  process Pi // i=0, 1, 2, 3, 4
  begin
    L1:
      ;
    P(fork[i]);

```

```

    P(fork[i +1] mod 5);
    ;
    V(fork[i]);
    V(fork[i +1] mod 5);
    goto L1;
end;
coend.

```

P(fork[4])

P(fork[0])

-
-
-

3.3.4

P V

```

                P V
            empty full          1 0 empty
full

```

```

var B : integer;
    empty: semaphore; /* */
    full : semaphore; /* */
    empty := 1; /* */
    full := 0; /* */
cobegin
process producer
begin
    L1:
    Produce a product;
    P(empty);
    B := product;
    V(full);
    Goto L1;
end;
process consumer
begin
    L2:
    P(full);
    Product := B;
    V(empty);
    Consume a product;

```

```

        Goto L2;
    end;
coend.

        P V
m        n        k
        empty    mutex    full    1    0
                k

var B : array[0..k-1] of item
empty: semaphore := k;    /*
full : semaphore := 0;    /*
mutex: semaphore := 1;    /*
in : integer := 0; x M

```

V

P
P

V

/
apple

orange

```

var
  plate : integer;
  sp: semaphore; /* */
  sg1: semaphore; /* */
  sg2: semaphore; /* */
  sp := 1; /* */
  sg1 := 0; /* */
  sg2 := 0; /* */
cobegin
  process father
  begin
    L1:      ;
    P(sp);
      plate;
    V(sg2);
    goto L1;
  end;
  process mother
  begin
    L2:      ;
    P(sp);
      plate;
    V(sg1);
    goto L2;
  end;
end;
process son
begin
  L3: P(sg1);
    plate ;
    V(sp);
      ;
    goto L3;
end;
process daughter
begin
  L4: P(sg2);
    plate ;
    V(sp);
      ;
    goto L4;
end;
coend.

```

3.3.5

reader-writer problem

Courtois 1971

2

4

F

1

3

mutex

rc

W

rc

```

var rc: integer;

```

```

W mutex: semaphore;
rc := 0; /* */
W := 1;
mutex := 1;
procedure read;
begin
  P(mutex);
  rc := rc + 1;
  if rc=1 then P(W)
  V(mutex);

  P(mutex);
  rc := rc - 1;
  if rc = 0 then V(W)
  V(mutex);
end;
procedure write;
begin
  P(W);
  ;
  V(W);
end;

cobegin
  process readeri;
  process writerj;
coend.
process readeri;
begin
  read;
end.
process writerj
begin
  write;
end.

```

/ /

Solaris reader/writer rw-enter()() rw-exit()()
 rw-tryenter()() rw-downgrade()(write lock read lock)
 rw-trygrade()(read lock write lock)

3.3.6

n

```

        3
        0      customers
        0      barbers
        0      mutex
                                1
                                waiting

var v a i t i n g : i n t e g e r ; /* */
    C H A I R S : i n t e g e r ; /* */
    c u s t o m e r s , b a r b e r s m u t e x : s e m a p h o r e ;
    c u s t o m e r s := 0 ; b a r b e r s := 0 ;
    v a i t i n g := 0 ; m u t e x := 1 ;
procedure barber ;
begin
while (TRUE) ; /* */
    P ( c u s t o m e r s ) ; /* */
    P ( m u t e x ) ; /* */
    v a i t i n g := v a i t i n g - 1 ; /* */
    V ( b a r b e r s ) ; /* */
    V ( m u t e x ) ; /* */
    c u t - h a i r ( ) ; /* */
end ;
procedure customer
begin
    P ( m u t e x ) ; /* */
    i f v a i t i n g < C H A I R S /* */
        begin
            v a i t i n g := v a i t i n g + 1 ; /* */
            V ( c u s t o m e r s ) ; /* */
            V ( m u t e x ) ; /* */
            P ( b a r b e r s ) ; /* */
            g e t - h a i r c u t ( ) ; /* */
        end ;
    e l s e V ( m u t e x ) ; /* */
end .
```

3.4

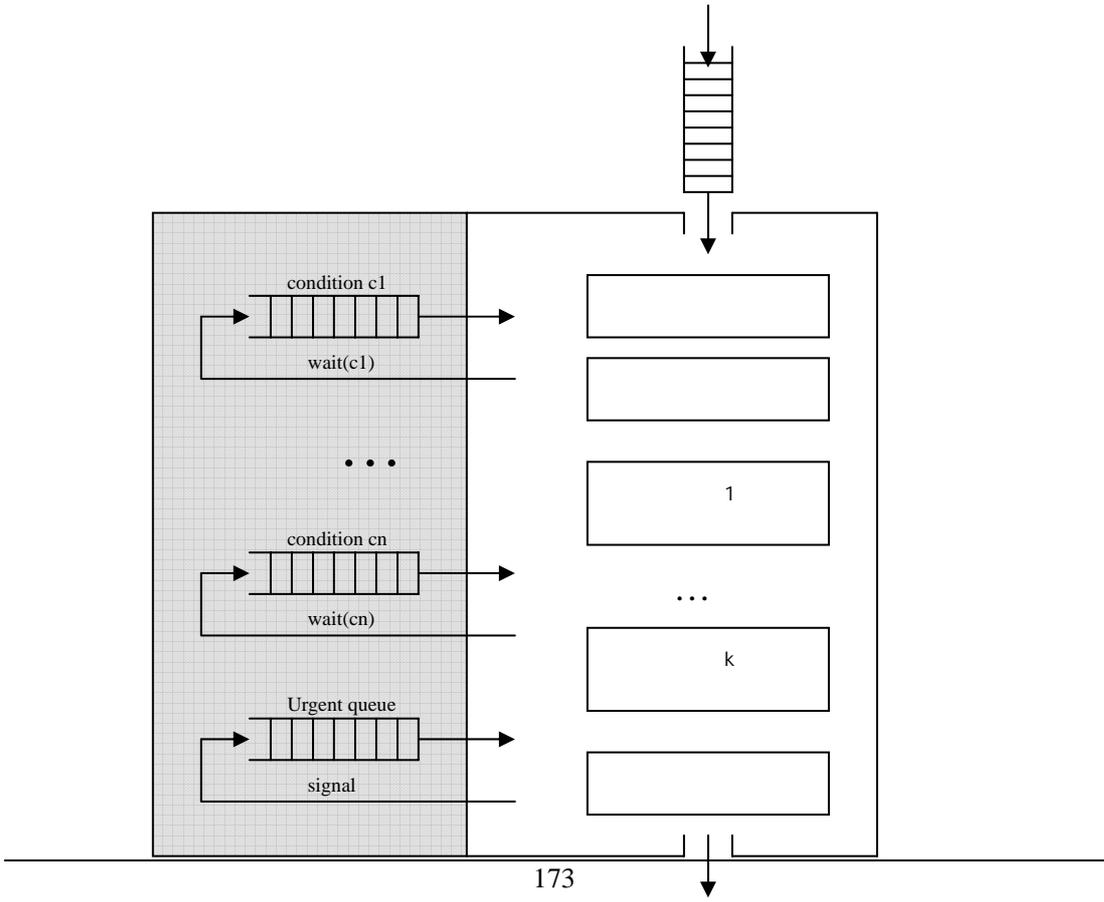
3.4.1

P V

1974	1975	Brinch Hansen	Hoare
—	monitor		
	”	”	”
			”

-
-
-

Pascal Modula-2 Java



```

TYPE < > = MONITOR
  < >;
  define < >;
  use < >;
  procedure < > < > ;
    begin
      < >;
    end;
  .....
  procedure < > < > ;
    begin
      < >
    end;
  begin
    < >
  end.

```

define

use

3-5

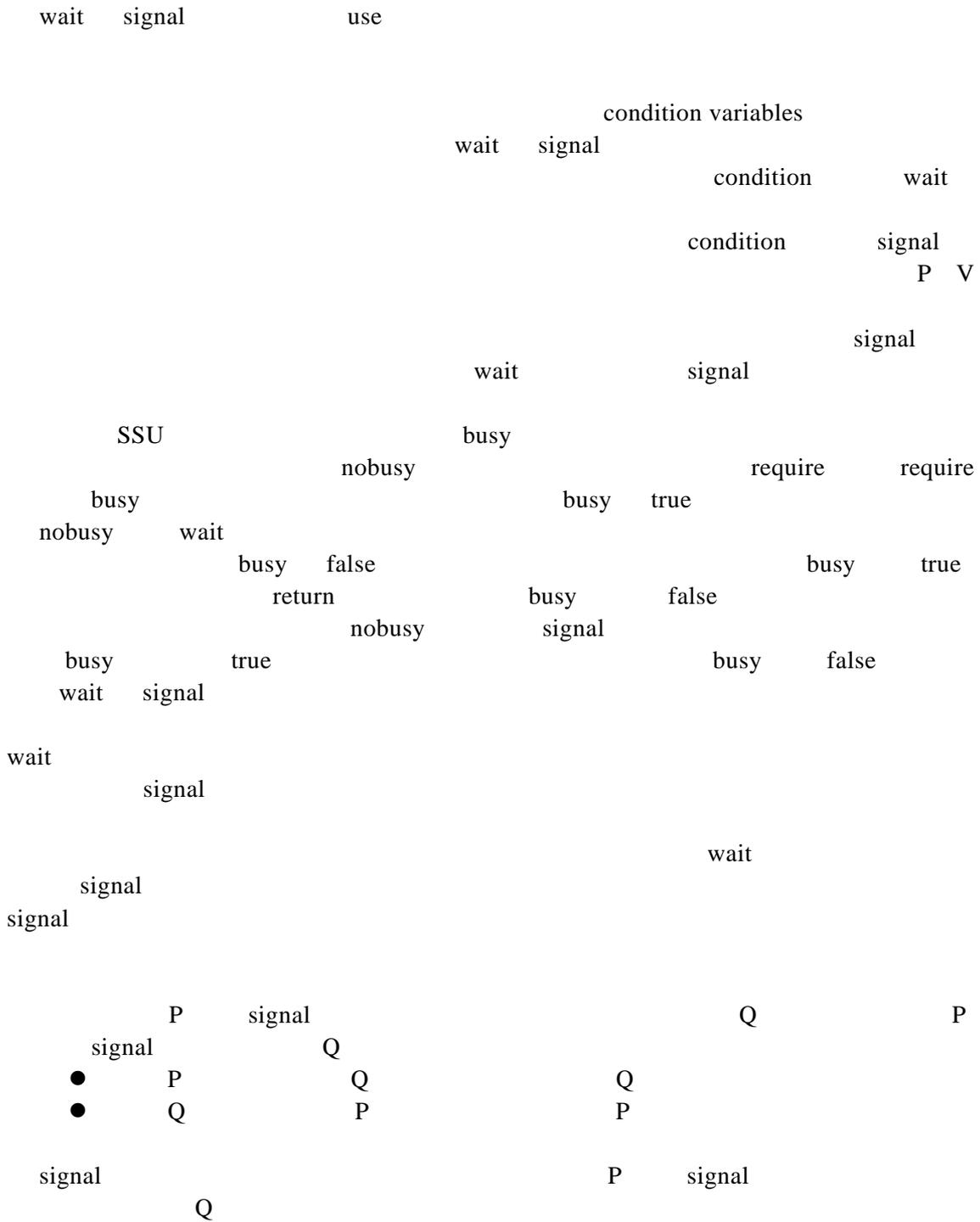
```

TYPE SSU = MONITOR
  var busy : boolean;
      nobusy : condition;
  define require, return;
  use wait, signal
  procedure require;
    begin
      if busy then wait(nobusy); /* */
      busy := true
    end;
  procedure return;
    begin
      busy := false
      signal(nobusy); /* */
    end;
  begin /* */
    busy := false
  end;

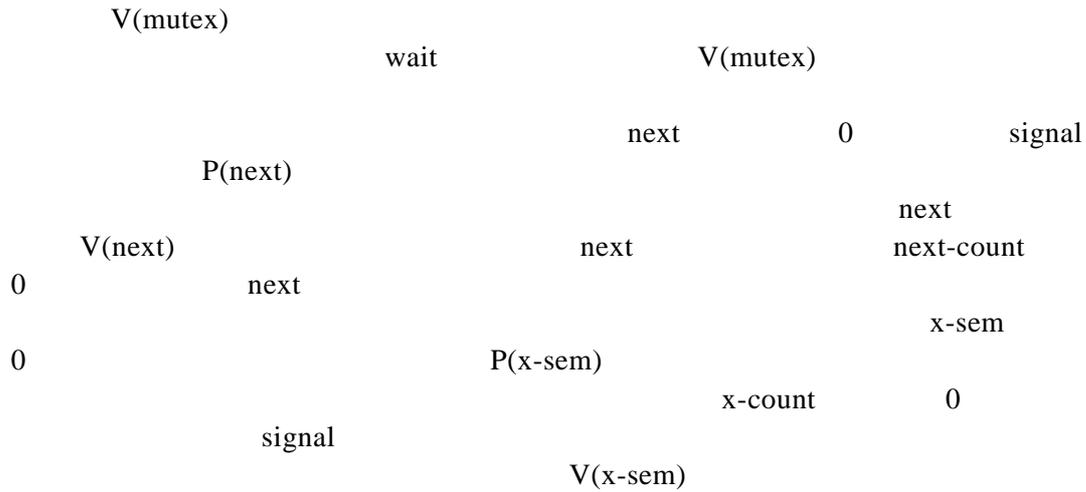
```

SSU

define



3.4.2 Hoare



```
TYPE interf = RECORD
```

```

  mutex: semaphore; /* */
  next: semaphore; /* signal */
  next_count: integer; /* next */

```

```
END;
```

```
wait signal
```

```
procedure wait (var x_sem semaphore, x_count: integer, IMinterf);
```

```
begin
```

```
  x_count := x_count + 1;
```

```
  if IMnext_count > 0 then V(IMnext); else V(IMmutex);
```

```
  P(x_sem);
```

```
  x_count := x_count - 1;
```

```
end;
```

```
procedure signal (var x_sem semaphore, x_count: integer, IMinterf);
```

```
begin
```

```
  if x_count > 0 then begin
```

```
    IMnext_count := IMnext_count + 1;
```

```
    V(x_sem);
```

```
    P(IMnext);
```

```
    IMnext_count := IMnext_count - 1;
```

```
  end;
```

```
end;
```

```
P(IMmutex);
```

```
< >;
```

```
if IMnext_count > 0 then V(IMnext);
```

```
else V(IMmutex);
```

```

        var state array[0..4] of (thinking hungry eating)
    i      state[i]=eating
state[(i-1)mod 5] eating      state[(i+1)mod5] eating
        var self array[0..4] of semaphore
    i
self[ ]      Hoare      P V
        condition      semaphore

```

```

TYPE dining-philosophers = MONITOR
    var state : array[0..4] of (thinking, hungry, eating);
        self : array[0..4] of semaphore;
        s-count : array[0..4] of integer;
    define pickup, putdown;
    use wait, signal
    procedure test(k : 0..4);
    begin
        if state[(k-1) mod 5] <>eating and state[k]=hungry
            and state[(k+1) mod 5] <>eating then begin
                state[k] := eating;
                signal(self[k], s-count[k], 1);
            end;
        end;
    procedure pickup(i : 0..4);
    begin
        state[i] := hungry;
        test(i);
        if state[i] <>eating then wait(self[i], s-count[i], 1);
    end;
    procedure putdown(i : 0..4);
    begin
        state[i] := thinking;
        test((i-1) mod 5);
        test((i+1) mod 5);
    end;
    begin
        for i := 0 to 4 do state[i] := thinking;
    end;

```

pickup

putdown

cobegin

```

    process philosopher-i
    begin
        .....
        P(IMutex);
        call dining-philosopher.pickup(i);
        if IMnext-count > 0 then V(IMnext);
            else V(IMutex);

            .....
            P(IMutex);
            call dining-philosopher.putdown(i);
            if IMnext-count > 0 then V(IMnext);
                else V(IMutex);

            .....
        end;
    coend.

```

3.4.3 Hanson

```

                                signal
                                signal
                                wait signal check release

```

- wait
- signal
- check
- release

```

                                check release
                                check
                                release

```

```

TYPE interf = RECORD
    intsem: condition; /* */
    count1 : integer; /* */
    count2 : integer; /* */
END; /* */

intsem          count1          count2

```

```

procedure wait(var s: condition; var IMinterf);
begin
  s := s + 1;
  IMcount2 := IMcount2 - 1;
  if IMcount1 > 0 then
    begin
      IMcount1 := IMcount1 - 1;
      IMcount2 := IMcount2 + 1;
      R(IMntsen);
    end;
  W(s);
end;

procedure signal(var s: condition; var IMinterf);
begin
  if s > 0 then
    begin
      s := s - 1;
      IMcount2 := IMcount2 + 1;
      R(s);
    end;
end;

procedure check(var IMinterf);
begin
  if IMcount2 = 0
  then IMcount2 := IMcount2 + 1;
  else
    begin
      IMcount1 := IMcount1 + 1;
      W(IMntsen);
    end;
end;

procedure release(var IMinterf);
begin
  IMcount2 := IMcount2 - 1;
  if IMcount2 = 0 and IMcount1 > 0 then
    begin
      IMcount1 := IMcount1 - 1;
      IMcount2 := IMcount2 + 1;
      R(IMntsen);
    end;
end;

```

signal

wait

s

W(s) s R(s)
s W(IM intsem) R(IM intsem)

-
-
-

1
F 1 2
3
rc wc R W

```

type read-writer = MONITOR
  var rc, wc : integer;
      R, W: condition;
  define start-read, end-read, start-writer, end-writer;
  use wait, signal, check, release;
procedure start-read;
begin
  check(IM);
  if wc>0 then wait(R, IM);
  rc := rc + 1;
  signal(R, IM);
  release(IM);
end;
procedure end-read;
begin
  check(IM);
  rc := rc - 1;
  if rc=0 then signal(W, IM);
  release(IM);
end;
procedure start-writer;
begin
  check(IM);
  wc := wc + 1;
  if rc>0 or wc>1 then wait(W, IM);
  release(IM);
end;
procedure end-writer;
begin
  check(IM);

```

```

vc := vc - 1;
if vc > 0 then signal (WIM);
else signal (RIM);
release(IM);
end;
begin
rc := 0; vc := 0; R := 0; W := 0;
end.

```

```

                                start-read start-write
                                end-read end-write

```

```

cobegin
  process reader
begin
  .....
  call read-writer.start-read;
  .....
  read;
  .....
  call read-writer.end-read;
  .....
end;
  process writer
begin
  .....
  call read-writer.start-write;
  .....
  write;
  .....
  call read-writer.end-write;
  .....
end;
coend.

```

2

orange

apple

put get

```

type FMSD = MONITOR
var plate : (apple, orange);
full : boolean;
SP, SS, SD : condition;

```

```

define put, get;
use wait, signal, check, release;
procedure put(var fruit: (apple, orange));
begin
  check(IM);
  if full then wait(SP, IM);
  full := true;
  plate := fruit;
  if fruit=orange
    then signal(SS, IM);
  else signal(SD, IM);
  release(IM);
end;
procedure get(var fruit: (apple, orange), x: plate);
begin
  check(IM);
  if not full or plate<>fruit
  then begin
    if fruit = orange
      then wait(SS, IM);
    else wait(SD, IM);
  end;
  x := plate;
  full := false;
  signal(SP, IM);
  release(IM);
end;
begin
  full := false; SP := 0; SS := 0; SD := 0;
end;

```

put

get

```

cobegin
  process father
  begin
    .....
    ;
    call FMSD.put(apple);
    .....
  end;
  process mother
  begin
    .....
    ;
    call FMSD.put(orange);
    .....
  end;
end;

```

```

    process son
begin
    .....
    call FMSD.get(orange, x);
        ;
    .....
end;
    process daughter
begin
    .....
    call FMSD.get(apple, x);
        ;
    .....
end;
coend;

```

3 monitor

```

type producer-consumer = MONITOR
    var B: array[0..k-1] of item /* */
        in, out: integer; /* */
        count: integer; /* */
        notfull, notempty: condition; /* */
    define append, take;
    use wait, signal, check, release;
procedure append(x: item);
begin
    check(IM);
    if count=k then wait(notfull, IM); /* */
        B[in] := x;
        in := (in+1) mod k;
        count := count+1; /* */
        signal(notempty, IM); /* */
    release(IM);
end;
procedure take(x: item);
begin
    check(IM);
    if count=0 then wait(notempty, IM); /* */
        x := B[out];
        out := (out+1) mod k;
        count := count-1; /* */
        signal(notfull, IM); /* */
    release(IM);
end;

```

```

begin
    in := 0; out := 0; count := 0;
end;
cobegin
    /* */
process producer;
    var x: item;
begin
    produce(x);
    append(x);
end;
process consumer;
    var x: item;
begin
    take(x);
    consume(x);
end;
coend.

```

3.5

PV

IPC InterProcess Communication

- signal
 - PV
 - pipeline
 - /
- shared memory
- shared file
- message passing
- PV

3.5.1

signal

OS

signal del+ctrl+c

PCB

UNIX

- 32 SIGCLD SIGHUP SIGKILL SIGCHLD SIGSTOP
- SIGBUS SIGSEGV SIGPWR SIGFPE
- SIGPIPE SIGSYS SIGILL
- SIGINT SIGQUIT
- delete break SIGTERM SIGALRM SIGUSR1 SIGUSR2
- SIGTRAP

UNIX UNIX alarm

proc user p-clktim P-sigign 32

SIGALAM 1 P-sig 32

1 U-signal[NSIG] 32

pid sig sig kill(pid,sig)

kill signal(sig,func) sig

function

function=1 sig function=0

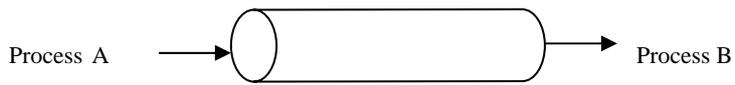
function (function

) signal

3.5.2

pipeline

3-6



3-6 pipe

pipe UNIX C UNIX
 () ()

•

i-node

i-node

•

SIGPIPE

•

i-node
 write read
 write

K

write
 write

•

UNIX

```
int pipe(files);
int files[2];
```

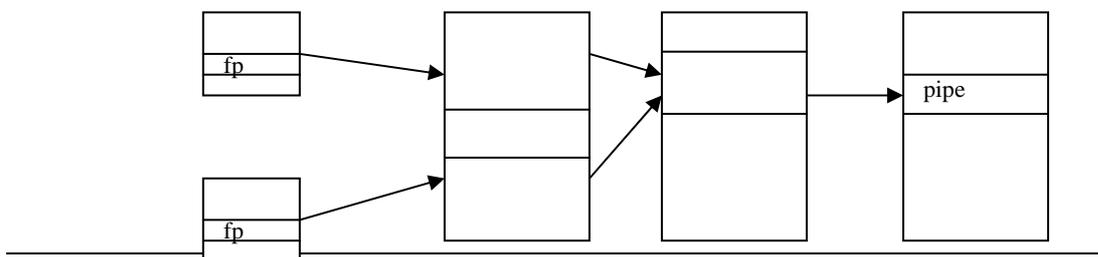
pipe

i-node

files[0] files[1]
 files[1]

pipe
 pipe

pipe()
 files[0]
 3-7



```

#include<stdio.h>
#define MSGSIZE 16
char *msg1="hello,world#1";
char *msg2="hello,world#2";
char *msg3="hello,world#3";
main( )
{
    char inbuf[MSGSIZE];
    int p[2],j,pid;
    /*open pipe*/
    if(pipe<0) {
        perror("pipe call");
        exit(1);
    }
    if((pid=fork( )<0) {
        perror("fork call");
        exit(2);
    }
    /*if parent, then close read file descriptor and write down pipe*/
    if(pid>0 {
        close(p[0]);
        write(p[1],msg1,MSGSIZE);
        write(p[1],msg2,MSGSIZE);
        write(p[1],msg3,MSGSIZE);
        wait((int*)0);
    }
    /*if child ,then close write file descriptor and read from pipe*/
    if (pid==0) {
        close(p[1]);
        for(j=0;j<3;j++)
            read(p[0],inbuf,MSGSIZE);
        printf("%s\n,inbuf");
    }
    }
    exit(0);
}

```

		UNIX
	FIFO	
UNIX	<pre> mknod(pipename,S_FIFO+rw,0) open(pipename,O_RDONLY) open(pipename,O_WRONLY) </pre>	<pre> write read pipe read write </pre>

who| sort | more

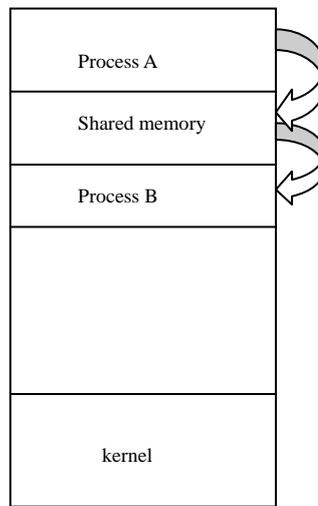
1000

10 100
10

AT&T UNIX

3.5.3

3-8



3-8

UNIX/Linux

- `shmget(key, size, permflags)`
 Parameters: `key`, `size`, `permflags`
- `shmat(shm-id, daddr, shmflags)`
 Parameters: `shm-id`, `daddr`, `shmflags`
- `Shmdt(memptr)`
 Parameter: `memptr`
- `Shmctl(shm-id, command, &shm-stat)`
 Parameters: `shm-id`, `command`, `&shm-stat`

3.5.4

1

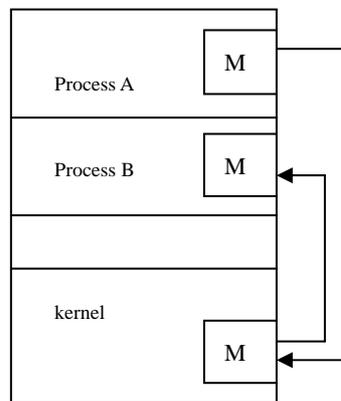
— — **message passing**

send receive

2

UNIX pipeline socket

3-9



3-9

1

send receive

- send P P
 - receive Q Q
- P Q

/
port

- send A
- receive A

" " " "

-
-

R() W()

```

type box=record
    size:integer;           /* */
    count:integer;         /* */
    letter:array[1..n] of message; /* */
    S1,S2:semaphore;      /* */
end

```

```

procedure send(varB:box,M:message)
    var i:integer;
    begin
        if B.count=B.size then W(B.s1);
        i:=B.count+1;
        B.letter[i]:=M;
        B.count:=i;
        R(B.S2)
    end;

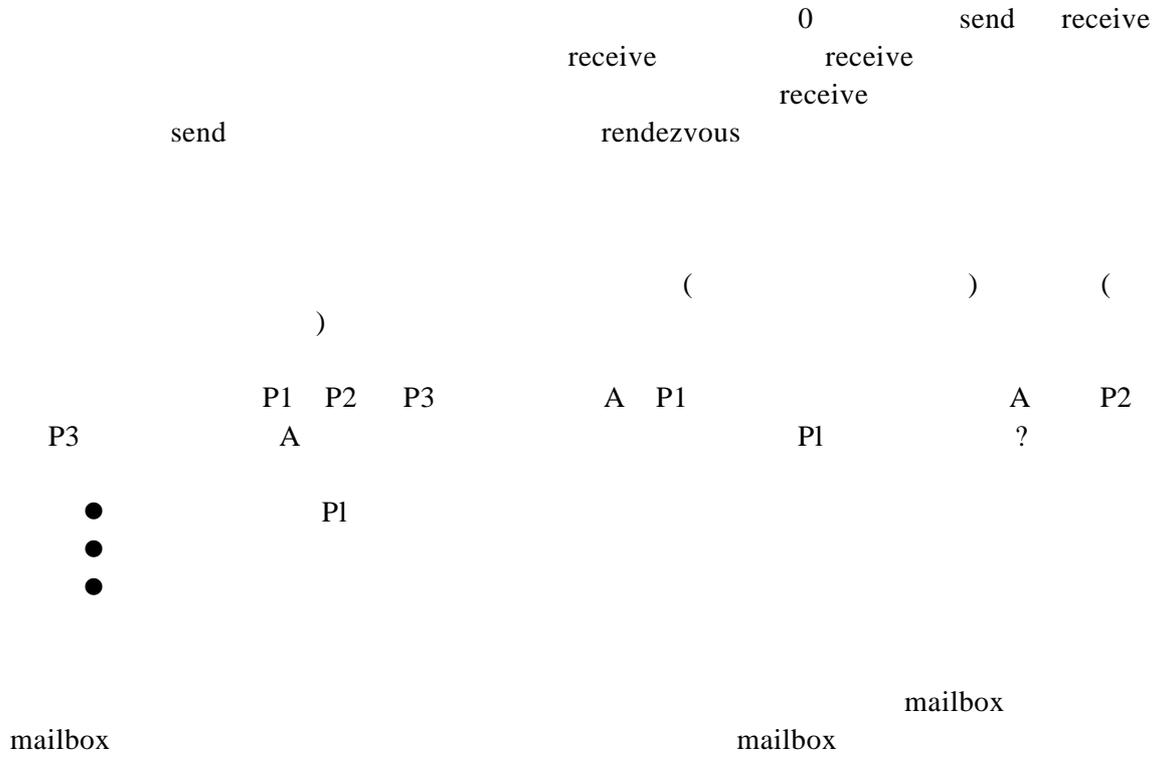
```

```

procedure receive(varB:box,x:message)
    var i:integer;
    begin
        if B.count=0 then W(B.s2);
        B.count:=B.count-1;
        x:=B.letter[1];
        if B.count > 0 then for i=1 to B.count do B.letter[i]:=B.letter[i+1];

```


3.5.5



1973

P.B.Hansan

RC4000

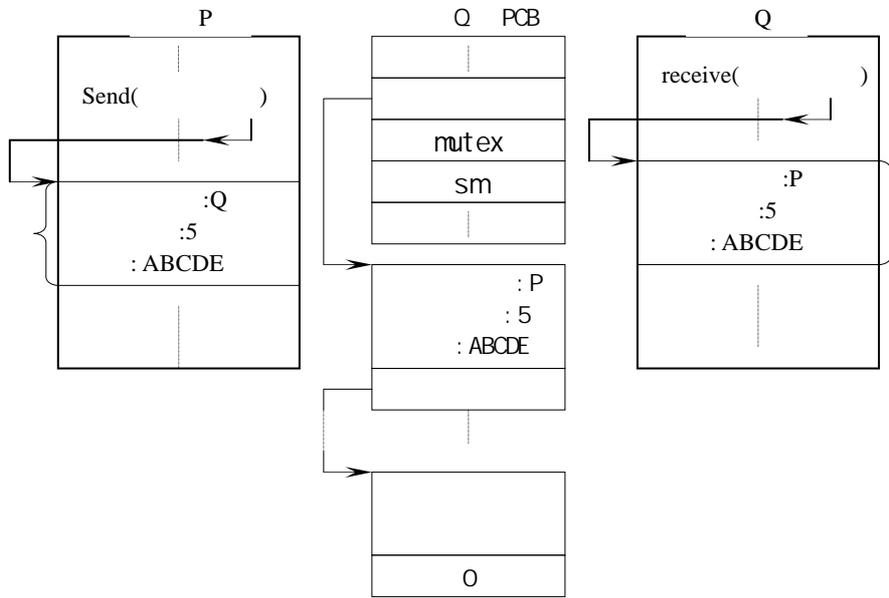
3-10

- sender
- size
- text
- next-ptr
- PCB
- mptr
- mutex
- sm

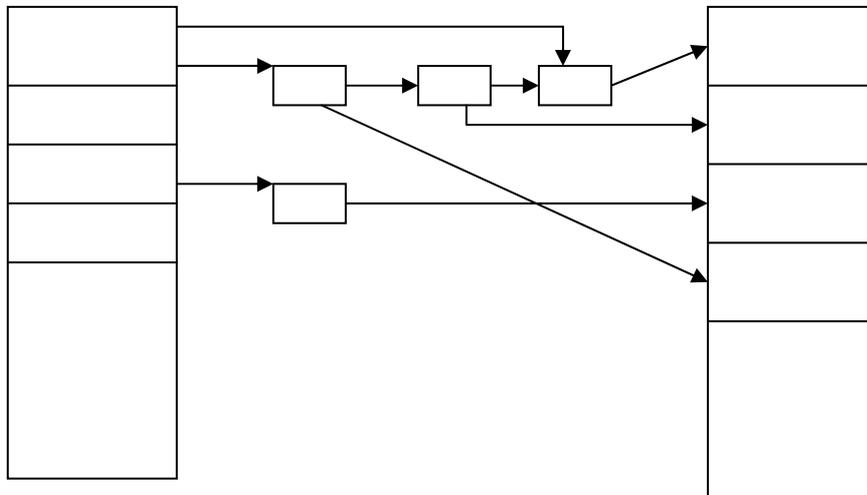
1

0

- send PCB P(mutex)
- V(sm) 1 V(mutex)
- receive P(sm) V(mutex) P(mutex)



3-10



3-11 Unix/Linux

- msgbuf
- 100
-

/

UNIX/Linux

- msgget
- msgsnd
- msgrcv
- msgctl

msgget

msgsnd

msgrcv
 =0 <0 >0 (=0
 >0 <0
)

msgctl

recei ve n 0 n

3.6

3.6.1

-
-
-

“ ”

P Q

PV r1 r2 s1 s2 r1 r2
Q1 Q2

s1 s2 1

Q1	Q2
.....
P(s1);	P(s2);
P(s2);	P(s1);
.....
r1 r2	r1 r2
.....
V(s1);	V(s2);
V(s2);	V(s1);
.....

Q1	Q2			Q1	P(s1)
P(s2)		Q2	P(s2)	Q1	P(s2)
	P(s1)			P	Q2

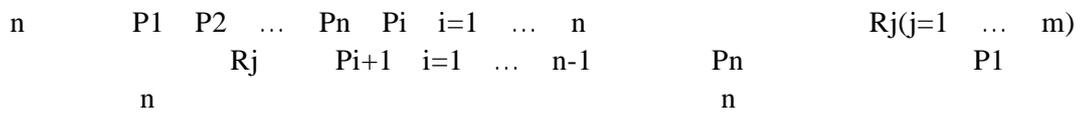
m	n	m < n · K
		m=5

n=5 k=2

			P1	P3	S3	
P2	S1	P2	P1	S1	P3	S2
P2	S2			S3		P3

avoidance)	(deadlock prevention)	(deadlock
	(deadlock detection and recovery)	

3.6.2



3.6.3

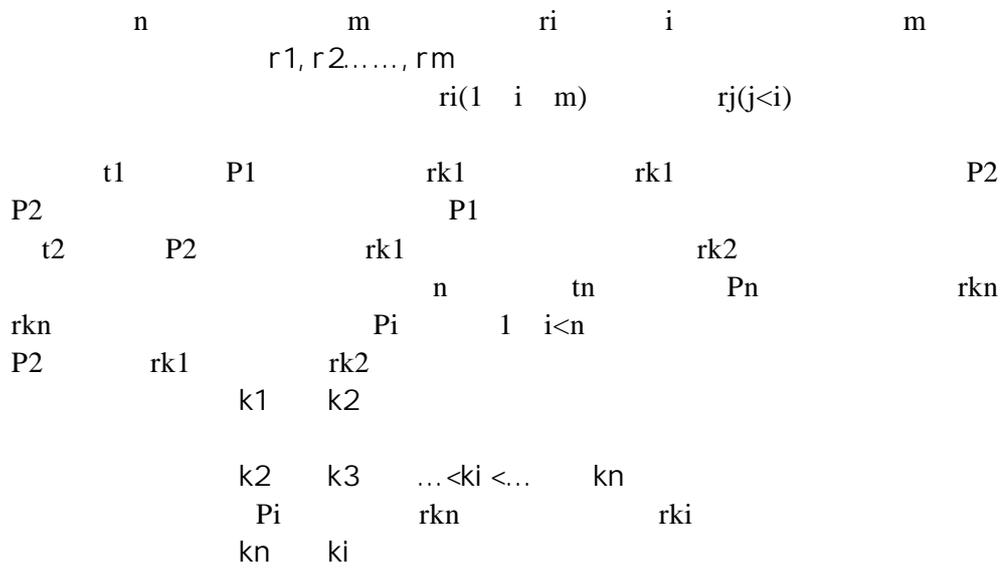
1

1971 Coffman

- mutual exclusion
- hold and wait
- no preemption
- circular wait

2

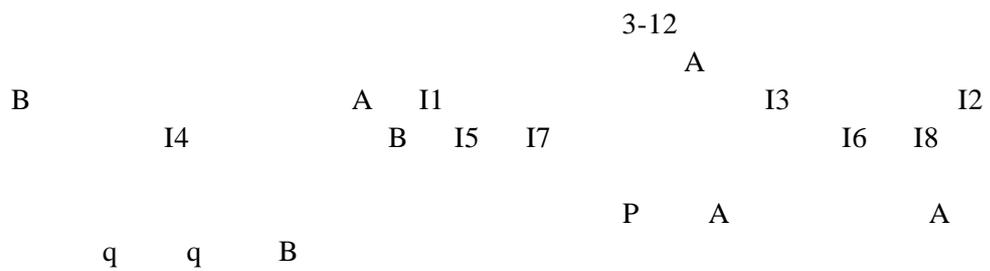
3

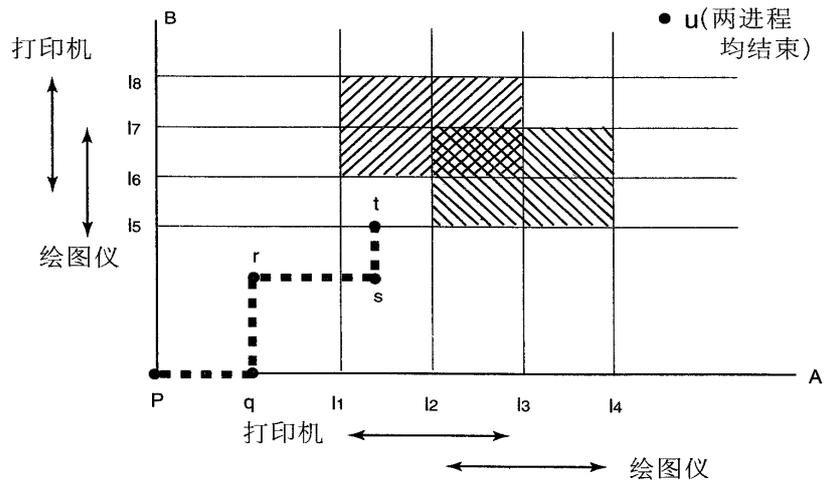


3.6.4

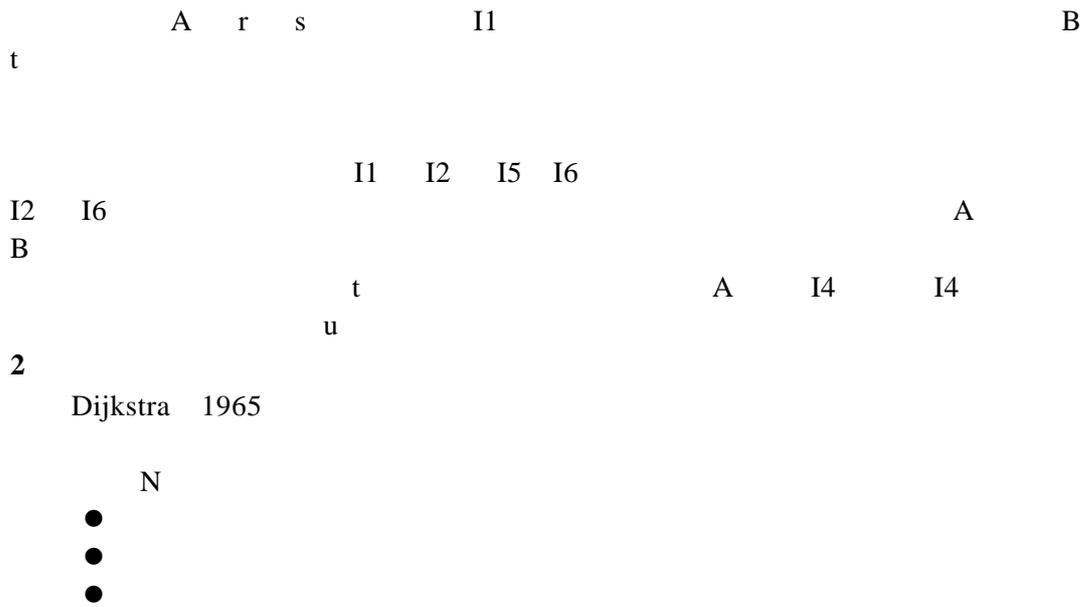
banker's algorithm

1





3-12



3-13 a 4

10

22 3-13 b

名字	已使用 最大	
	↓	↓
Andy	0	6
Barbara	0	5
Marvin	0	4
Suzanne	0	7

可用: 10
(a)

名字	已使用 最大	
	↓	↓
Andy	1	6
Barbara	1	5
Marvin	2	4
Suzanne	4	7

可用: 2
(b)

名字	已使用 最大	
	↓	↓
Andy	1	6
Barbara	2	5
Marvin	2	4
Suzanne	4	7

可用: 1
(c)

3-13

(a)

(b)

(c)

3-13 b

Marvin
Suzanne Barbara
Barbara

4

3-13 b

3-13 c

4

3-14

5

3-14

E					E		P	A
	5	6	3	4	2	2	CD-ROM	P
							CD-ROM	
1								A
2			A					

3

	进程	磁带机	绘图仪	打印机	CD-ROM
A	3	0	1	1	
B	0	1	0	0	
C	1	1	1	0	
D	1	1	0	1	
E	0	0	0	0	

已分配的资源

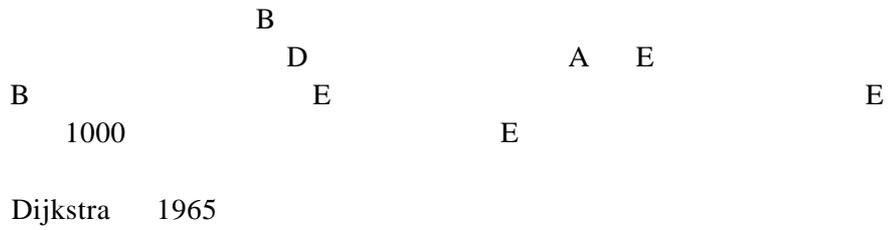
	进程	磁带机	绘图仪	打印机	CD-ROM
A	1	1	0	0	
B	0	1	1	2	
C	3	1	0	0	
D	0	0	1	0	
E	2	1	1	0	

仍需要的资源

E = (6342)
P = (5322)
A = (1020)

3-14

1



4

n m

● -- m

- Resource = (R_1, R_2, \dots, R_m)
- m
- Available = (V_1, V_2, \dots, V_m)
-
- C_{ij} Pi Rj

$$\text{Claim} = \begin{pmatrix} C_{11} & C_{12} & \dots & C_{1m} \\ C_{21} & C_{22} & \dots & C_{2m} \\ \dots & \dots & \dots & \dots \\ C_{n1} & C_{n2} & \dots & C_{nm} \end{pmatrix}$$

- Allocation = A_{ij} Pi Rj

$$\text{Allocation} = \begin{pmatrix} A_{11} & A_{12} & \dots & A_{1m} \\ A_{21} & A_{22} & \dots & A_{2m} \\ \dots & \dots & \dots & \dots \\ A_{n1} & A_{n2} & \dots & A_{nm} \end{pmatrix}$$

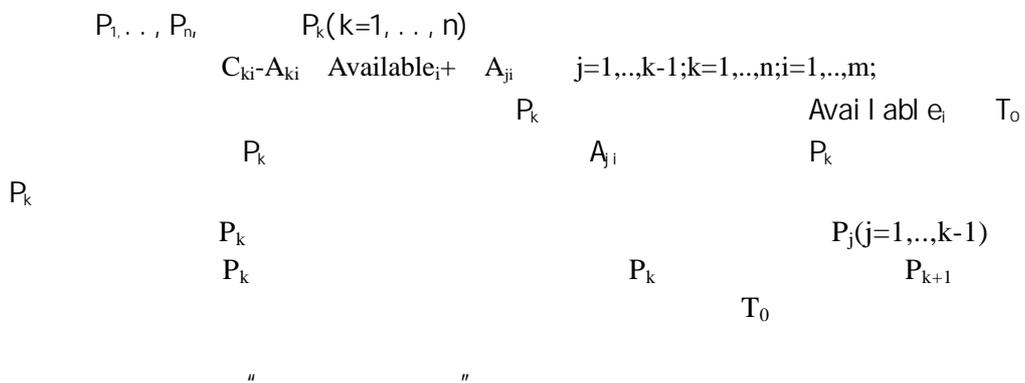
- $R_i = V_i + \sum_{k=1}^n A_{ki} \quad i=1, \dots, m, k=1, \dots, n;$
- $C_{ki} \leq R_i \quad i=1, \dots, m, k=1, \dots, n;$
- $A_{ki} \leq C_{ki} \quad i=1, \dots, m, k=1, \dots, n;$

Ri

$$R_i = C_{(n+1)i} + \sum_{k=1}^n A_{ki} \quad i = 1, \dots, m, k = 1, \dots, n;$$

" "

TO



(banker's algorithm)

```
type state= record
    resource,available:array[0...m-1] of integer;
    claim,allocation:array[0...n-1,0...m-1] of integer;
end
/*
if alloc[i,*]+request[*]>claim[i,*] then <error> /*
else
    if request[*]>available[*] then <suspend process.>
    else /*
        < define newstate by:
        allocation[i,*]:=allocation[i,*]+request[*]
        available[*]:=available[*]-request[*] >
        end;

        if safe(newstate) then
            < carry out allocation>
        else
            <restore original state>
            < suspend process>
        end
    end
end
/*

function safe(state:s):boolean;
var currentavail:array[0...m-1] of integer;
    rest:set of process;
begin
    currentavail:=available;
    rest:={ all process };
    possible:=true;
while possible do
    find a Pk in rest such that
        claim[k,*]-allocation[k,*] < currentavail;
    if found then
        currentavail:=currentavail+allocation[k,*];
        rest:=rest-[Pk];
    else
        possible:=false;
    end
end;
safe:= (rest=null);
end.
```

1
2

2

3

3 P_i
 $allocation[i,*]:=allocation[i,*]+request[*]$
 $available[*]:=available[*]-request[*]$

4 5

P_i

5

- $currentavail$ possible
- $currentavail:=available, possible:=true$
- $possible:=true$ rest $claim[k,*]-allocation[k,*]$
- $currentavail$ P_k
- $currentavail:=currentavail+allocation[k,*]$ P_k
- $rest:=rest-[P_k]$ $possible:=false$
- rest

B 5 C $\{P_0, P_1, P_2, P_3, P_4\}$ A B C A 10
 T_0

Process	allocation			claim			available		
	A	B	C	A	B	C	A	B	C
P_0	0	1	0	7	5	3	3	3	2
P_1	2	0	0	3	2	2			
P_2	3	0	2	9	0	2			
P_3	2	1	1	2	2	2			
P_4	0	0	2	4	3	3			

$C_{ki} - A_{ki}$

Process	$C_{ki} - A_{ki}$		
	A	B	C
P_0	7	4	3
P_1	1	2	2
P_2	6	0	0
P_3	0	1	1
P_4	4	3	1

1 T_0

T_0 $\{P_1, P_3, P_4, P_2, P_0\}$

T_0

	currentavail	$C_{ki} - A_{ki}$	allocation	currentavail+allocation	possible
--	--------------	-------------------	------------	-------------------------	----------

	A	B	C	A	B	C	A	B	C	A	B	C	
P ₁	3	3	2	1	2	2	2	0	0	5	3	2	TRUE
P ₃	5	3	2	0	1	1	2	1	1	7	4	3	TRUE
P ₄	7	4	3	4	3	1	0	0	2	7	4	5	TRUE
P ₂	7	4	5	6	0	0	3	0	2	10	4	7	TRUE
P ₀	10	4	7	7	4	3	0	1	0	10	5	7	TRUE

2 P₁

P₁ 1 A 2 C

- request₁(1,0,2) C_{k₁}-A_{k₁}(1,2,2)
 - request₁(1,0,2) Available(3,3,2)
- P₁ Available Allocation C_{k_i}-A_{k_i}

process	allocation			C _{k_i} -A _{k_i}			available		
	A	B	C	A	B	C	A	B	C
P ₀	0	1	0	7	4	3	2	3	0
P ₁	3	0	2	0	2	0			
P ₂	3	0	2	6	0	0			
P ₃	2	1	1	0	1	1			
P ₄	0	0	2	4	3	1			

? P₁

	currentavil			C _{k_i} -A _{k_i}			allocation			currentavil+allocation			possible
	A	B	C	A	B	C	A	B	C	A	B	C	
P ₁	2	3	0	0	2	0	3	0	2	5	3	2	TRUE
P ₃	5	3	2	0	1	1	2	1	1	7	4	3	TRUE
P ₄	7	4	3	4	3	1	0	0	2	7	4	5	TRUE
P ₀	7	4	5	7	4	3	0	1	0	7	5	5	TRUE
P ₂	7	5	5	6	0	0	3	0	2	10	5	7	TRUE

{P₁, P₃, P₄, P₀, P₂}

P₁

3 P₄

P₄

- request₄(3,3,0) C_{k₄}-A_{k₄}(4,3,1)
- request₄(3,3,0) > Available(2,3,0)

P₄

4 P₀

P₀

- request₀(0,2,0) C_{k0}-A_{k0} (7,3,1)
 - request₄(0,2,0) Available(2,3,0)
- P0

	allocation			C _{ki} -A _{ki}			available		
	A	B	C	A	B	C	A	B	C
P0	0	3	0	7	2	3	2	1	0
P1	3	0	2	0	2	0			
P2	3	0	2	6	0	0			
P3	2	1	1	0	1	1			
P4	0	0	2	4	3	1			

P0 0,2,0

P0

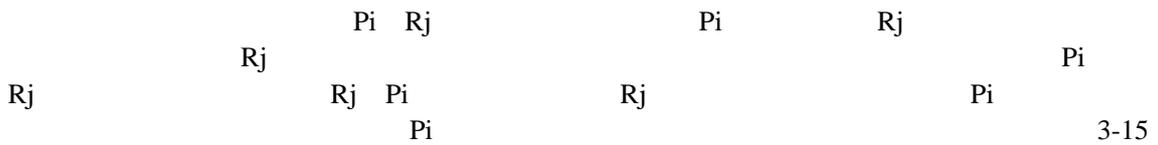
3.6.5

1

" "

-

-



(5) possible[k]=false k=1,...,n possible[k]=true
 Pk

$$A[b_{ij}] = \begin{pmatrix} b_{11} & b_{12} & & b_{1n} \\ b_{21} & b_{22} & & b_{2n} \\ & & & \\ b_{n1} & b_{n2} & & b_{nn} \end{pmatrix}$$

warshall " " warshall A [b_{ij}] b_{ij} [b_{ij}]

for k:=1 to n do
 for i:=1 to n do
 for j:=1 to n do
 b_{ij} := b_{ij} or (b_{ik} and b_{kj})

Pk b_{ij} i j b_{ik} b_{kj} i P_i

 k j warshall A A

b_{ii}=1 (i=2.....,n)

P₃ P₁ P₂

$$A[b_{ij}] = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}$$

j b₁₂ b₂₃ b₃₁ 1 k=1 i i

 b₃₁=1 b₁₂=1 b₃₂ 1 k=2 i j

i j (b₁₂ b₂₃) b₃₁ 1 (b₃₂ b₂₃) b₃₃ 1 k=3

 b₁₁ b₂₂ 1

P₁ P₂ P₃

-
-
-

?

CPU

-
- checkpoint
-

3.7 Windows 2000/XP

3.7.1 Windows 2000/XP

Windows 2000/XP

- CreateMutex API
- OpenMutex
- ReleaseMutex
- CreateSemaphore API
- OpenSemaphore
- ReleaseSemaphore API
- CreateEvent
- OpenEvent
- SetEvent PluseEvent
- ResetEvent

WaitForMultipleObjects

WaitForSingleObjecs

Windows 2000/XP
API

```

CRITICAL_SECTION          API
InitializeCriticalSection(          ) EnterCriticalSection(
          ) TryEnterCriticalSection(
0) LeaveCriticalSection(          )
DeleteCriticalSection(          )

```

API
) InterlockedCompareExchange()
 InterlockedExchangeAdd() InterlockedDecrement(1)
) InterlockedIncrement(1)

3.7.2 Windows2000/XP

Windows2000/XP (signal)

(1)SetConsoleCtrlHandler
 GenerateConsoleCtrlEvent
 CTRL_C_EVENT
 SetConsoleCtrlHandler
 CTRL+C
 CTRL+C
 (2)signal raise
 UNIX

Windows2000/XP (shared memory)

Windows2000/XP
 (file mapping)
 CreateFileMapping(
) OpenFileMapping(
 MapViewOfFile(
) FlushViewOfFile(
) UnmapViewOfFile(
 CloseHandle(
 Windows2000/XP (pipe)

I/O
 CreatePipe
 ReadFile WriteFile
 C/S
 ConnectNamePipe
 CreateNamePipe
 CallNamePipe
 ReadFile WriteFile(
) ReadFileEx
 WriteFileEx(
 Windows2000/XP (mailslot)

C/S
 CreateMailslot(
) GetMailSlotInfo(
) SetMailslotInfo(
) ReadFile(
) CreateFile(
 WriteFile(
 Windows2000/XP (socket)

C/S
 Windows2000/XP TCP/IP
 "Winsock" BSD
 API

3.8 Linux

Linux

```

struct semaphore {
    atomic_t count;
    int waking;
    struct wait-queue *wait;
};

count
0 count 1( ) count ( ) sema-init
waking wait-queue up
down( P ) up(
V ) down-interruptible wake-up wake-up-interruptible Down up
down-interruptible
0
wake-up wake-up-interruptible
TASK-INTERRUPTIBLE

```

```

1 1 Linux semary
semid-ds
static struct semid-ds *smeary[SEMMNI]

struct sem {
    int semval;
    int sempid;
};

struct sem semval 0 sys-semctl
Sempid pid
struct semid-ds {
    struct ipc-perm sem-perm;
    -kernel-time-t sem-otime; /*
    -kernel-time-t sem-ctime; /*
    struct sem *sem-base;
    struct sem-queue *sem-pending;
    struct sem-queueq **sem-pending-last;
    struct sem-undo *undo;
    unsigned short sem-nsems; /*
};
semid-ds
sem-base " "
SEMMSL sem-pending
sem-pending-last
undo
struct sem-queue { /* sem-queue
    struct sem-queue *next; /*

```

```

struct sem-queue **prev; /*
struct wait-queue *sleeper; /*
struct sem-undo *undo; /* sops
int pid; /*
int status; /*
struct semid-ds *sma; /* struct sem-ds
struct sembuf *sops; /*
int nsops; /*
int alter; /*
};
struct sem-undo { /* undo sem-undo
struct sem-undo *proc-next; /* undo
struct sem-undo *id-next; /* undo
int semid; /* semid-ds
short *semadj; /*
};
semget semop semctl Semget
semctl semop ( )
0 Linux
0
sem-queue sleeper sem-queue
sem-queue sem-pending sem-pending-last Linux sem-pending
sem-queue sem-queue
Linux
semid-ds task-struct sem-undo

```

3.9

CH2

(Dekker)

Hoare

Dijkstra

CPU

P V

Hansen

—

P V

PV

/

P

V

- 1.
- 2.
- 3.
- 4.
- 5.

6.
 7.
 8.
 9. Dekker
 10. Peterson
 11.
 12.
 13. P V
 14. P V
 15.
 16.
 17.
 18. Hanson Hoare
 19. P V ?
 20.
 21.
 22. pipeline
 23. ?
 24.
 25.
 26. ? ?
 27.
 28.
 29. ? ?
 30. P V
 31. /
 32. 20 65 3
 33. ?
 34. ?
 35. 8 N 3 N
 36. ?
 37.
 38. ?
 39. m n x (l x m)

40. (n m x)
 ?
 41. ?
 42.

43. P V
 44.
 45.

	1/4

1 R M P

1 K
 2 K

P V

2 n
 (1)
 (2) m m n

3 P1 P2 S1 S2 0

P1 P2 x y z

<p>P1 begin y:=1; y:=y+3; V(S1); z:=y+1; P(S2); y:=z+y end.</p>	<p>P2 begin x:=1; x:=x+5; P(S1); x:=x+y; V(S2); z:=z+x; End.</p>
---	--

4
V 2 100 1 P

5 P1 P2 P1 P2

6 P1 P2 Wait Signal

7 P V

8 4 1 2 3

9 S P V S S 0 S=0 S<0

10 Process P Process Q
begin begin
A C
B D
C end
end

11 1
2

12 (1)
(2)

13 1
2

14 (1) (2) P V

15 Bakery Algorithm Lamport 1974

16

var choosing array [0...n - 1] of boolean

```

number array [0...n-1] of integer
repeat
  choosing [i] =true
  number[i] = 1 max ( number[0] ... number[n-1])
  choosing[i] =false
  for j = 0 to n-1 do
  begin
  while choosing [j] do {nothing}
  while number[j] = 0      number [j] = j      number[i] = i
  do {nothing}
  end
  {critical section}
  number [i]:=0
  {remainder}
forever
  choosing = number = false = 0 = i = i
  a b c d a c or a=c b d
  1 2 3

```

17 (patil 1971)

(1) P

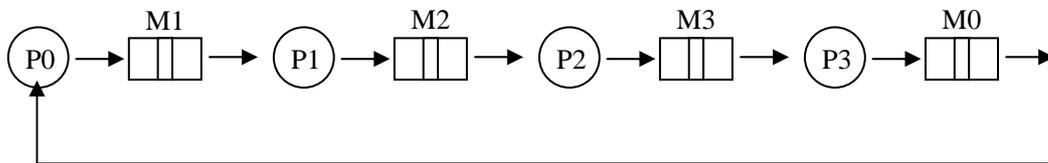
V (2)

18 $P_i \ i=0...3$ $M_j \ j=0...3$

P_i M_i $M_{(i+1) \bmod 4}$ $M_0 \ M_1 \ M_2 \ M_3$

3 3 2 2 M_0 PV

$P_i \ i=0...3$



19 $P_i \ Q_j \ R_k$ $P_i \ Q_j$ M_1

$buf1 \ Q_j \ R_k$ M_2

$buf2 \ P_i$ $buf1 \ Q_j$

$buf2 \ R_k$ $P \ V$

20 $P \ Q$ $P \ 1$

W Q 1
 INPUT OUTPUT I/O Delay
 seconds
 21 m n m n m n
 ?
 22 N M / M
 M+N
 23 100
 P B A
 200 ?

24 Jurassic m n

n P V m
 n
 25 k 1 2 ... k file
 K 1) P V 2)

26 Available=(1 1 2)

	Claim			Allocation		
,	R1	R2	R3	R1	R2	R3
P1	3	2	2	1	0	0
P2	6	1	3	5	1	1
P3	3	1	4	2	1	1
P4	4	2	2	0	0	2

(1) Cki-Aki?

(2) ?

(3) P₁ request1(1 0 1)

(4) P₁ P₀ request0(1 0 1)

(5) P₀ P₂ request0(0 0 1)

27 A B C D 4 P0 P1 P2 P3 P4

Process	Allocation				Claim				Available			
	A	B	C	D	A	B	C	D	A	B	C	D
P ₀	0	0	3	2	0	0	4	4	1	6	2	2
P ₁	1	0	0	0	2	7	5	0				
P ₂	1	3	5	4	3	6	10	10				
P ₃	0	3	3	2	0	9	8	4				
P ₄	0	0	1	4	0	6	6	10				

(1)

(2) P1 request1(1 2 2 2)

28

Available=(1 0 2 0)

$$\text{Need} = \begin{pmatrix} 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 2 \\ 3 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 2 & 1 & 1 & 0 \end{pmatrix} \quad \text{Allocation} = \begin{pmatrix} 3 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

(1)

(2) request2(0 0 1 0)

(3) request5(0 0 1 0)

29 150 P1 70, 25 P2

60, 40 P3 60, 45
(1)P4 P4 60, 25 (2)P4 P4

30 60, 35 X Y (1)

X Y (2) -N<X -Y <M N M
P V X Y

31 A B m A B

A B n(n<m)
P V

32 A1 A2 ... An1 m B1 B2 ... Bn2

(1)

(2) B1 B2 ... Bn2

(3) M

33 PV R1 3 R2 4 P1 P2 P3 P4 4

	R1	R2	R1	R1	R2	R1
(1)						
(2)						—
34						PV

35	-		(1)	P V	(2)
36		R1	R2		

```

        c1:=1-c2
    until c2<>0;
    Critical Section;    (/*      */)
    c1:=1
    until false
end;
procedure p2;    (/*      p2 */)
begin
    repeat
        Remain Section 2;
        repeat
            c2:=1-c1
        until c1<>0;
        Critical Section;    (/*      */)
        c2:=1
        until false
    end;
begin    (/*      */)
    c1:=1;
    c2:=1;
    cobegin
        p1;p2    (/*      p1, p2      */)
    coend
end.
42                                ?
repeat
    key =true
    repeat
        swap(lock key)
    until key=false
        Critical Section;    (/*      */)
    lock =false
    other code
    until false
43                                1
const n=50
var tally    integer
procedure total()
var count    integer
begin
    for count =1 to n do tally =tally+1
end

```

```

begin (/*main program*/)
  tally =0
  cobegin
    total() total()
  coend
  writeln(tally)
end
(1)          tally
(2)          tally
44
var blocked array[0 1] of boolean
  turn 0 1
procedure P(id integer)
begin
  repeat
    blocked[id] =true
    while turn = id do
      begin
        while blocked[1-id] do skip
        turn =id
      end
      {Critical Section}
      blocked[id] =false
      {remainder}
    until false
  end
begin
  blocked[0] =blocked[1] =false
  turn =0
cobegin
  P[0] P[1]
coend
End
45          P1 P2 P3

```

```

          P1
          P2
          P3
1)      2)      P V
46          A B          A          B

```

```

47      1)      P V      2)
.      1)      P V      2)      ,
48
49  P1 P2 P3      F P1 F      P2 F      P3 F
      F      (1)      P V      (2)
50  100      10
      10
      P V
51      10000
Monitor, Monitor
52      2m      2n
(1)      (2)
      (3)
      P V
53  10      3
      1
54      free[index]      n
(index=0 ... n-1) free[index]=true      index      free[index]=false
      index      acquire      release
55 AND      P      " AND"
" P V      SP SV(Simultaneous P V) SP(s1 s2 s_n) VS(s1 s2
s_n)
      procedure SP(var s1 ,s_n semaphore)
      begin
      if s1>=1 & & s_n>=1 then begin
      for i:= 1 to n do
      s_i:= s_i-1;
      end
      else begin
      { s_i<1 s_i
      SP      };
      end

```

```

        procedure VP(var s1 ..., sn semaphore)
        begin
            for i:=1 to n do begin
                si:=si+1;
                { si };
            end
        AND
56    AND SP SV -
57    AND SP SV
58    AND SP
    AND SV ?
59    ( )
    AND SP(s1,t1,d1;...;sn,tn,dn)
SV(s1,d1;...sn,dn)
    procedure SP(s1,t1,d1 ... sn,tn,dn)
        var s1,...,sn:semaphore;
            t1, ...,tn:integer;
            d1,...,dn:integer;
        begin
            if s1>=t1& ...&sn>=tn then begin
                for i :=1 to n do
                    si:=si-di;
                end
            else
                { si<ti si SP };
            end
        end
    procedure SV(s1,d1;...sn,dn)
        var s1,...sn:semaphore;
            d1,...dn:integer;
        begin
            for i:=1 to n do begin
                si:=si+di;
                { si };
            end
        end
    end
    ti di

```

- 60
- SP(s,d,d)
 - SP(s,1,1)
 - SP(s,1,0)

61 -

CH4

GB



Pentium

Windows2000/XP

Linux

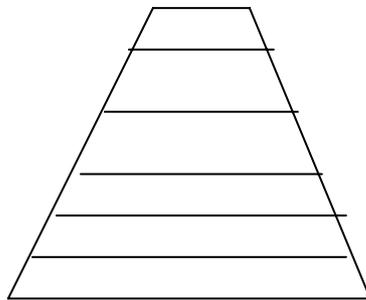
Intel

4.1

4.1.1

4-1

7



4-1

word

0.1

0.55

1μ s

100

40GB

60ns

CPU

50%

128MB

1TB

15ns

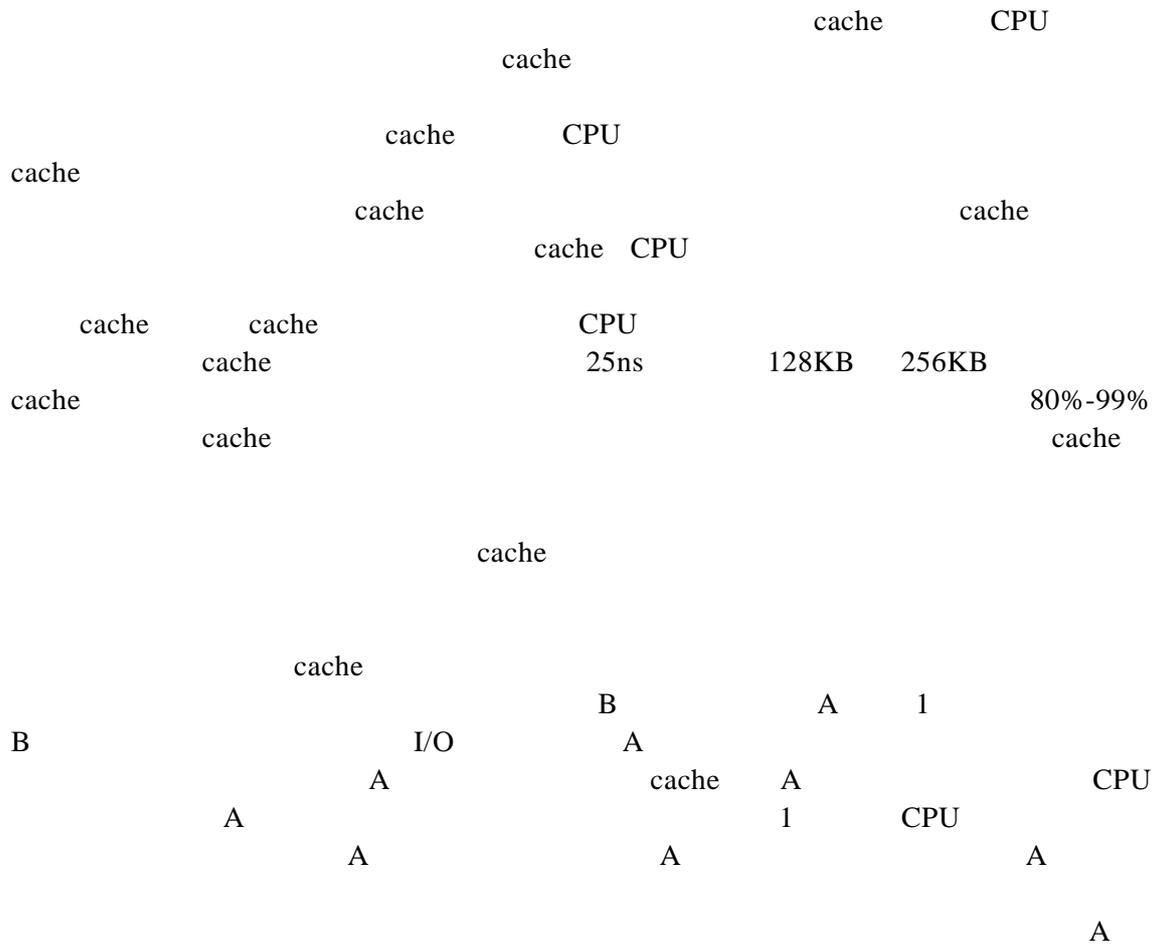
512KB

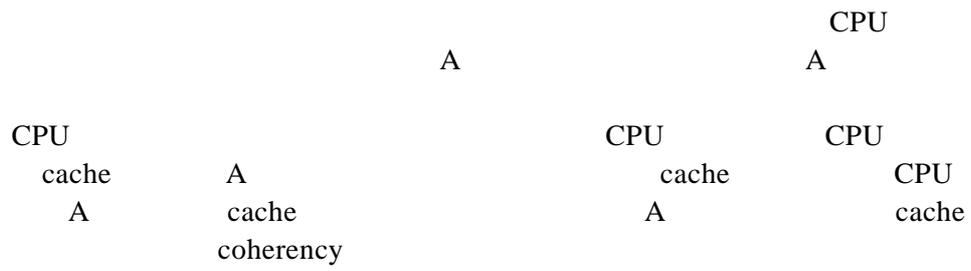
μ s

μ s

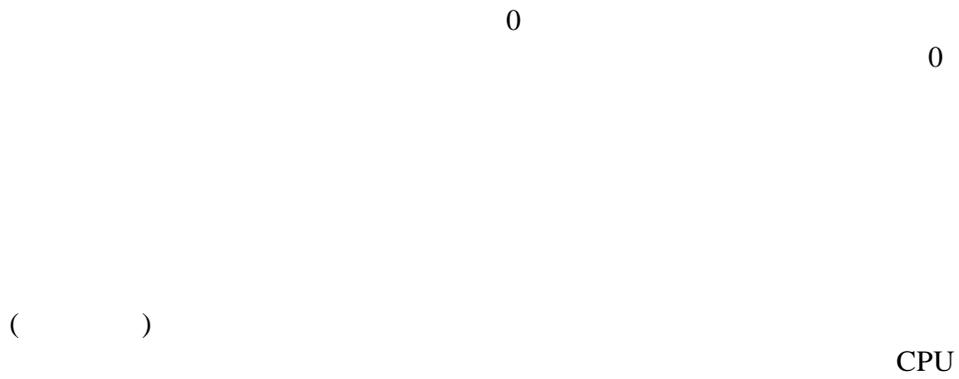
/

4.1.2 caching





4.1.3



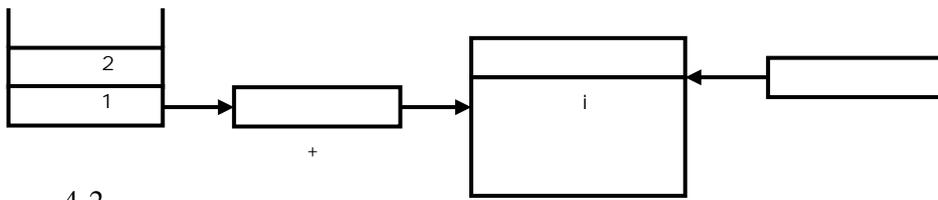
4.2

4.2.1

fence register

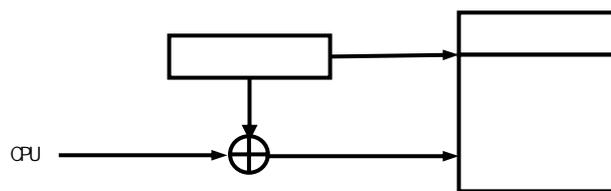
4-2

()



CPU

4-3



•

•

•

70

IBM7094

FORTRAN

IBM1130

MIT CISS cromemco CDOS Digital Research
 Dyhabyte CP/M DJS0520 0520FDOS

4.2.2

8K
1 8K
2 16K
3 16K
4 16K
5 32K
6 32K

4-4

fixed partition

OS/MFT Multiprogramming with a Fixed Number of Tasks " " IBM

4-5

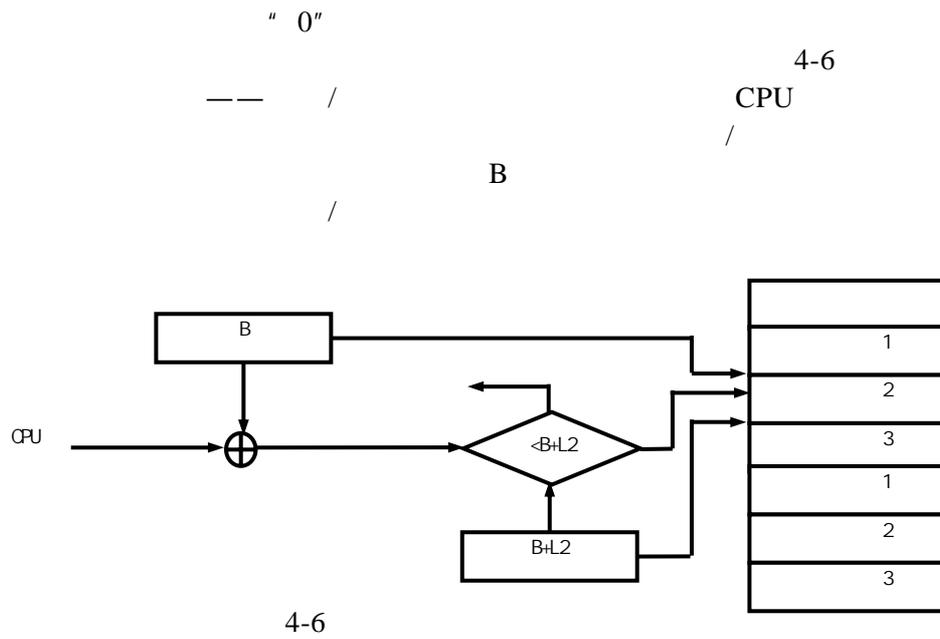
1	8K	8K	0
2	16K	16K	Job1
3	32K	16K	0
4	48K	64K	0
5	64K	32K	Job2
6	96K	32K	0

4-5

" 0" " 0" " 0" " 0"

Job1 Job2 4-5 2 5

4-4



CPU

			4-5	Job1	Job2
10	18	16	32		20

4.2.3

1

variable partition

IBM

OS/MVT Multiprogramming with a Variable Number of Tasks

1 (first fit)

" "

2) (next fit)

3 (best fit)

4 (worst fit)

5) (quick fit)
n

2KB

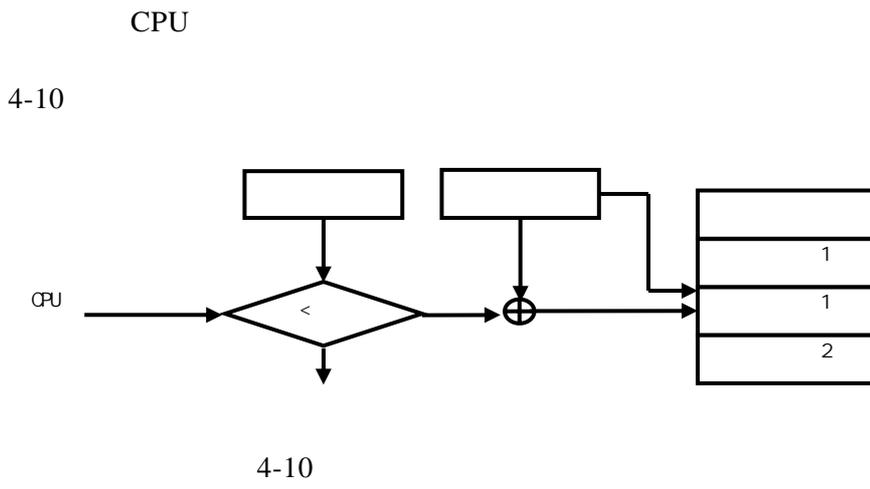
4KB

9KB

8KB

8KB

2

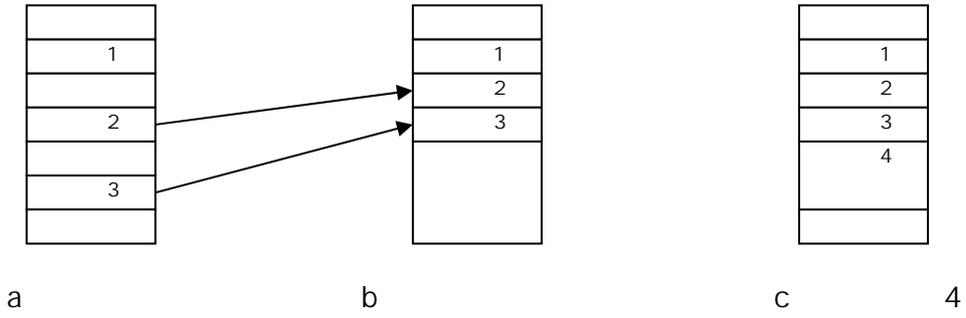


PSW

CDC6600

3

/



4-11

4-11

()

1 i xK
 2 xK 5
 3 xK i
 4
 5 xK /
 6 i /

A B

A
B

A B

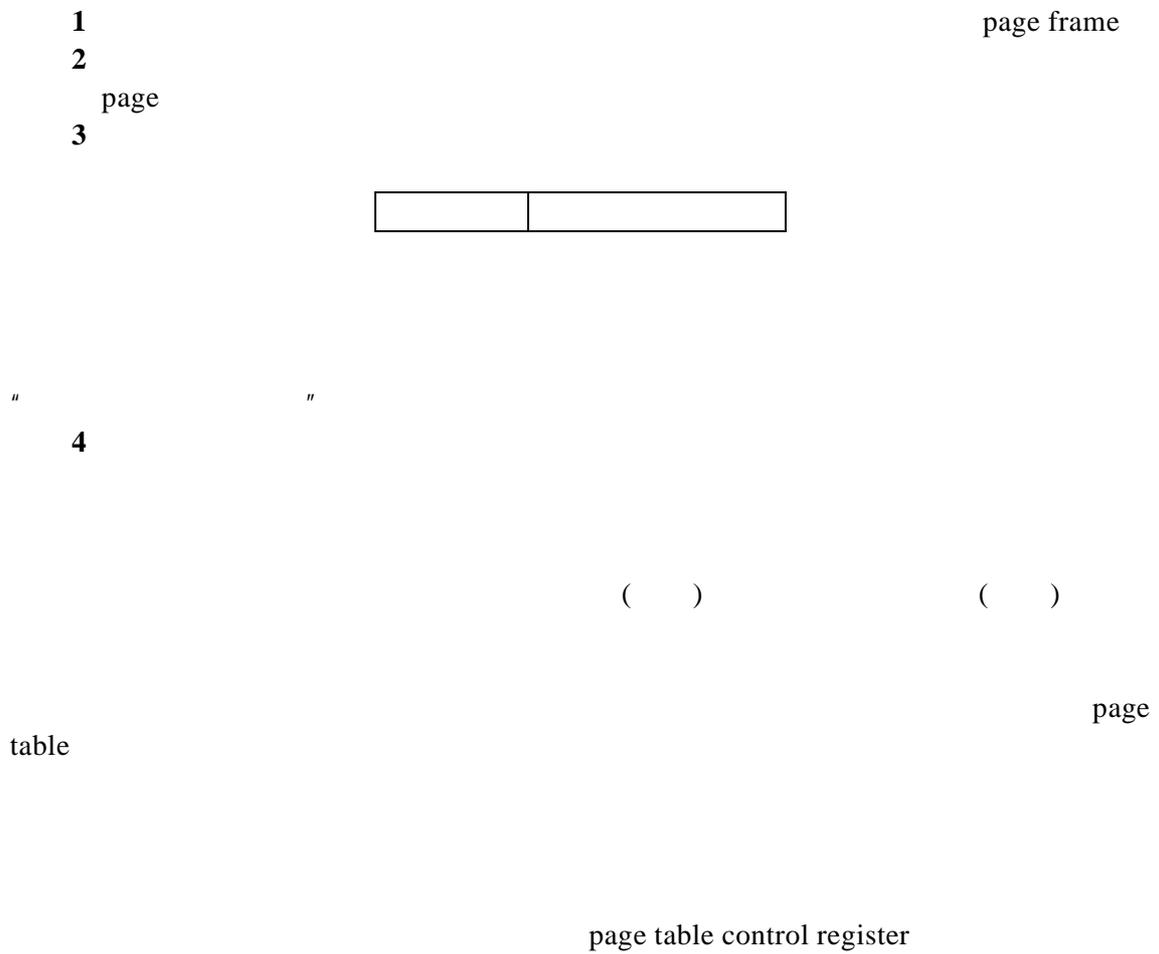
n

$$1/2 + 1/3 + \dots + 1/n \cdot V$$

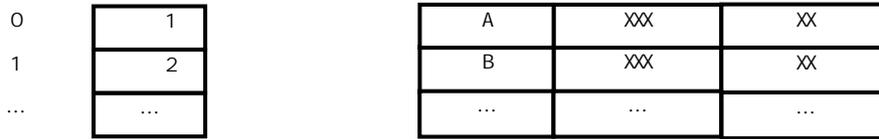
/n× V

4.3

4.3.1



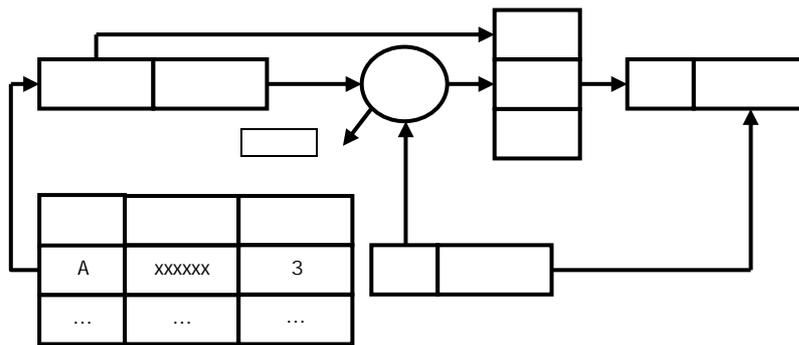
× +



4-12

4-13

×



4-13

CPU

4.3.2

/

/

MMU()

associative

memory

TLB Translation Lookaside Buffer

Intel 80486

32

MMU

" "

(hit ratio)

100%

0

100

20

32

$$\frac{100 \times 90\%}{200} \times 90\% = 100 + 100 \times \frac{1 - 90\%}{200} = 128$$

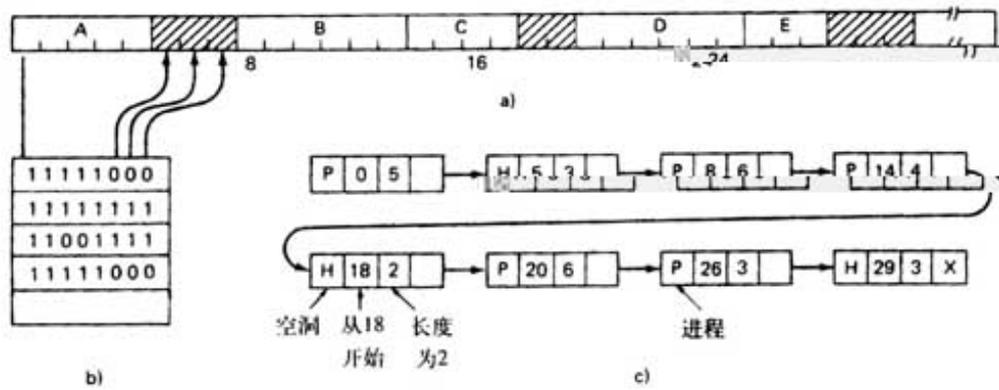
CPU

4.3.3

0/1

/

4-14



4-14

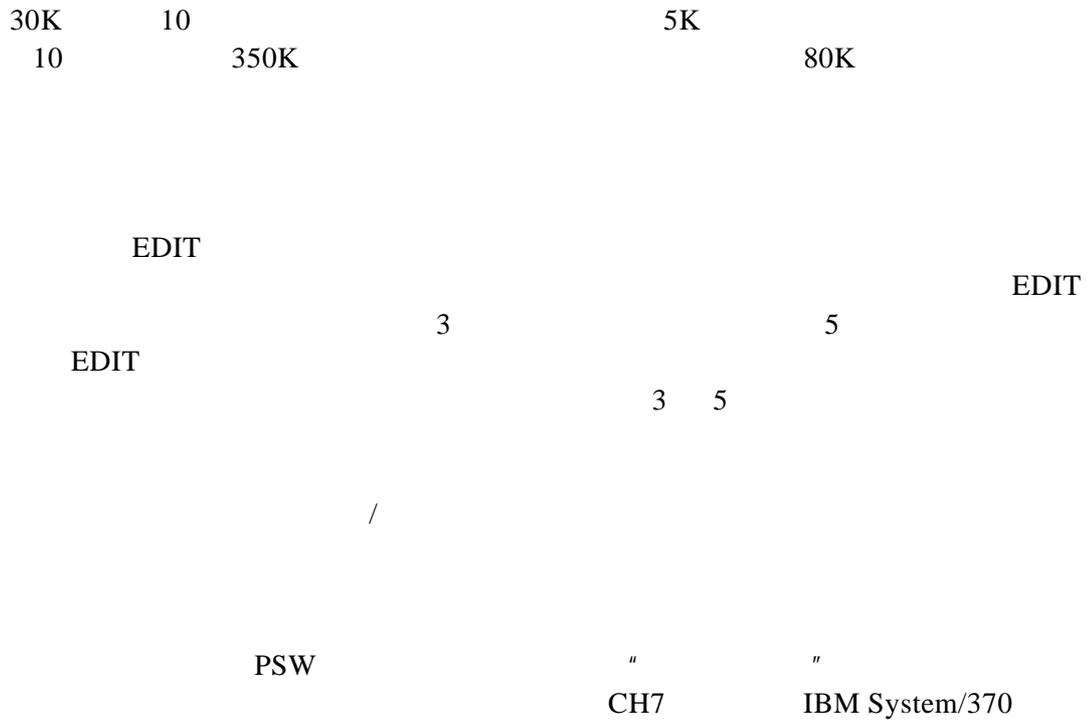
" 0"

" 0"

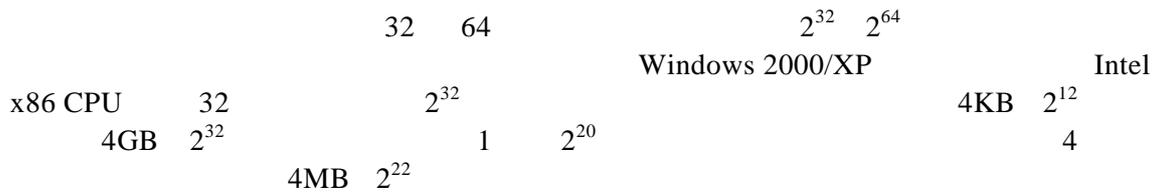
4-14 c
H

P

4.3.4



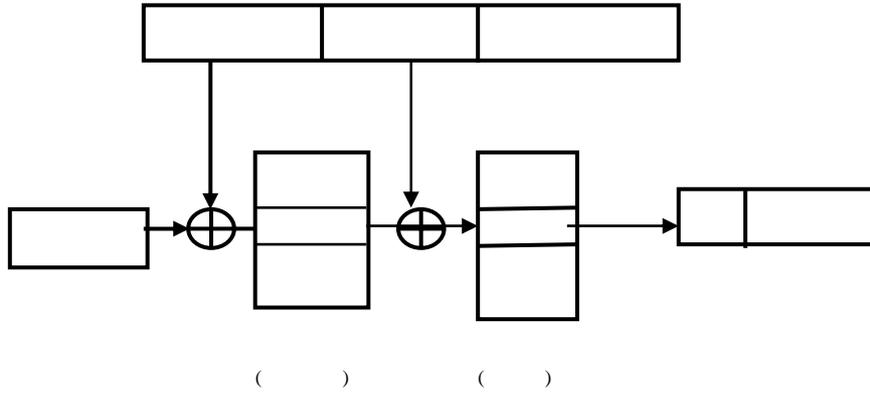
4.3.5



page directory table

4-15

CPU



4-15

) " " " " (

" " " " "

64

SUN

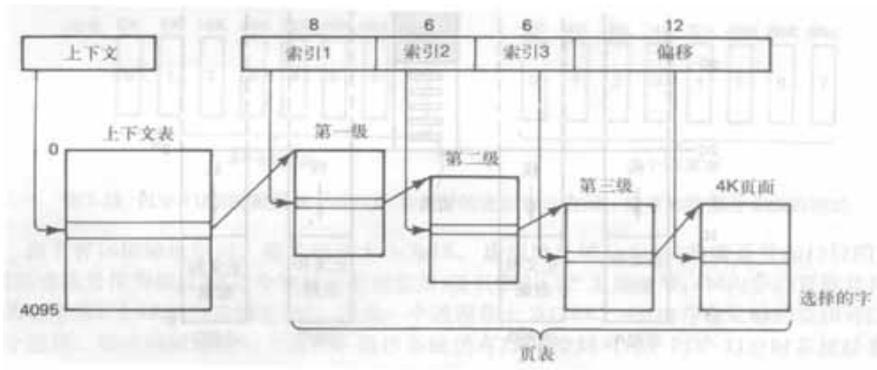
SPARC

4-16

4096

CPU

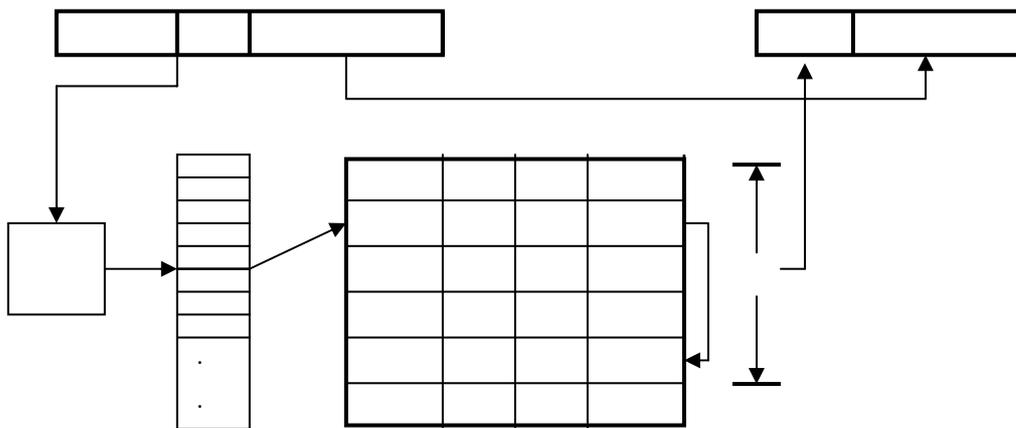
MMU



4-16 SPARC

4.3.6

	IBM AS/400	Mac OS	IPT	Inverted Page
Table	IPT			
	4-17			
MMU			MMU	IPT
	IPT			



4-17

IPT IPT
128KB 128KB

IPT

128MB

1KB

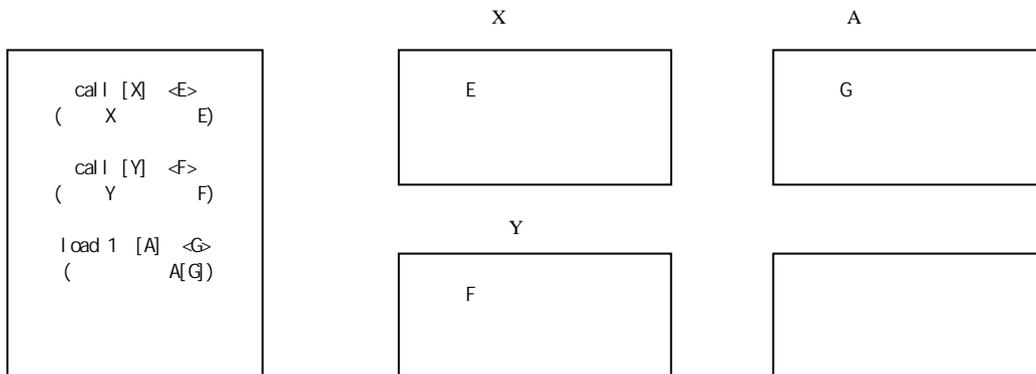
4.4

4.4.1

0

4-18

" 0"



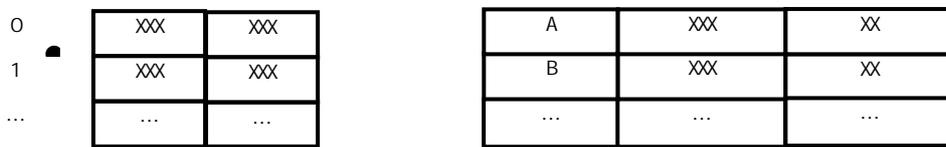
4-18

4.4.2



--

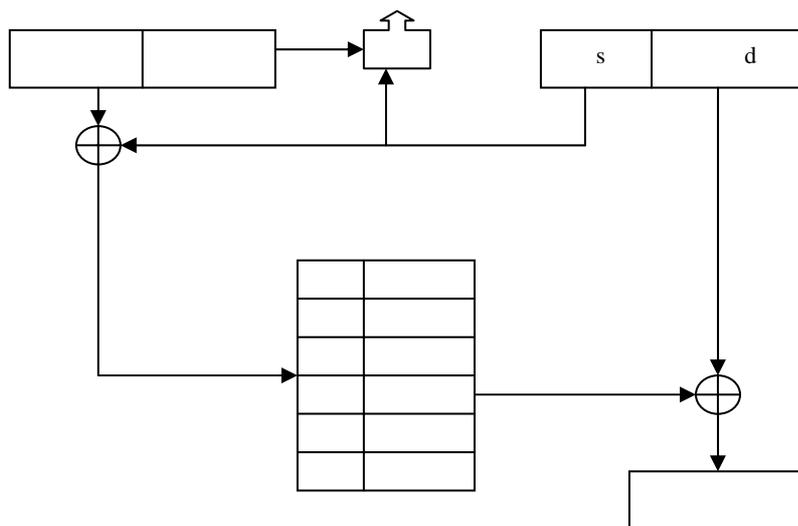
4-19



4-19

/

4-20



4-20

4.4.3

/

4.4.4

()

()

4.5

4.5.1

"

"

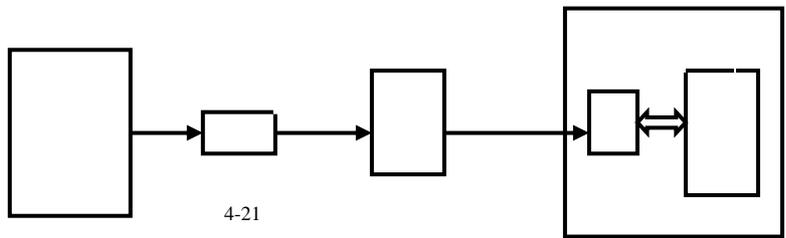
" **virtual memory** (Fotheringham 1961)

"

"

"

4GB 32 20 4GB Windows2000/XP 1MB Intel pentium
 4-21



locality 1968 P.Denning Knuth **principle of**
 Tanenbaum Huck Fortran

70 60 MULTICS Atlas 60
 IBM

I/O

4.5.2

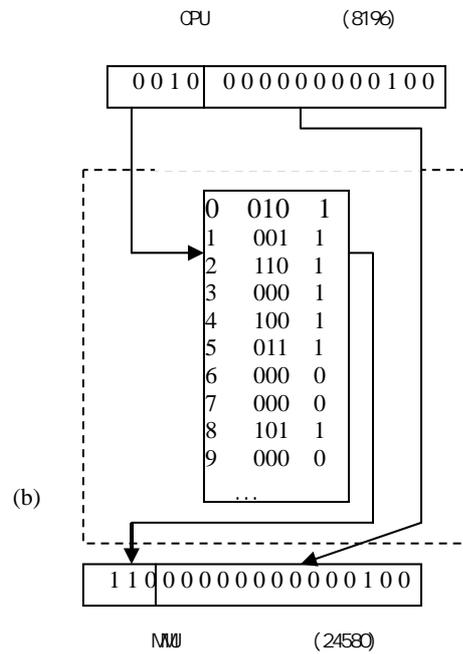
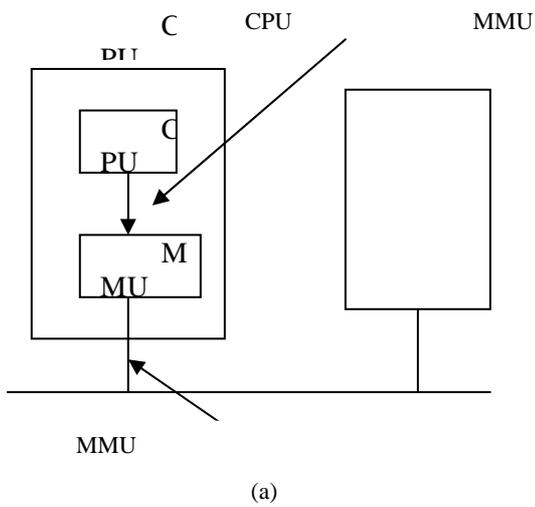
1

MMU Memory Management Unit

MMU

4-22(b)

4-22(a)



4-22(a) MMU
(b) 16 4KB MMU

- MMU
- TLB MMU TLB
- TLB
- TLB TLB
- TLB TLB MMU MMU
- MMU
- MMU

- MMU 4-22(b) 8196(
 - 0010000000000100) MMU 16 4
- 12
 " " 1
 12 15 (" " 0
 1100000000000100)

2

paging

demand

4-23

4-23

0

CPU

()

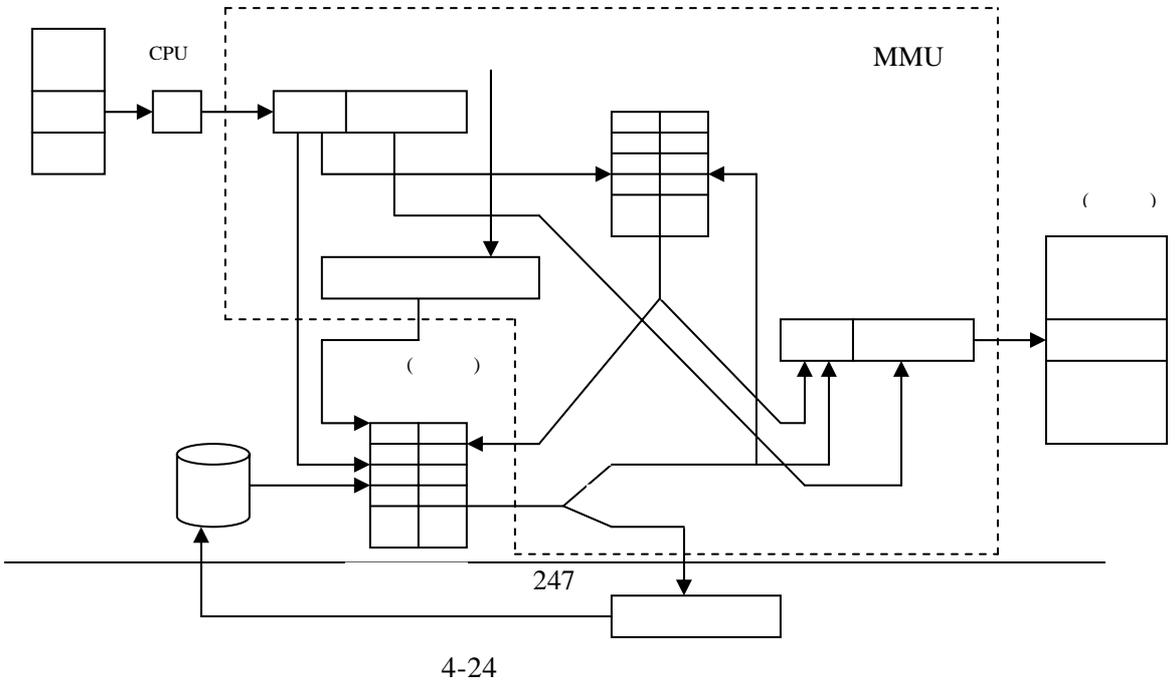
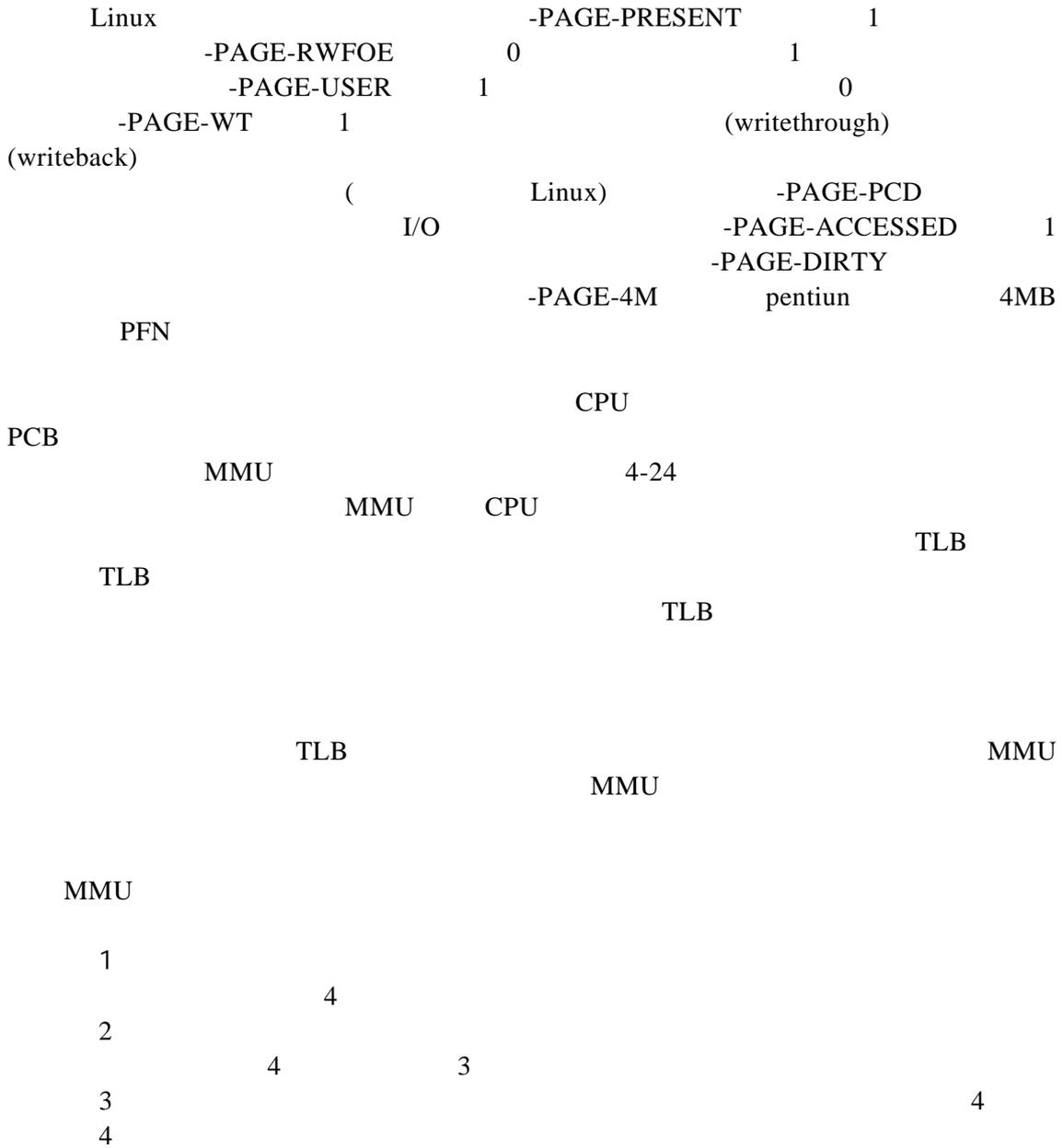
()

Modified

Referenced

()

()



VM/370 Honeywell 6180 MULTICS UNIVAC IBM/370 70/64 VS/1 VS/2
VMS

3

demand paging

prepaging

I/O

I/O

50%

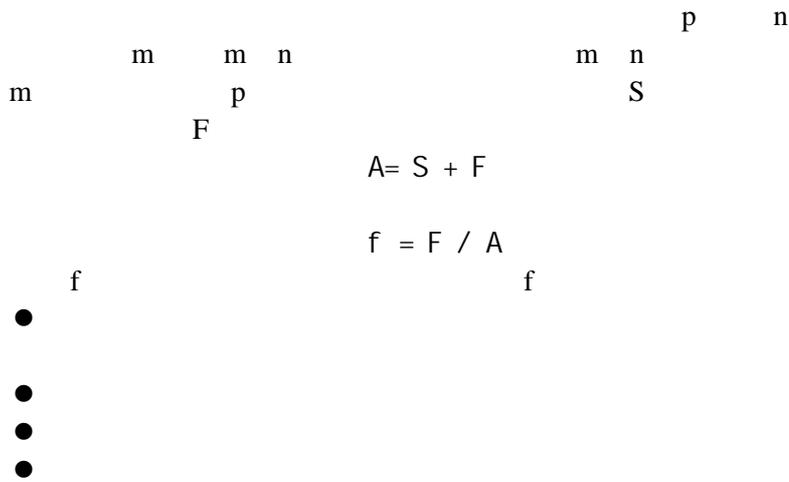
I/O

CPU

4

I/O

Thrashing



128 × 128 " 0"
128

```

Var A array[1..128] of array [1..128] of integer;
For j = 1 to 128
do for i = 1 to 128
do A[i][j] = 0
[i][j]

```

128 × 128

1

```

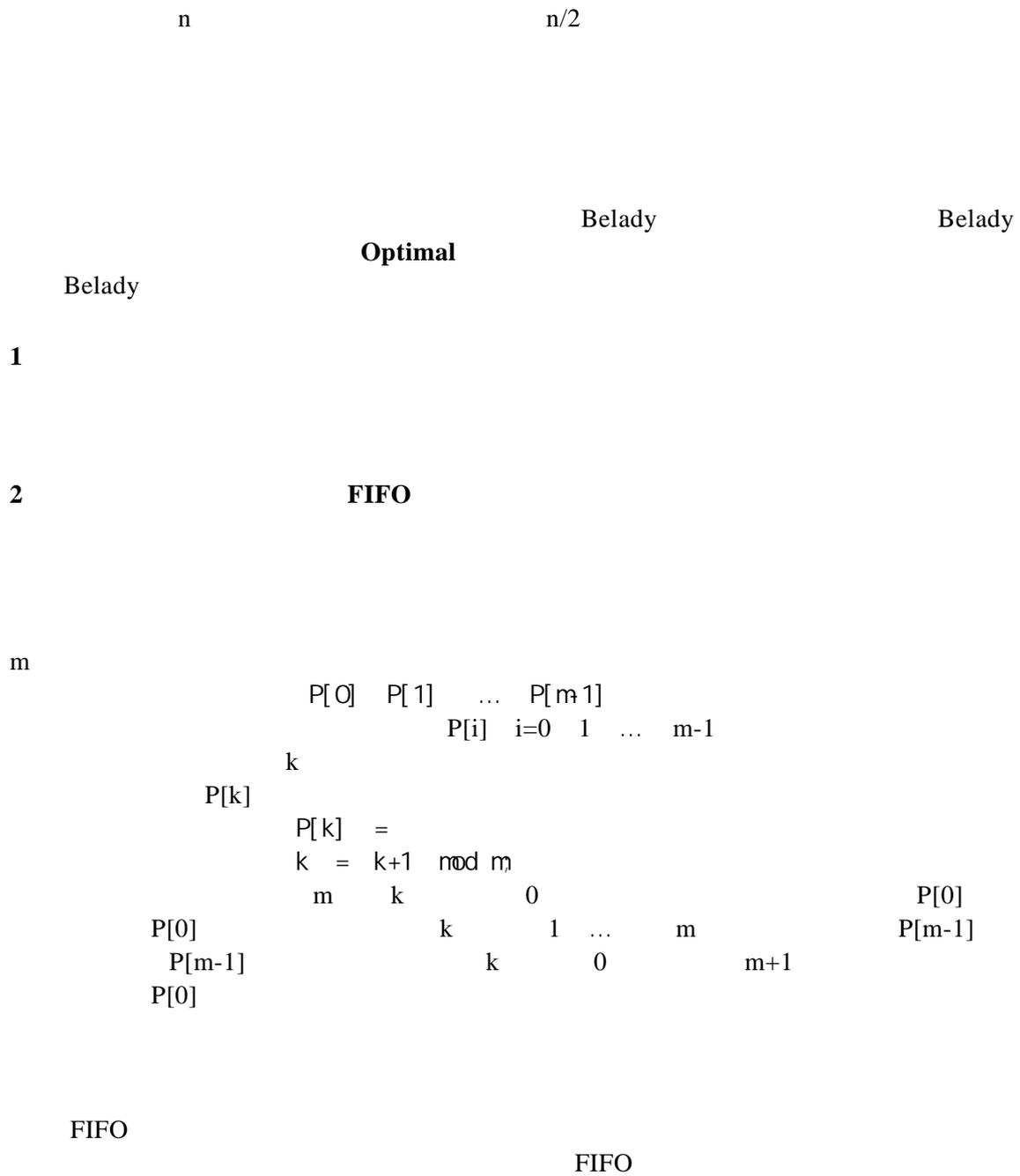
Var A array[1..128] of array[1..128] of integer;
for i = 1 to 128

```

```

do for j = 1 to 128
  do A[i][j] = 0
    128 - 1

```



2 3

3

LRU Least Recently Used

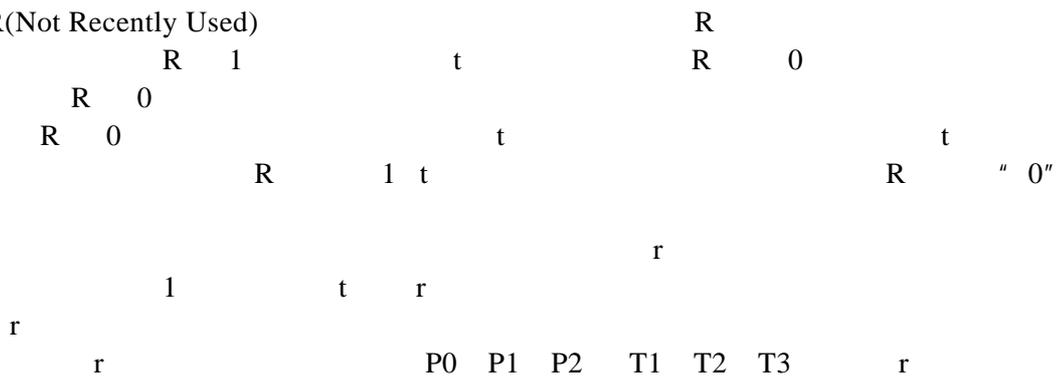
4 3 0 4 1 1

2 3 2

4	4	
3	4 3	
0	4 3 0	
4	3 0 4	
1	0 4 1	3
1	0 4 1	
2	4 1 2	0
3	1 2 3	4
2	1 3 2	

LRU

NUR(Not Recently Used)



	T1	T2	T3
P0	1000	0100	1010
P1	1000	1100	0110
P2	0000	1000	0100
T3		P2	P0
P1		T3	T2

T1 P2

" "

LFU Least Frequently used

t

" "

4

FIFO

" "

FIFO

FIFO

FIFO

" " 0
" " 1

" " 0

second chance

5

Clock Policy
FIFO

Clock

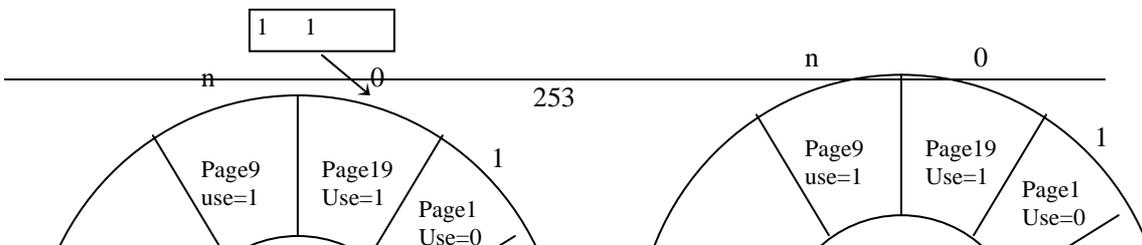
" "

-
-
-
-

" " 1
" " 1
" " 0
" " 0
" " 1
" " 0

4-25

page727 page45 2 Clock
page45 " " 1 " " 0
page191 3 " " 0
page556 4 " " 0 page556 page727
page727 " " 1 page13 5



```

                                "      "      "      "
1                                r=0  m=0
2                                r=1  m=0
3                                r=0  m=1
4                                r=1  m=1

1                                r=0  m=0
"      "      r=0  m=0
2      1                                r=0  m=1
                                "      "  r
0      3      2                                "      "
r      0      1      2

```

Macintosh

UNIX SVR4

clock

clock

" " 0 " " 1
 " " 0
 " " 1
 " " 0
 clock

Opt FIFO LRU Clock

3 2 5 2 4-26 5 2 3 2 1 5 2 4 5
 Opt 3 page1
 page2 page4
 LRU 4
 page3 page1 page2 page4
 FIFO

6

page2 page3 page1 page5 page2 page4
 Clock 5
 * " " 1
 page5 " " 1
 page2 page5 page2 page3 page1 " " 0
 page2 page3 page2 page3
 page4 page4 page1 page3
 page2 " " 1 page2 " " 1 page4
 " " 0 page3 " " 1 page5 " " 1 page4
 FIFO Opt Clock LRU

O	P	T	F(1)	F(2)	F(4)
2	2	2	2	4	2
	3	3	3	3	3
			5	5	5

L U F(3) F

	2	3	2	2	5	2	4	5	3	2	5
			3	1	1	5	2	4	5	3	3

2

3

4

512 4K

S
S >> P

e

I/O

P

S/P

P/2

Se/P

$$= \left(\frac{Se}{P} + \frac{Se}{P} + \frac{P}{2} \right) = \frac{Se}{P} + \frac{P}{2}$$

P/2) P (Se/P) =Se/P+P/2 P O

$$-Se/P^2 + 1/2 = 0$$

P = $\frac{2Se}{P}$

S = 128KB e = 8B

1448 48 4KB 2048

2KB 512B 4KB 2048

4096	VAX	512	Atlas	512	512	48	IBM370	2048
4096	Pentium	4096	IBM AS/400	4	MIPS R4000	4096	Motorola 68040	Macintosh
7			CPU		MMU			16

(2)

CPU

/

1968 P.J. Denning working

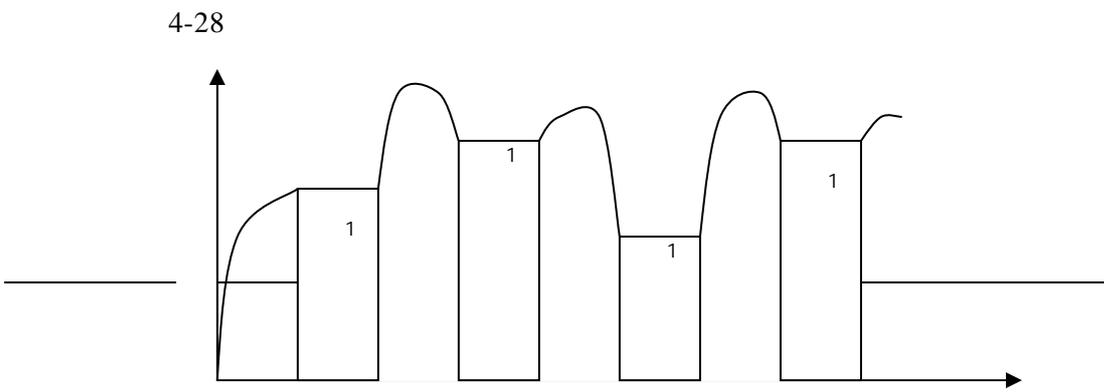
Denning

set

"

"

"

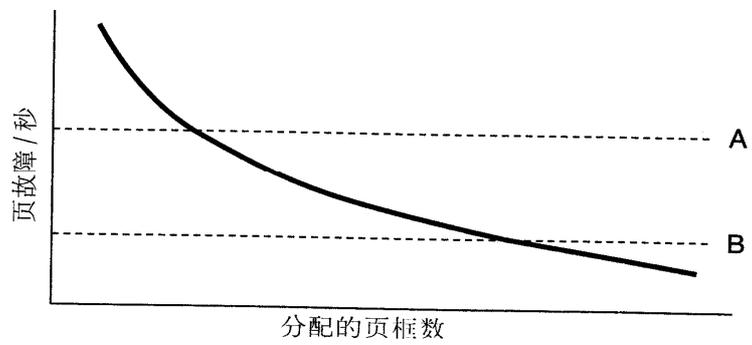


t t- t
W t " " "
W t W t

	2	3	4	5
24	24	24	24	24
15	15 24	15 24	15 24	15 24
18	18 15	18 15 24	18 15 24	18 15 24
23	23 18	23 18 15	23 18 15 24	23 18 15 24
24	24 23	24 23 18	*	*
17	17 24	17 24 23	17 24 23 18	17 24 23 18 15
18	18 17	18 17 24	*	*
24	24 18	*	*	*
18	18 24	*	*	*
17	17 18	*	*	*
17	17	*	*	*
15	15 17	15 17 18	15 17 18 24	*
24	24 15	24 15 17	*	*
17	17 24	*	*	*
24	24 17	*	*	*
18	18 24	18 24 17	*	*

W 4-29 W t

W 4-29 W t *



4-30

Page Fault Frequency PFF

LRU A

4-30 B

 PFF

A

(3)

UNIX/Linux

4.5.3

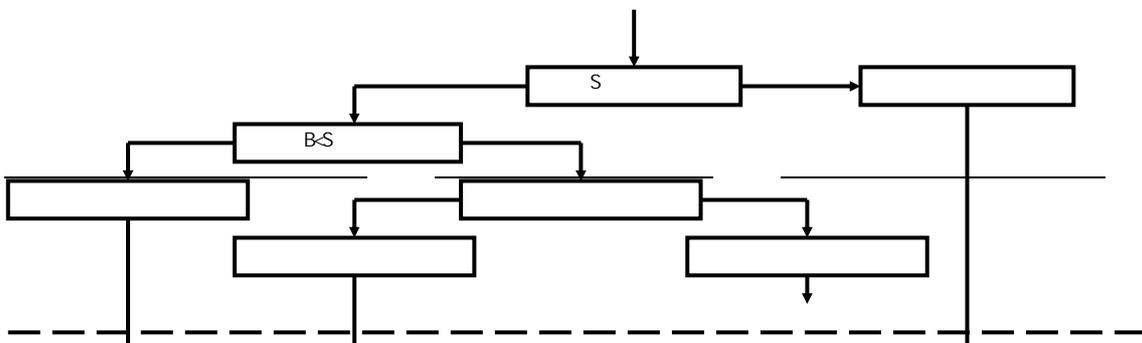
4-32

4-32

- 00 01 11
- 00 01 11
- 0 1
- 00 01 11
-
-

" "

4-33



1 " " 1
1
0 1 1

4.5.4

1
2

3

4

5 d' p s d d' $d'=p \times$ $+d$

6

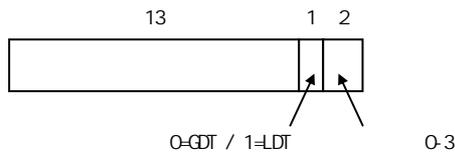
p' d s p s $d'8$ p
 s s' p s

8086 CPU DOS
 386 CPU
 Windows 8086 Linux

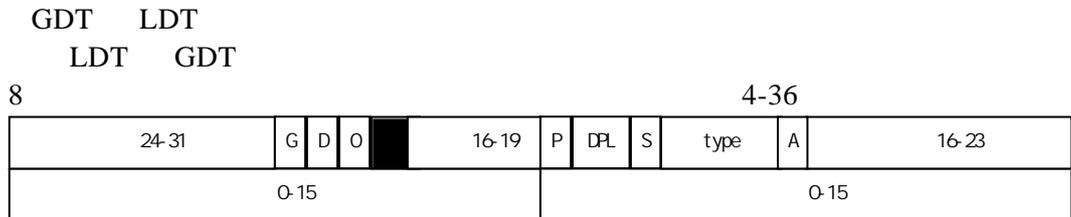
4.6.1 Intel x86/Pentium

Intel x86/Pentium descriptor table
 LDT local
 GDT global descriptor table
 LDT
 GDT

Pentium selector 6 DS
 CS 4-35 DS
 16 16k GDT 8192 LDT GDT
 LDT 8192 LDT GDT LDT



4-35 Pentium



4-36 Pentium

- 32 286 24-31 32 286
- 24 20
- G G=1 Pentium G=0 4KB
- D D=1 2²⁰ 2²⁰ 4KB=2³² D=0 16
- P P=1 P=0

- Dpl 2 2 0—3 0 1
2 3 Windows 95 0 3
- S S=1
S=0
- type 3
- A

4.6.2 Intel x86/Pentium

Intel 286 MSW 1
0
386 CR0 0 PE PE 1
PE 0 CR0
31 PG PG 1 CR3
32 PG 0
32
(1) CR0 PE =1 PG =1
8192 2²⁰ 4KB 2³²
(2) Intel " "

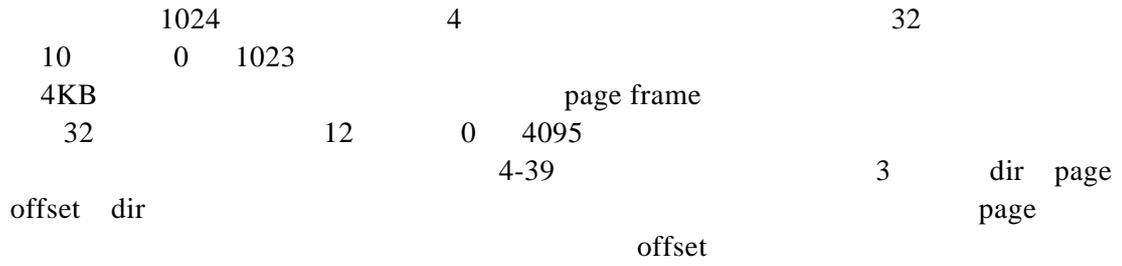
•

LDT GDT

•

P

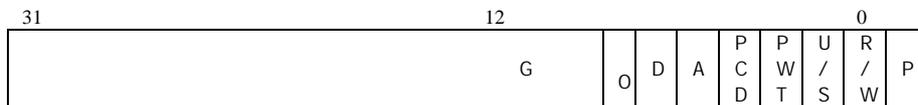
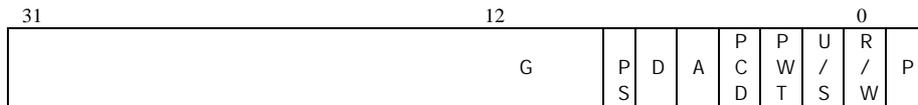
0



dir(10)

(32)

4-40 32 20
 /



4-40

- G
- D

0

- A / accessed
- A=0

- PCD cache
 - PWT cache
 - P P=1
 - U/S / User/Supervisor U/S=0
 - R/W / Read/Write R/W=1 R/W=0
- OS 1024 4K 4M

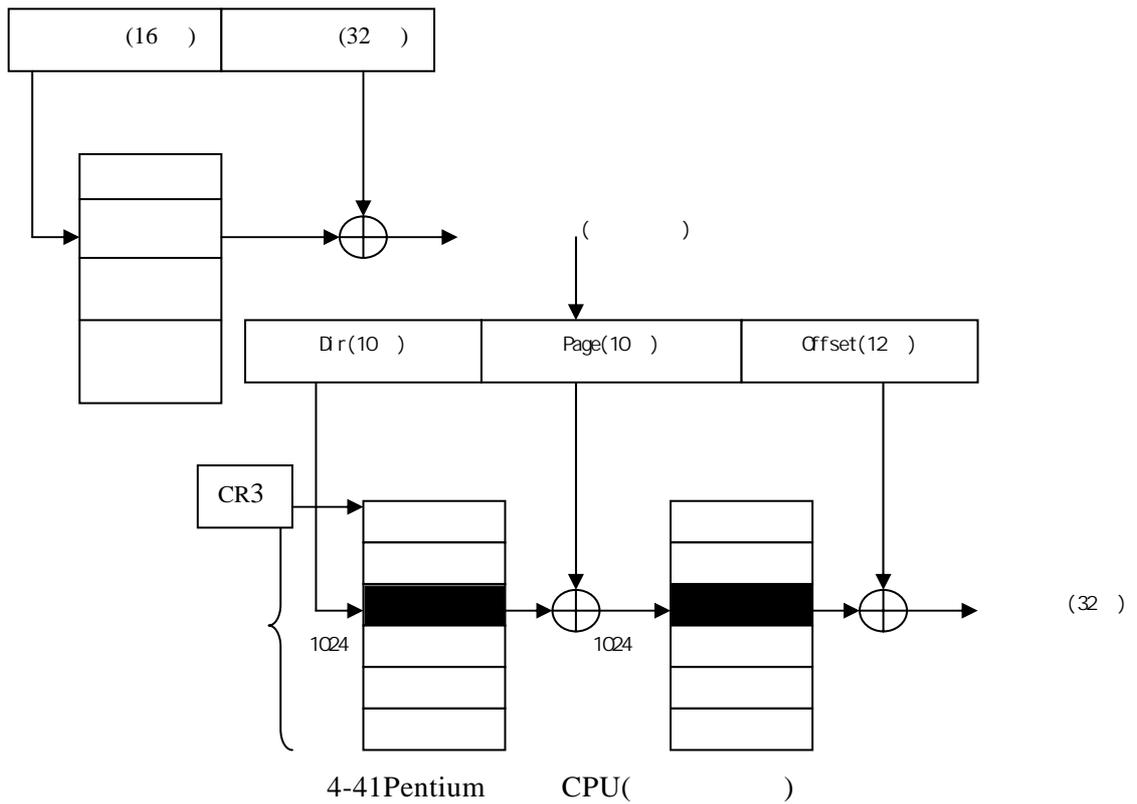
Intel x86/Pentium
4-41

Intel x86/Pentium
dir/page

TLB

Intel x86/Pentium

80286



4.7 Windows 2000/XP

4.7.1

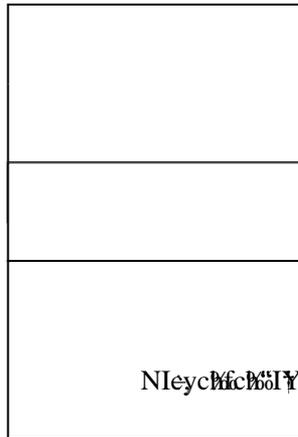
Windows 2000/XP
Memory Manager

Ntoskrnl.exe

VMM Virtual
Windows

Windows 2000/XP " " 386
 32 4GB 2³²B 4 42
 2GB 2GB Windows
 3GB 1GB

FFFFFFFFH



{550c14c1-4000-4293-a959-964182c22341}

4.7.2

Windows 2000/XP

-
-
-

1

Windows 2000/XP

4GB

"

" VAD Virtual Address Descriptor

VAD

VAD

Self-balancing Binary Tree

4 43

Windows 2000/XP
VAD

VAD VAD

2

" " section object Win32 " "

-
-
-

4 44

-

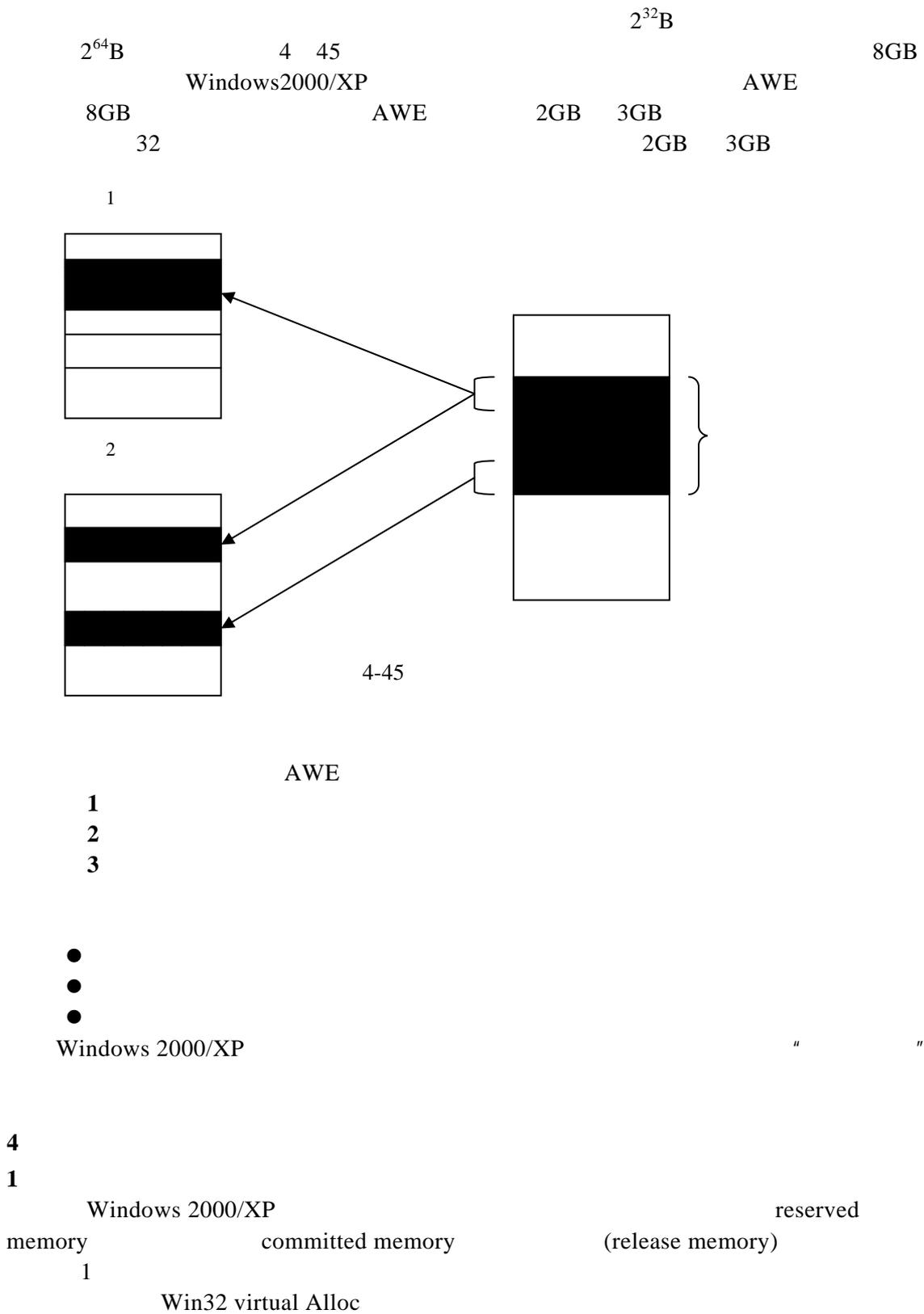


4 44

3

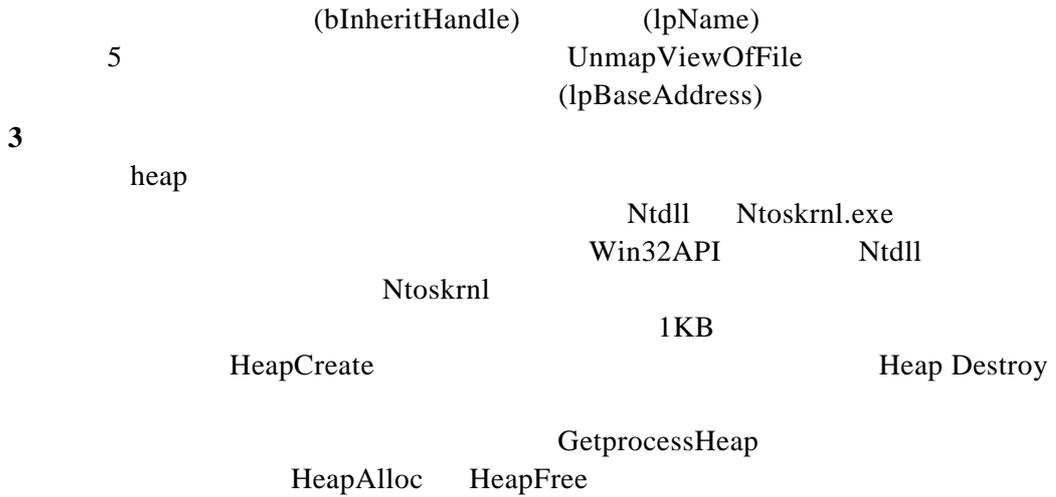
2^{32}B Windows 2000/XP 2^{64}B AWE Address

Windowing Extension

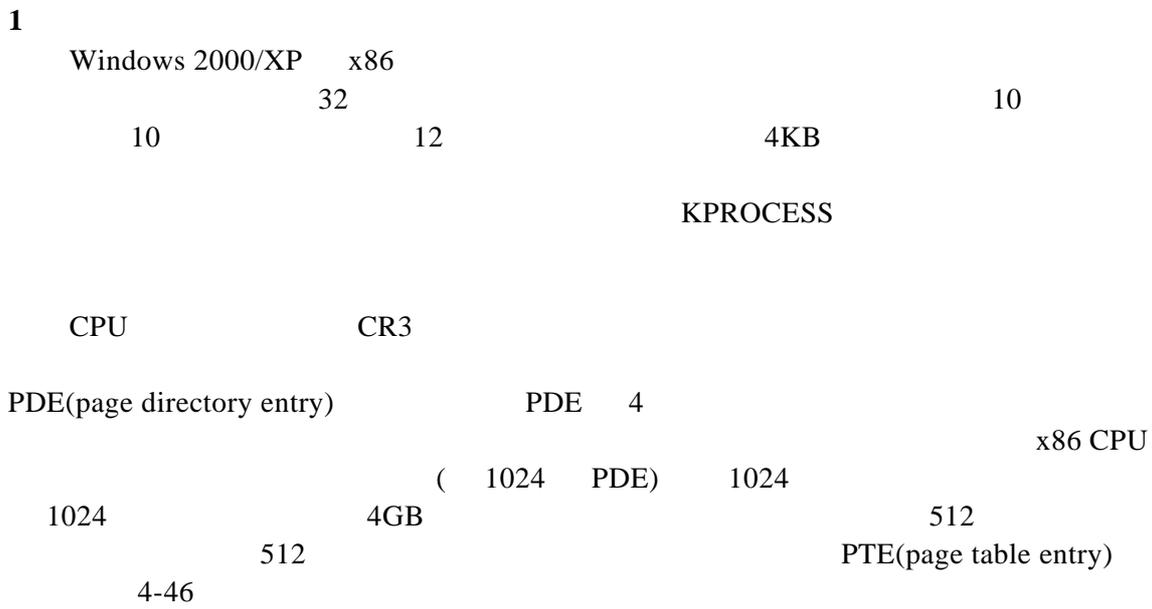


IpAddress

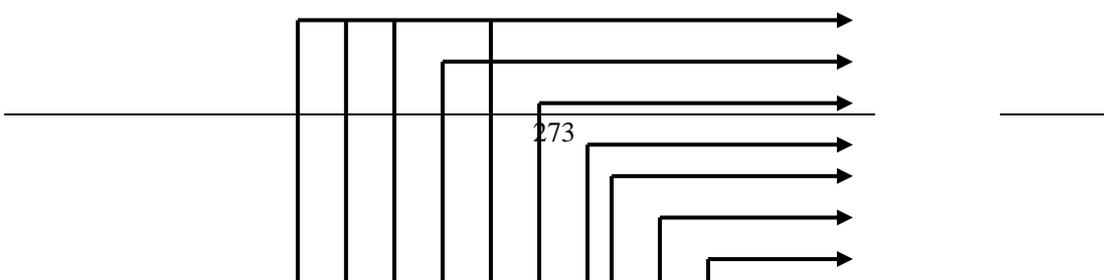
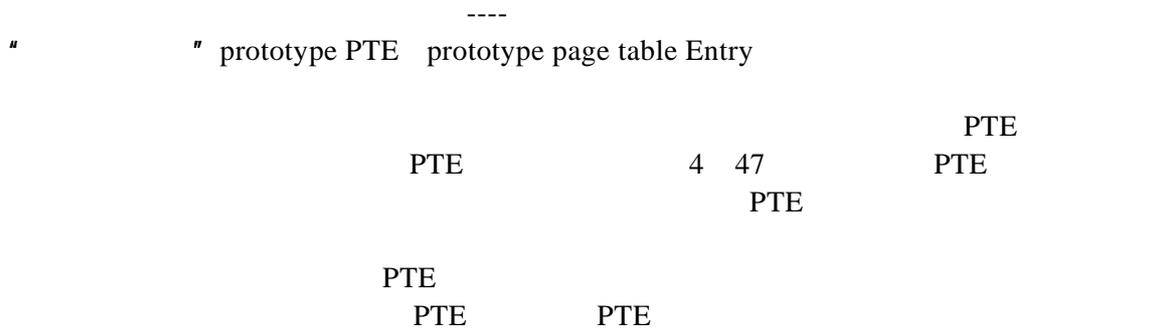
NULL



4.7.3



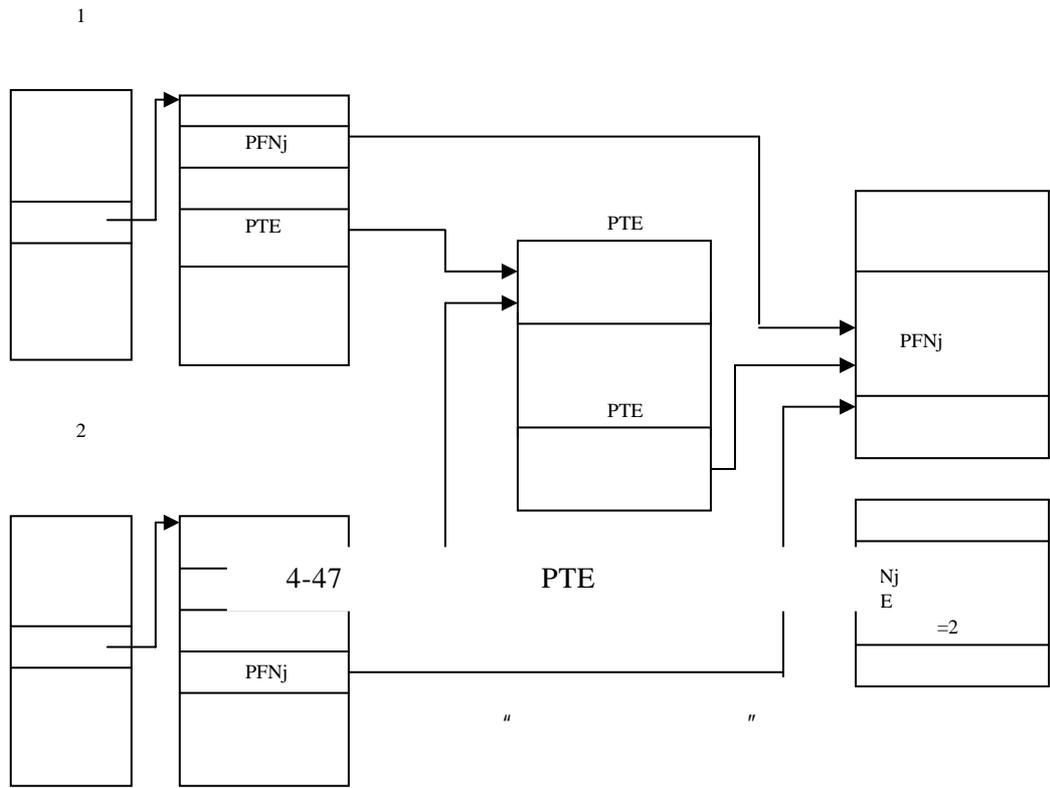
2



PTE 32

PTE

PTE

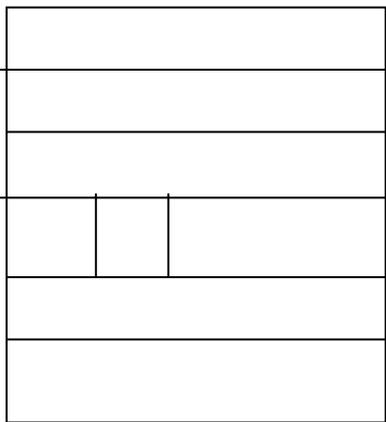


- PTE active/valid
- transition

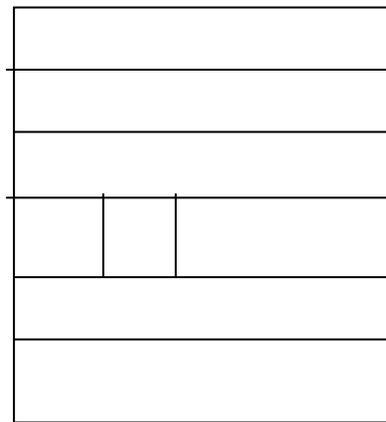
- modified-no-write
- demand zeropage
- page file
- mapped file

3

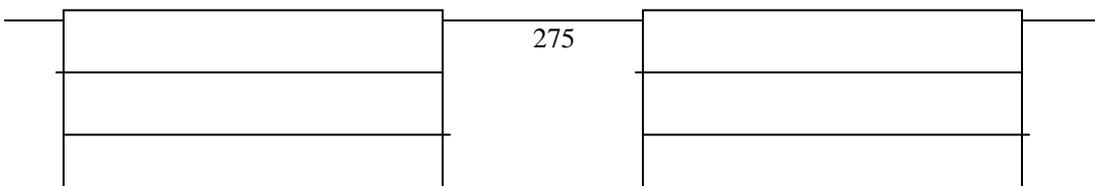
Database	PFN	Page Frame Number
Windows 2000/XP		
4-48		
●		
●		0
● PNF		
●		PTE



PFN



PFN

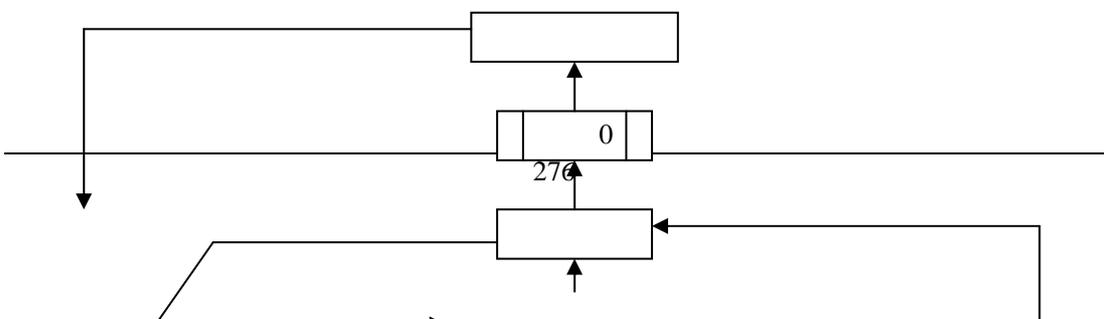


PFN		8	
1	Valid		
2	transition		
		I O	
3	Stand by		
4	Modified		
5	modified no write		
6	free		(0)
7	zeroed		
8	bad		

4 49

4

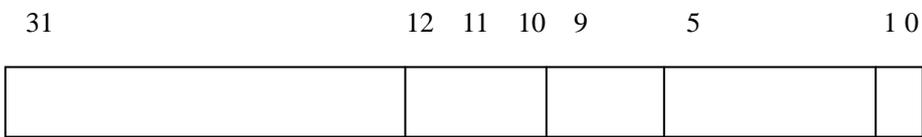
| | > > >



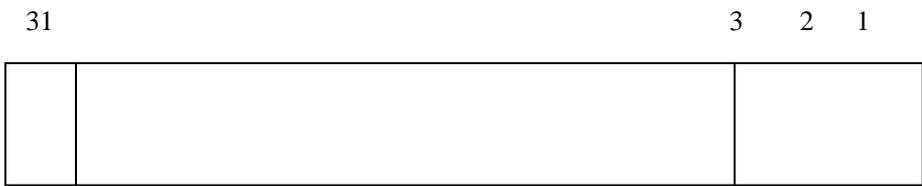
	MmMinimum Freepages	"	"
	"	"	
	PTE		
	PTE		
5	Windows 2000/XP		
		19MB	I O
	4		8
	Windows 2000/XP		FIFO
			FIFO
	Windows 2000/XP	"	
	"		
	64MB	50	32MB Windows 2000/XP server
			345

1
 Windows 2000/XP
 Win32 Set process working
 Set

6
 Windows 2000/XP 16
 4-50 a 4-50 b



(a)



b

4-50

7
 1

-
- Win32

2

allocation granularity
 64KB

x86 4KB 18KB
 20KB

3

Windows 2000/XP
 4

-

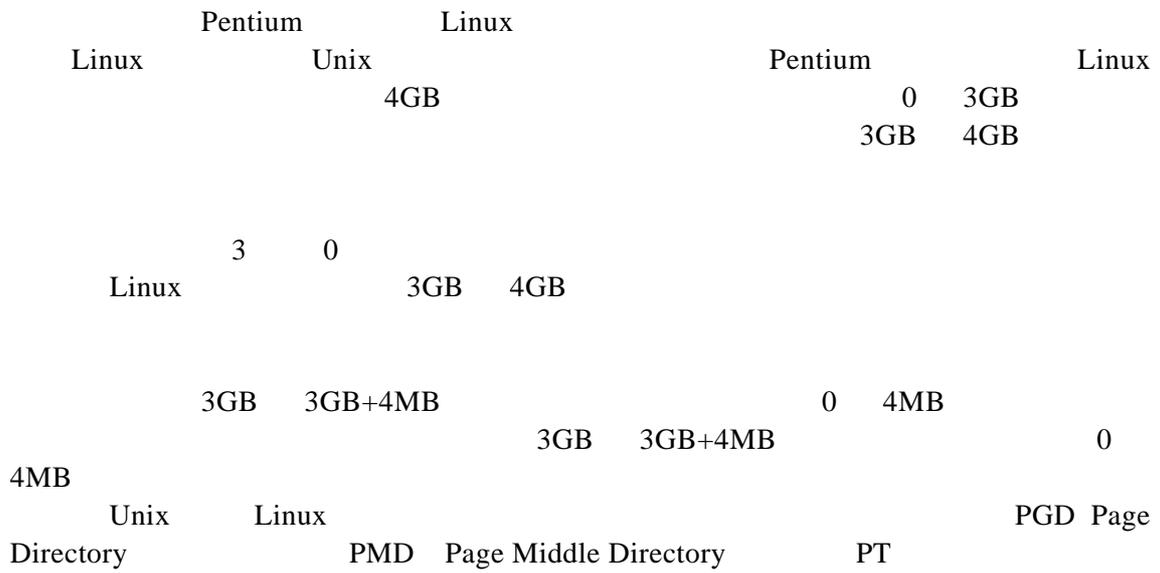
-
-
-

Access Control List

ACL
ACL

4.8 Linux

4.8.1 Linux



4GB

4MB

Linux

vma virtual memory

area

vma

vma

vm_area_struct

```

struct vm_area_struct{
  struct mm_struct *vm_mm      /* */
  unsigned long vm_start      /* */
  unsigned long vm_end        /* */
  pgprot_t vm_page_prot
  unsigned short vm_flags
  /*          vma  AVL      *
  shrot vm_avl_height
  struct vm_area_struct *vm_avl_left
  struct vm_area_struct *vm_avl_right
  /*          vma          *
  struct vm_area_struct *vm_next
  /*                                          *
  struct vm_area_struct *vm_next_share
  struct vm_area_struct *vm_prev_share
  /*          *
  struct vm_operations_struct vm_ops      *      open  close*
  unsigned long vm_offset      /*          *
  struct inode *vm_inode      /*      inode  NULL*
  unsigned long vm_pte
}

```

vma

Linux vma

PCB task_struct
mm_struct

mm_struct

vma

mmap

vm_next

vma

pgd

Linux

" "

4.8.3

linux

men-map

free-area-init()

men-map

men-map-t

men-map-t

```

typedef struct page{
  struct page *next *prev      /* */
  struct inode *inode
  unsigned long offset

```

```

/*
    struct page *next_hash /*page hash */
    atomic_t count /* */
    unsigned flags /* */
    unsigned dirty /* */
    unsigned age /* */
    struct wait_queue *wait
    struct page *prev_hash /*pager hash */
    struct buffer_head *buffers /* */
    unsigned long swap_unlock_entry
    unsigned long map_nr /* mem_map */
}mem_map_t

mem_map bitmap
free_area_init()
NR_MEM_LISTS 6
end_mem-start_mem/PAGE_SIZE/20+3
20
end_mem-start_mem/PAGE_SIZE/21+3
21 1 2
i end_mem-start_mem/PAGE_SIZE/2i+3
2i 1 1
Linux free_area NR_MEM_LISTS
free_area_struct
struct free_area_struct{
    struct page *next *prev /* next prev struct page */
    unsigned int *map /* bitmap */
}
static struct free_area_struct free_area[NR_MEM_LISTS]
bitmap free_area
0 2i free-area i Linux buddy
2i 0 i<NR_MEM_LISTS
_get_free_pages() free_pages()
2i free_area i
i+1
free_area
free_area
bitmap i NR_MEM_LISTS
bit 1
change_bit() bitmap
bitmap free_area
free_area

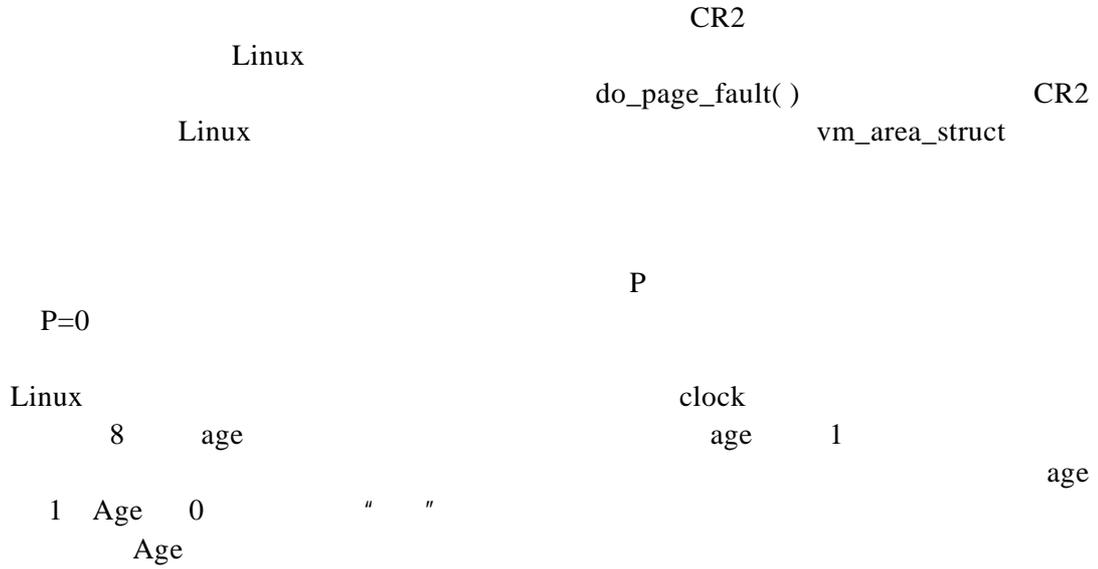
```

4.8.4

vmalloc() vfree()

```
3GB      3GB  high_memory  HOLE_8M      vmlist
          high_memory
HOLE_8M  3GB  high_memory
          8MB  "      "      vmlist
          0   3GB      :
```

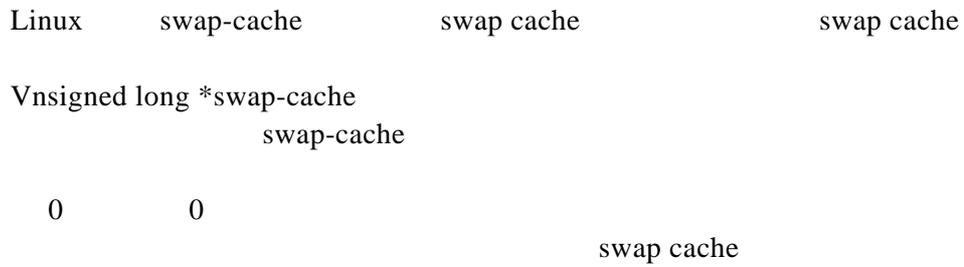

3



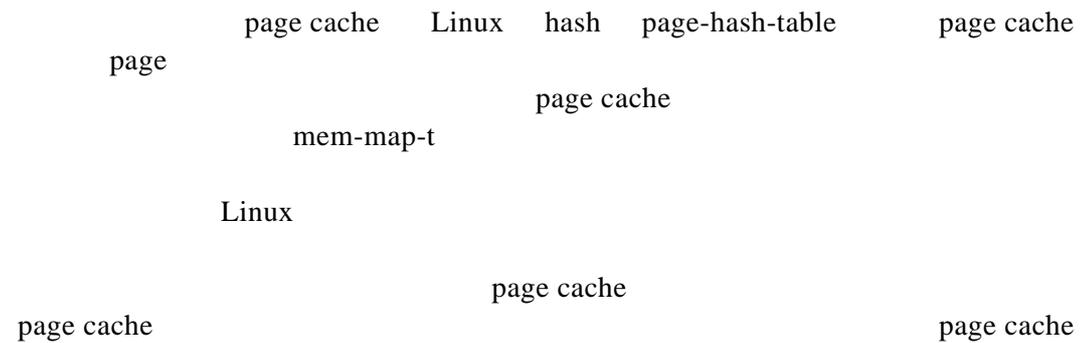
4.8.7

Linux kmalloc cache swap cache page cache

1 swap cache



2 page cache



3 kmalloc cache

kmalloc() kfree()
kmalloc cache

4.9

,

Optimal LRU
LRU

Clock LFU

Optimal

-
-
-
-

()

,

- 1
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 13
- 14
- 15
- 16

- 17
- 18
- 19
- 20

?

?

/

/

25 (1) (2)
 26 " "
 27
 28
 29
 30
 31
 32
 33
 34 FIFO LRU ()
 35 CPU20% 97.7%
 36 50% CPU (1)
 37 (2) (3) goto (4) (5)
 ?
 1
 1 2 3 4 2 1 5 6 2 1 2 3 7 6 3 2 1 2 3 6
 FIFO OPT LRU 3 4 5 6
 2 5
 (1)1 4 3 1 2 5 1 4 2 1 4 5
 (2)3 2 1 4 4 5 5 3 4 3 2 1 5
 FIFO LRU
 3 FIFO OPT LRU
 (1)2 3 2 1 5 2 4 5 3 2 5 2
 (2)4 3 2 1 4 3 5 4 3 2 1 5
 (3)1 2 3 4 1 2 5 1 2 3 4 5
 3 4
 4 10K 4K 20K 18K 7K 9K
 12K 15K (1)12K 10K 9K (2)12K 10K 15K 18K

5
 212K 417K 112K 426K (1) 100K 500K 200K 300K 600K
 first-fit best-fit worst-fit
 ?(2) ?

6 32 9 11

7 5 A B C D A B E A B C D E
 FIFO 3 4

8 Ans

Cns Bns n-1

/n m-1 /m cache cache

9 20ns cache 60ns cache

cache 12ms 0.9 0.6 60ns cache (ns)

10 (1) 1.2 80%
 (2)

11

0	219	600
1	2300	14
2	90	100
3	1327	580
4	1952	96

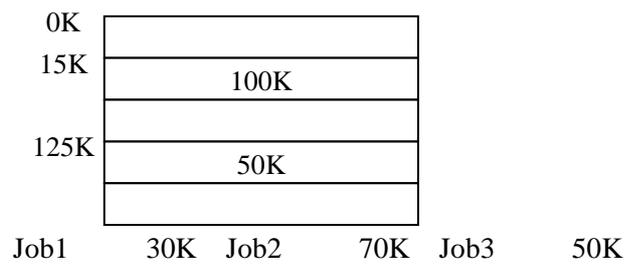
1 [0 430] 2 [3 400] 3 [1 1] 4 [2 500] 5 [4

42]

12 24 2¹⁸B 100(

1KB)

13



14 8 16 2048

15 2F6AH 0 1 2 16 4096 15 18 21

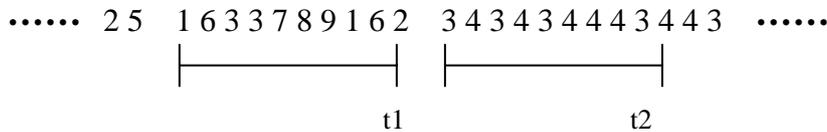
16 VAR A ARRAY[1 100,1 100] OF integer LRU 3 200

```
1
A
FOR i 1 TO 100 DO
  FOR j 1 TO 100 DO
    A[i,j] 0
  B
  FOR j 1 TO 100 DO
    FOR i 1 TO 100 DO
      A[i,j] 0
    A B
```

17 48 32 8KB

? ?

18



19 9 t1 t2 CPU (1)CPU 13% 97% (2)CPU 87% 3% (3)CPU 13% 3% 20 10 0 1 2 3 32 1KB 16KB 8 7 4 10

0AC5H 1AC5H 21 4 R D

()

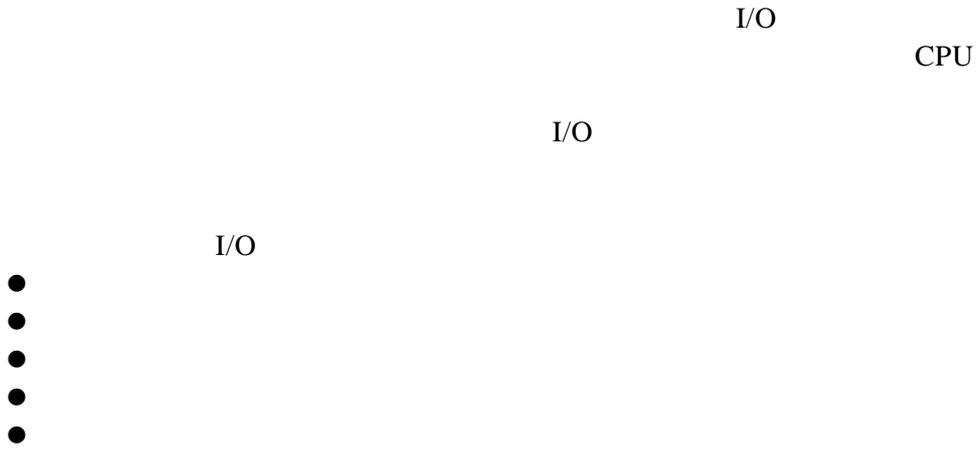
page	loaded	last ref	R	D
0	126	279	0	0
1	230	260	1	0
2	120	272	1	1
3	160	280	1	1

21 FIFO LRU CLOCK ?

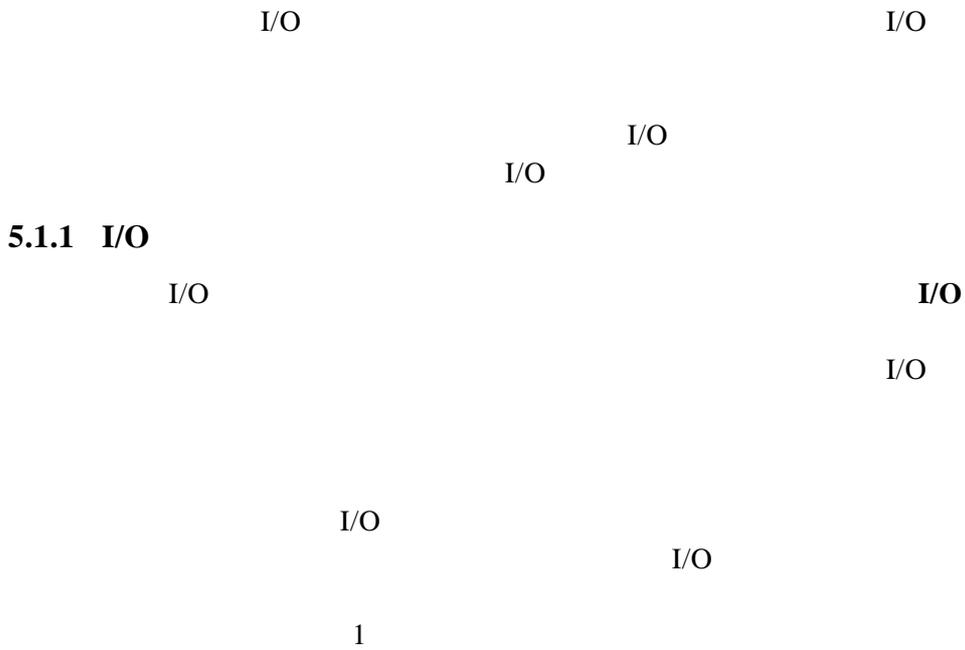
22

```
for (i=0;i<20,i++)
```


CH5



5.1 I/O

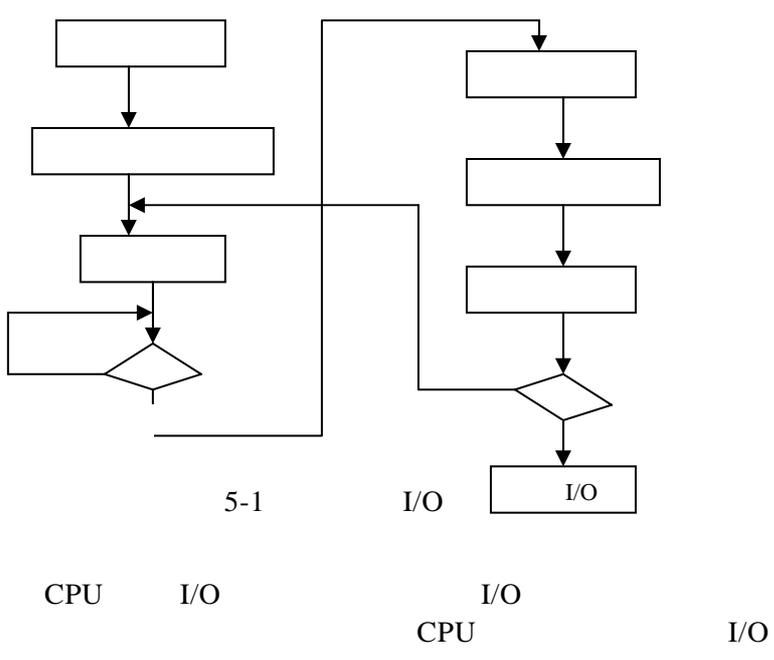
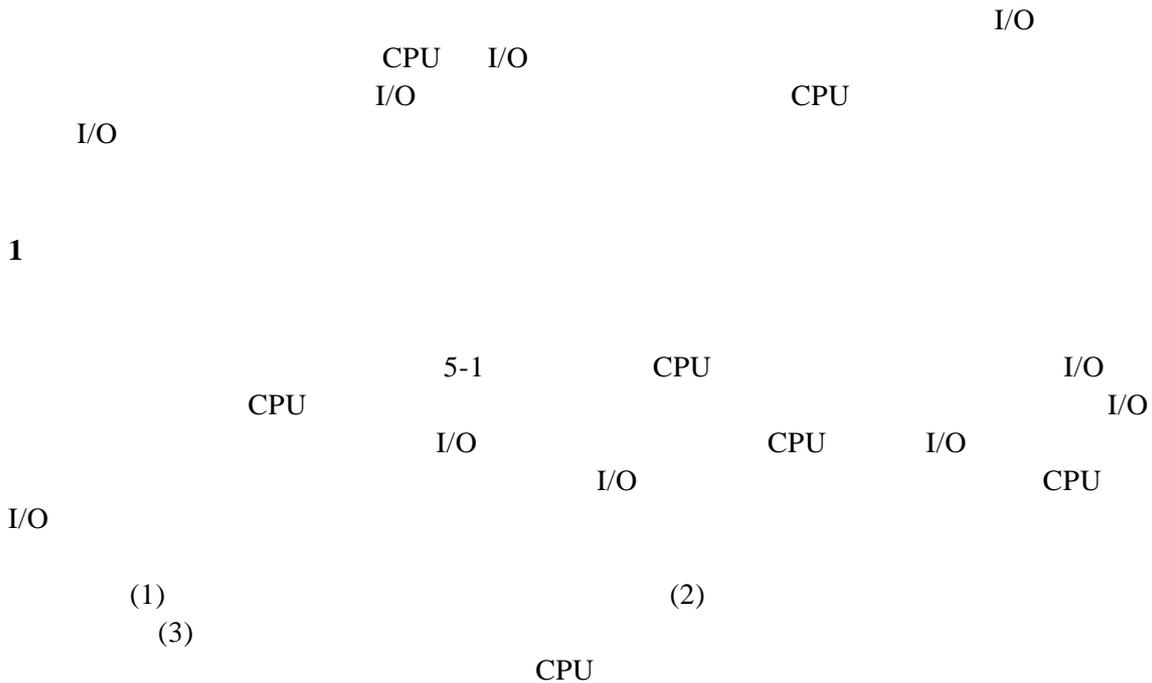


-
-
-
-

KB

I/O

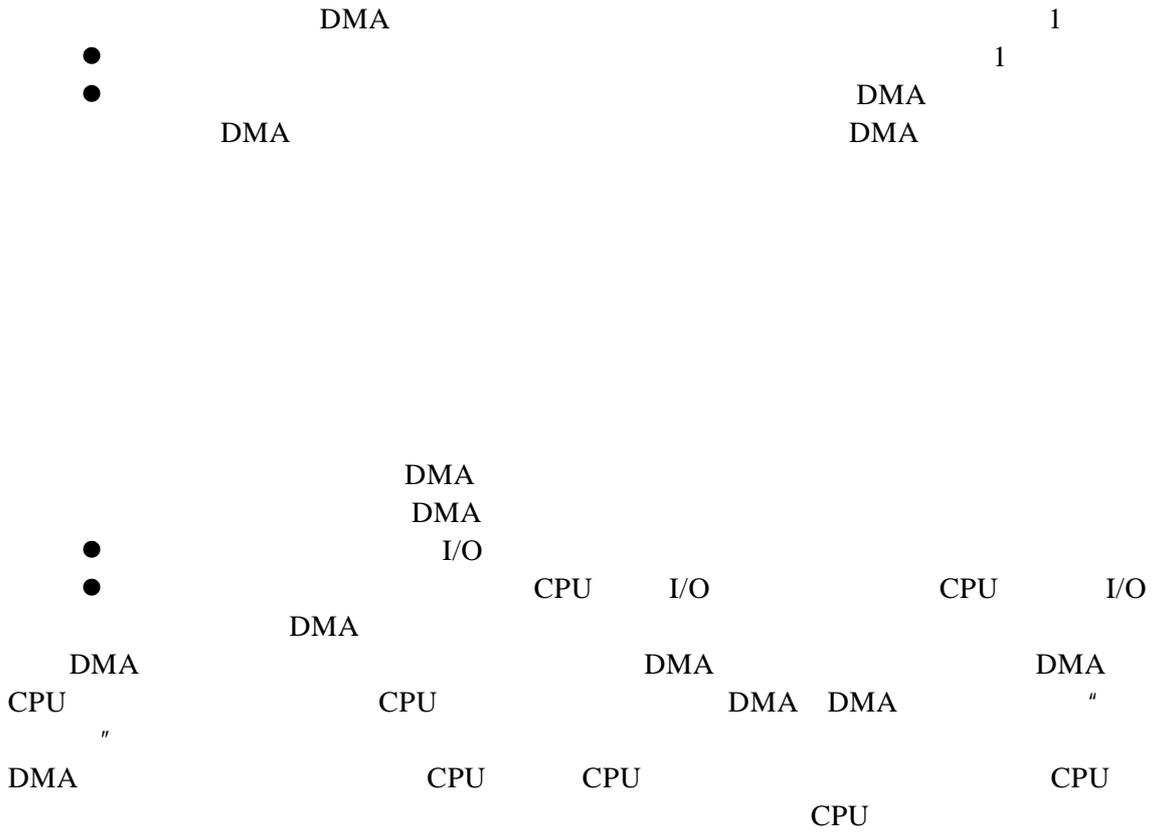
5.1.2 I/O



CPU

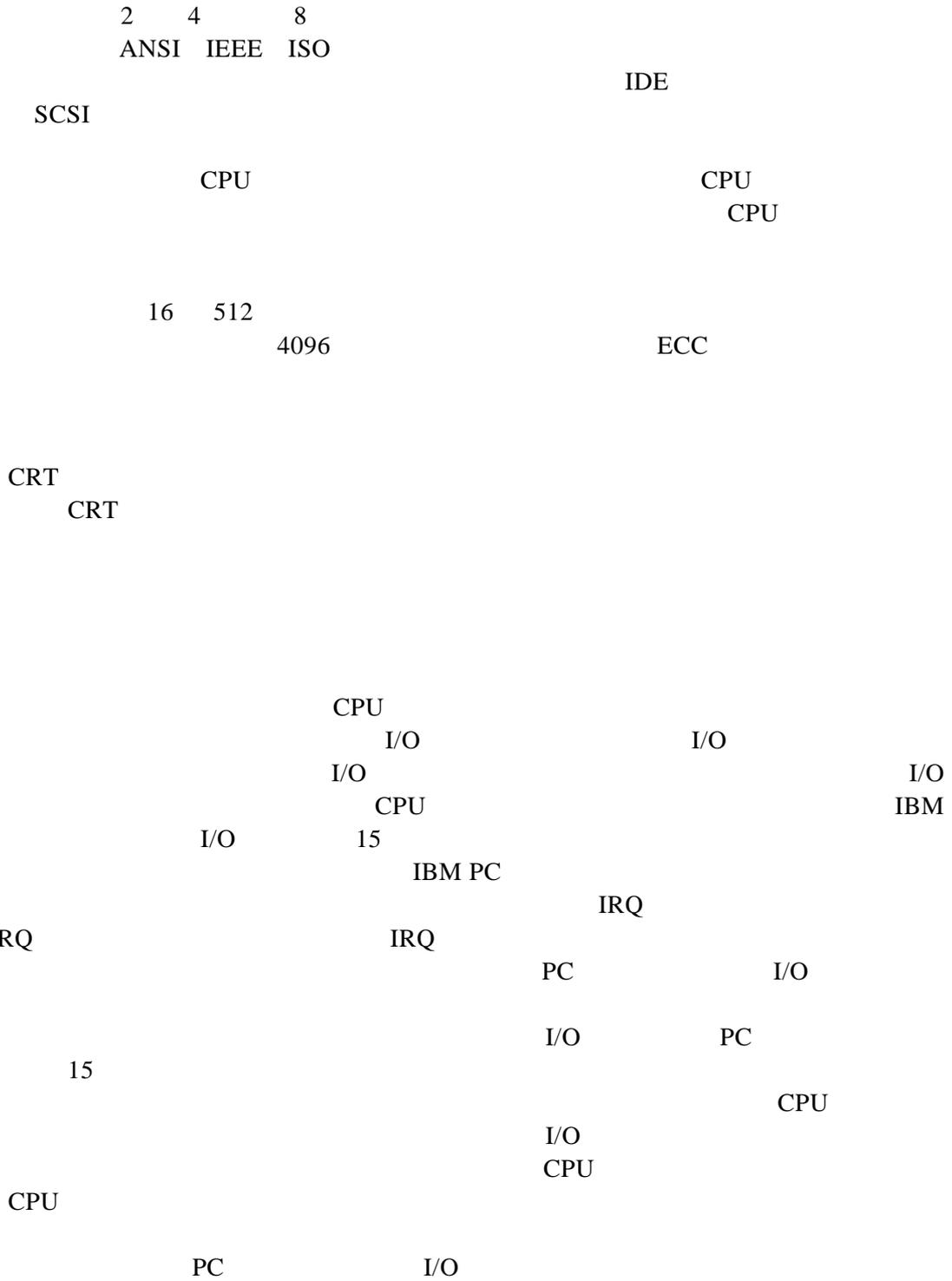
CPU

CPU

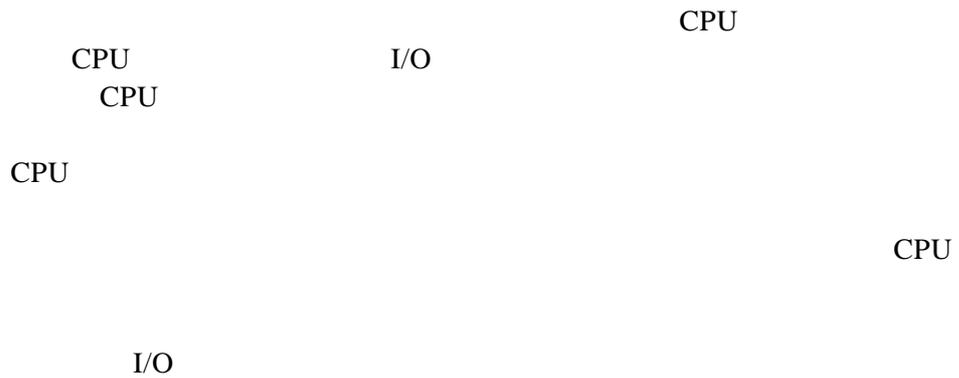


5.1.3

I/O

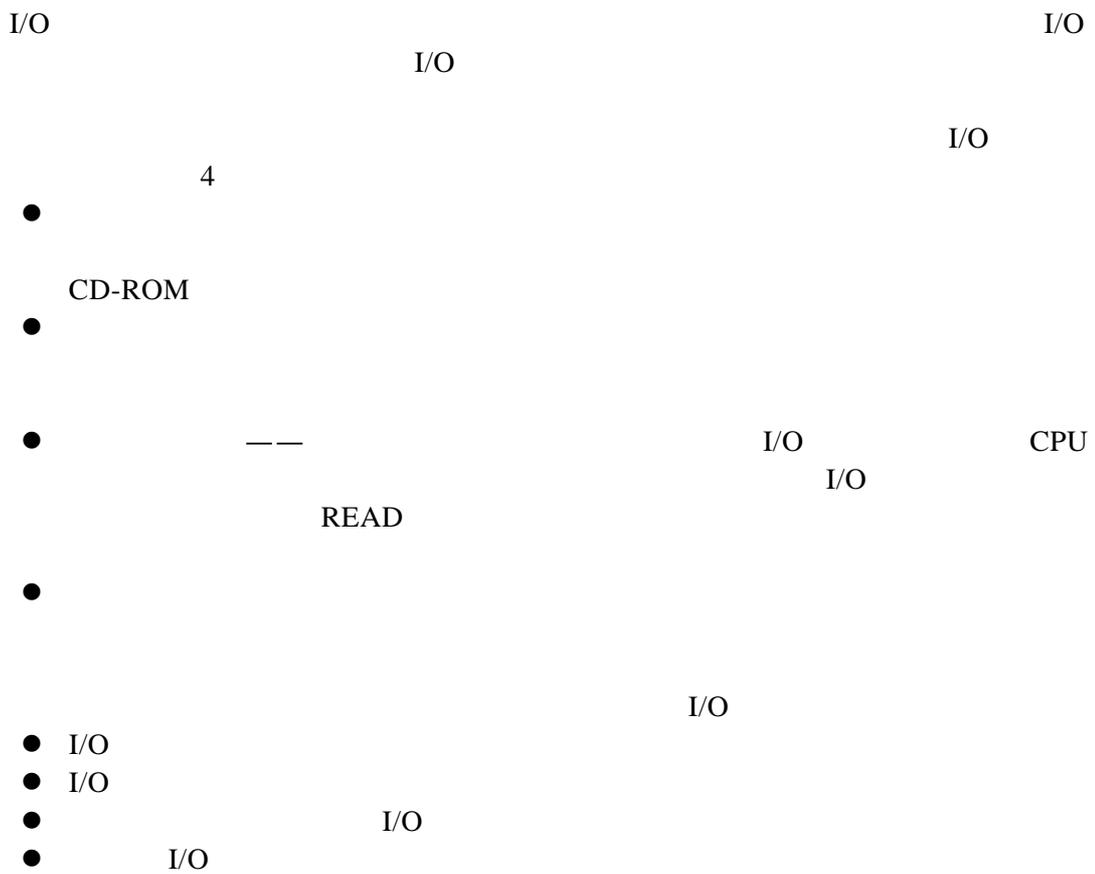


I/O	I/O		
	040 - 043	0	8
	060 - 063	1	9
	1F0—1F7	14	118
	3F0-3F7	6	14
LPT1	378—37F	7	15
COM1	3F8—3FF	4	12
COM2	2F8-2FF	3	11

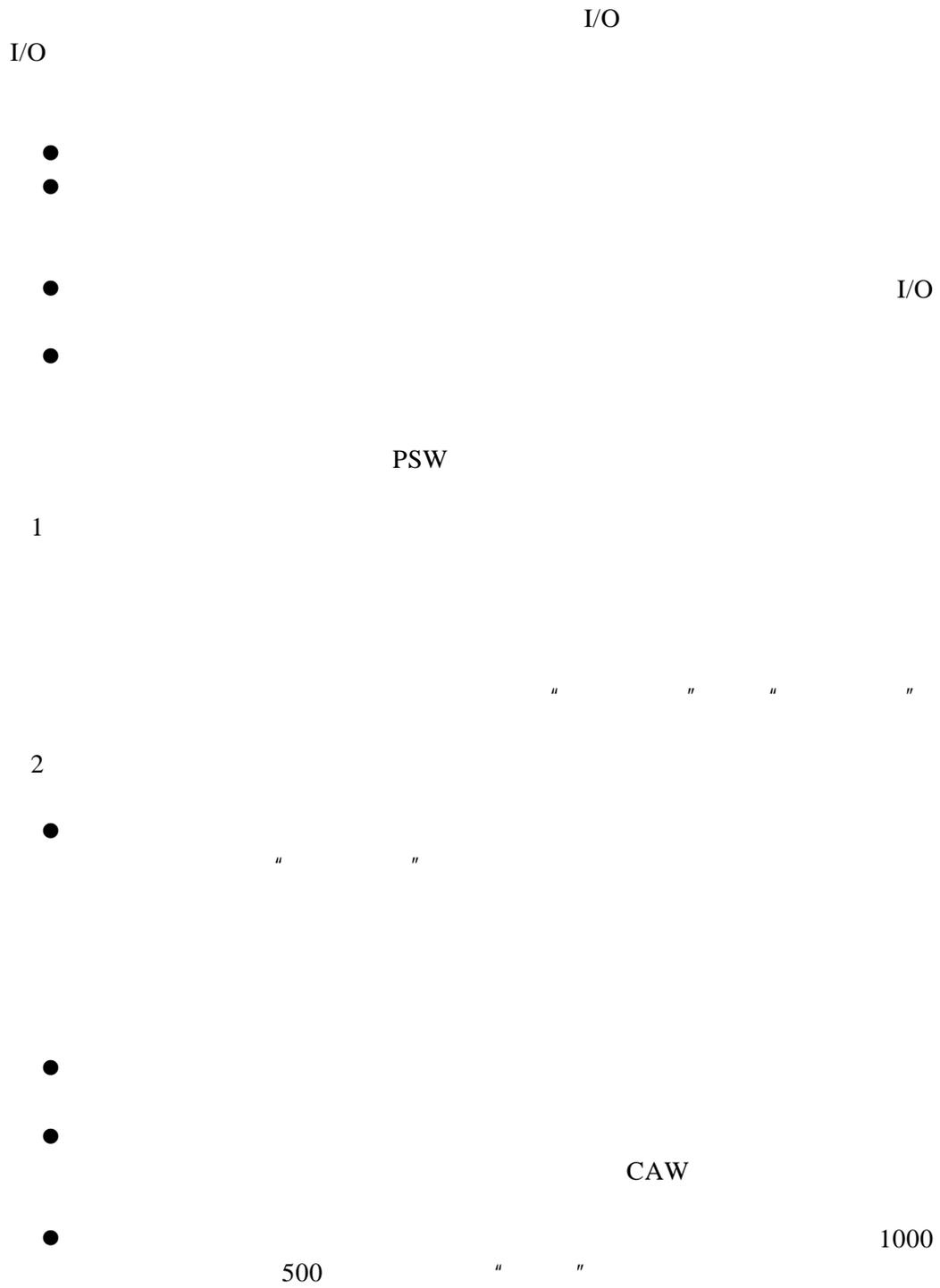


5.2 I/O

5.2.1 I/O



5.2.2 I/O



3

"

"

"

"

4

"

"

5.2.3

I/O

I/O

n

I/O

5.2.4

I/O

I/O



(1)I/O

/dec/tty00
number
minor device number

i-

I/O

UNIX
major device

i-

UNIX
rwx
(2)

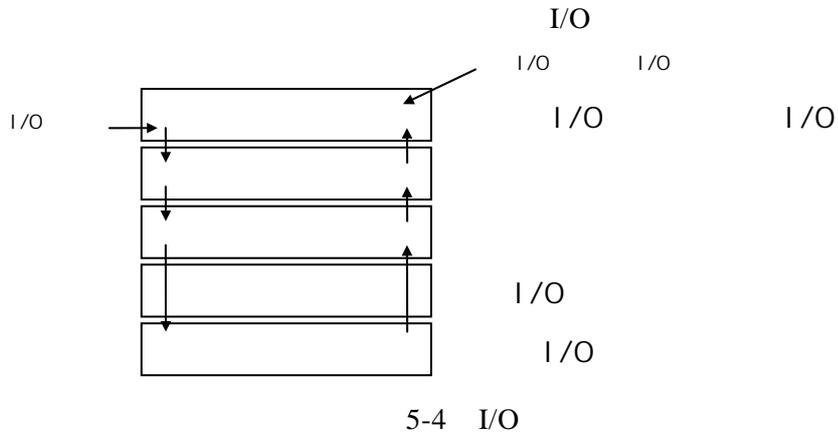
I/O

I/O

5.2.5 I/O

(1)
 I/O
 () I/O
 C
 count = write fd buffer nbytes
 write
 I/O

C printf
 write ASCII
 scanf printf I/O I/O
 2 spooling
 I/O I/O spooling
 5-4 I/O



5.3 I/O

CPU CPU I/O

5.3.1

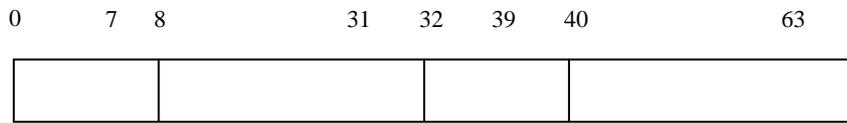
1
 I/O I/O
CCW Channel Command Word I/O
 I/O

I/O

I/O

IBM370

5-5



5-5 IBM370

●

●

" "

●

32 36

32 33 0 01

32 1

34 1

35 1 36 1

●

0

2

CCW	X'02'	inarea	X'40'	80
CCW	X'02'	*	X'50'	80
CCW	X'02'	inarea +80	X'40'	80
CCW	X'02'	*	X'50'	80
CCW	X'02'	inarea +160	X'40'	80

inarea DS CL240

*

X'50' " "

3

I/O	CAW	Channel
Address Word	CSW	Channel Status Word

8

I/O IBM

-
-
-
-

8

5.3.2 I/O I/O

IBM	I/O	I/O	I/O	I/O	I/O	I/O	I/O	Start
I/O SIO	I/O	Test I/O	TIO	Test Channel	TCH	I/O	I/O	I/O
Halt I/O	HIO		Halt Device	HDV				
	I/O	SIO X'00E'						
	0	0E	CPU I/O			I/O		
SIO			0			1	CAW	
		CPU						
CSW		2						
	3							
	I/O						I/O	CPU
				CAW				
	CPU	I/O						I/O
●	I/O							
●								
●			I/O					
●		I/O						
		IBM						

START

```

BALR 11 0
USING * 11
SSM = X'00' /*
LA 8 READ0
ST 8 CAW
SIO X'0182' /*
BC 7 *-4 /*
TIO X'0182'
BC 7 *-4 /*
LOOP LA 8 READ1

```

```

ST 8 CAW
SIO X'0182' /* 1
BC 7 *-4
TIO X'0182'
BC 7 *-4 /*
LA 8 PRINT1
ST 8 CAW
TIO X'00E'
BC 7 *-4 /* 2
SIO X'00E' /* 1
LA 8 READ2
ST 8 CAW
SIO X'0182' /* 2
BC 7 *-4
TIO X'0182'
BC 7 *-4 /*
LA 8 PRINT2
ST 8 CAW
TIO X'00E' /*
BC 7 *-4 /* 1
SIO X'00E' /* 1
B LOOP

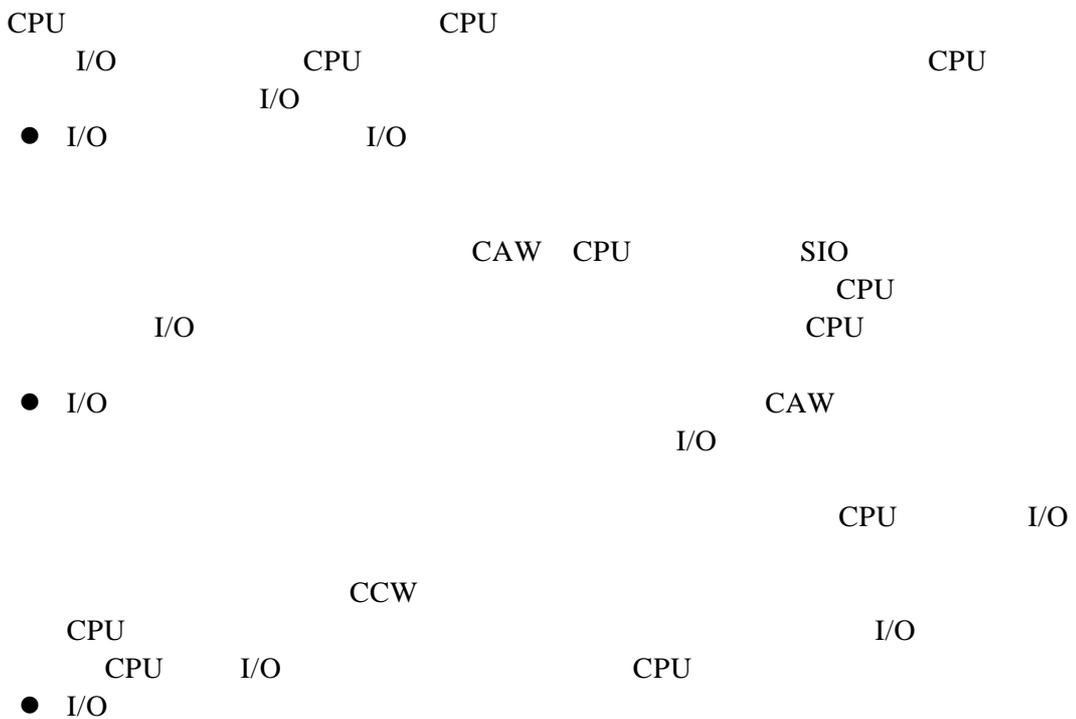
```

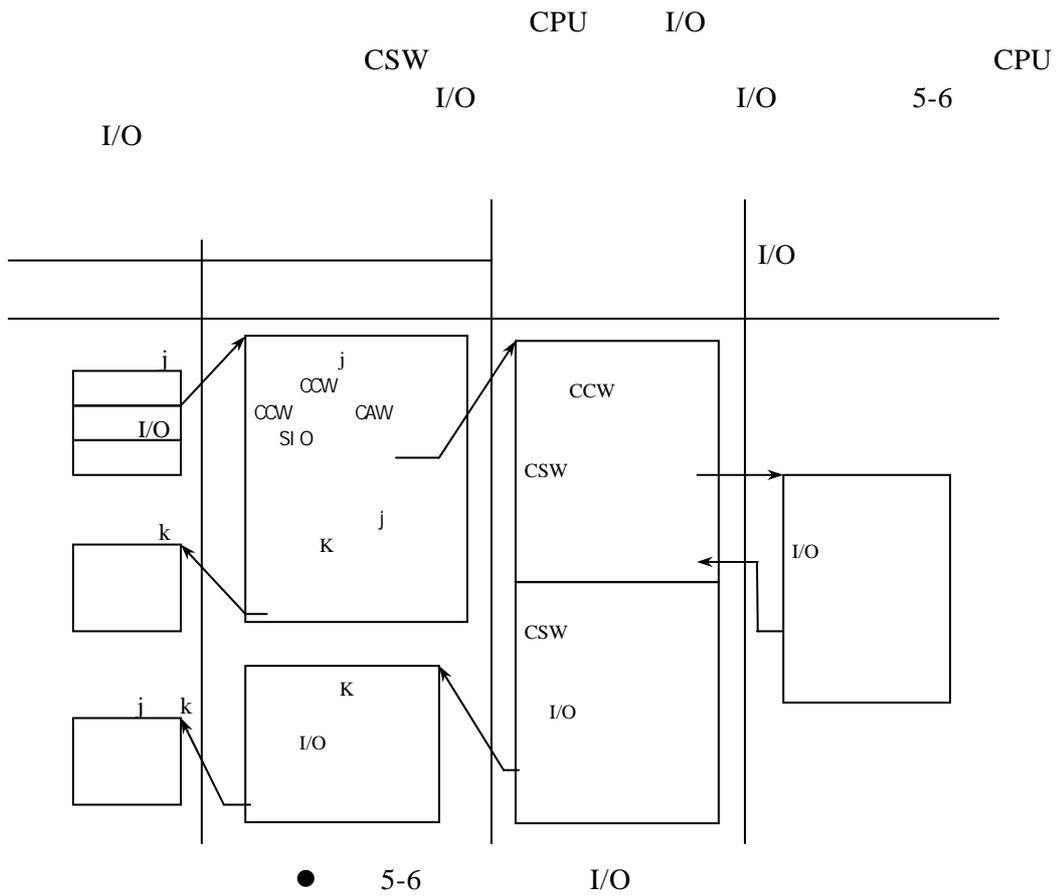
```

READ0 CCW X'07' * X'20' 1
READ1 CCW X'02' BUFFER1 X'00' 512
READ2 CCW X'02' BUFFER2 X'00' 512

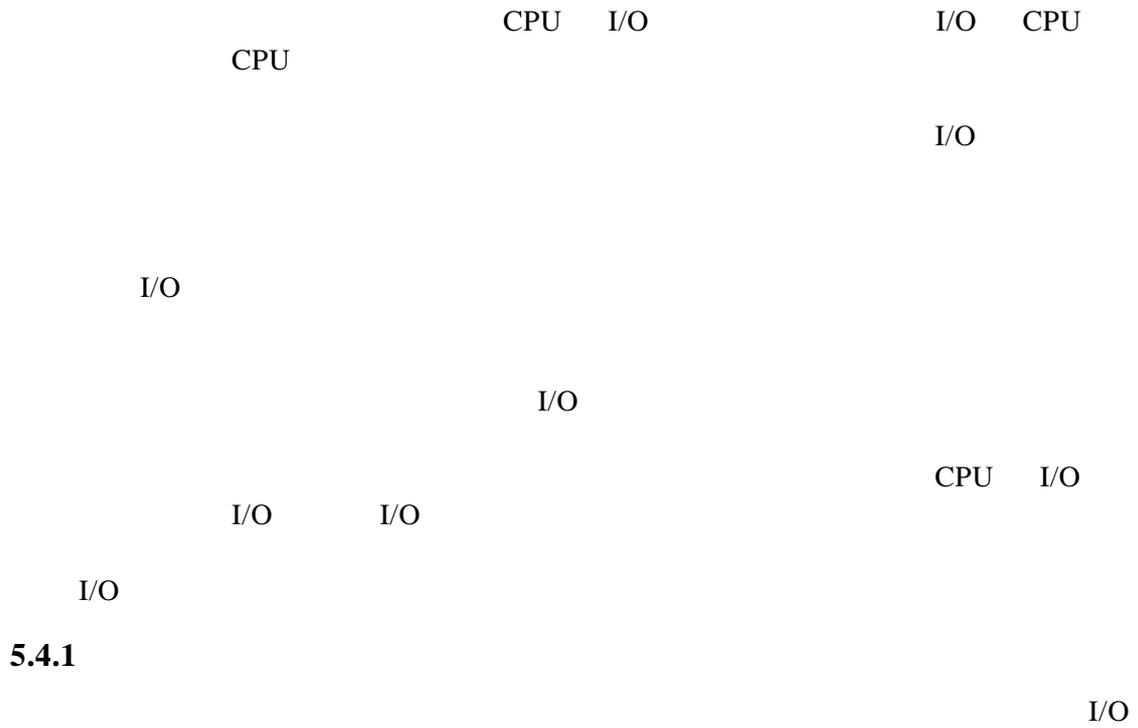
```

5.3.3 I/O





5.4



5.4.1

T

M

C

T+C

max[C

T]+M M C T

I/O

5.4.2

I/O

buffer swapping

			1	1
2	2	1	1	
2			2	

CPU I/O

UNIX 15 512

100 8

" " " " " " "

CPU

I/O

5-7 UNIX I/O

5-7 UNIX I/O

c-cf / c-cc

c-cf c-cl c-cf c-cf

0 e f c-cc c-cf c-cf c-cf

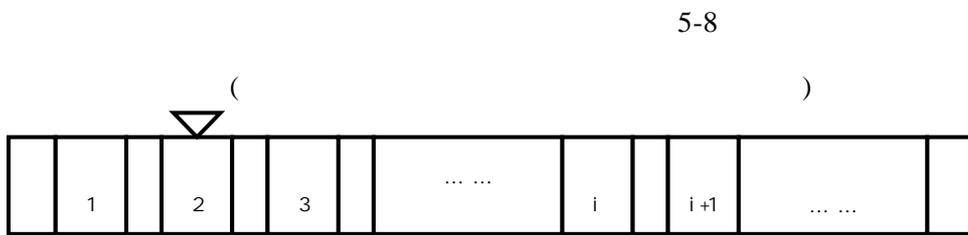
0 c-cl c-next c-cf c-cl

I/O

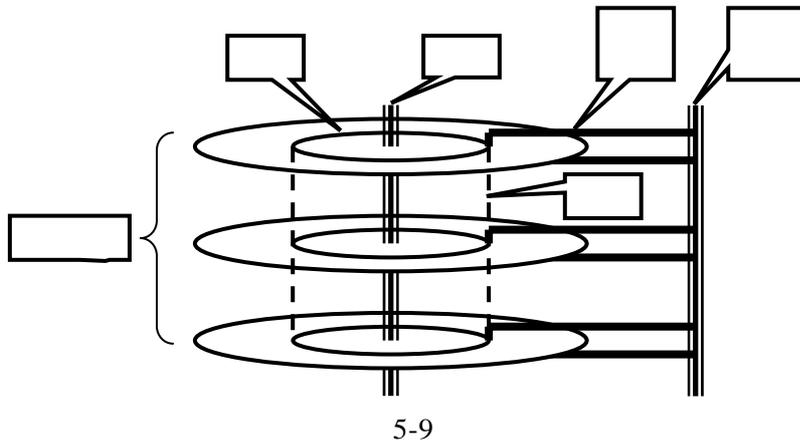
c-cl

5.5

5.5.1



5-8



5-9

" 20

"

10

" "

5.5.2

4

- | | |
|---|---|
| 1 | 4 |
| 2 | 3 |
| 3 | 2 |
| 4 | 1 |

● 1 4 3 2 1
 1/2 1/4 4 3/4
 3 1/2+1/4+3× 3/4=3 60

● 2 1 2 3 4
 ● 3 1/2+1/4+3× 1/4=1.5 30 4 1
 2 3 3 20

1	1	2
2	1	3
3	1	1
4	6	3
5	4	2
1 5	2 4	
1		5

5.5.3

B J 10 10 A

1	A
2	B
3	C
4	D
5	E
6	F
7	G
8	H
9	I
10	J

4 A 20 D
 B +2 10 A +4 10 A +9× [16 A +2
+4] 214

1	A

2	H
3	E
4	B
5	I
6	F
7	C
8	J
9	G
10	D

[2 A B A +10×
 10 10 3
 × 4]=70

5.5.4

1 A 8 20 A 1
1 5 10 " " 1 1
 n " " n

5.5.5

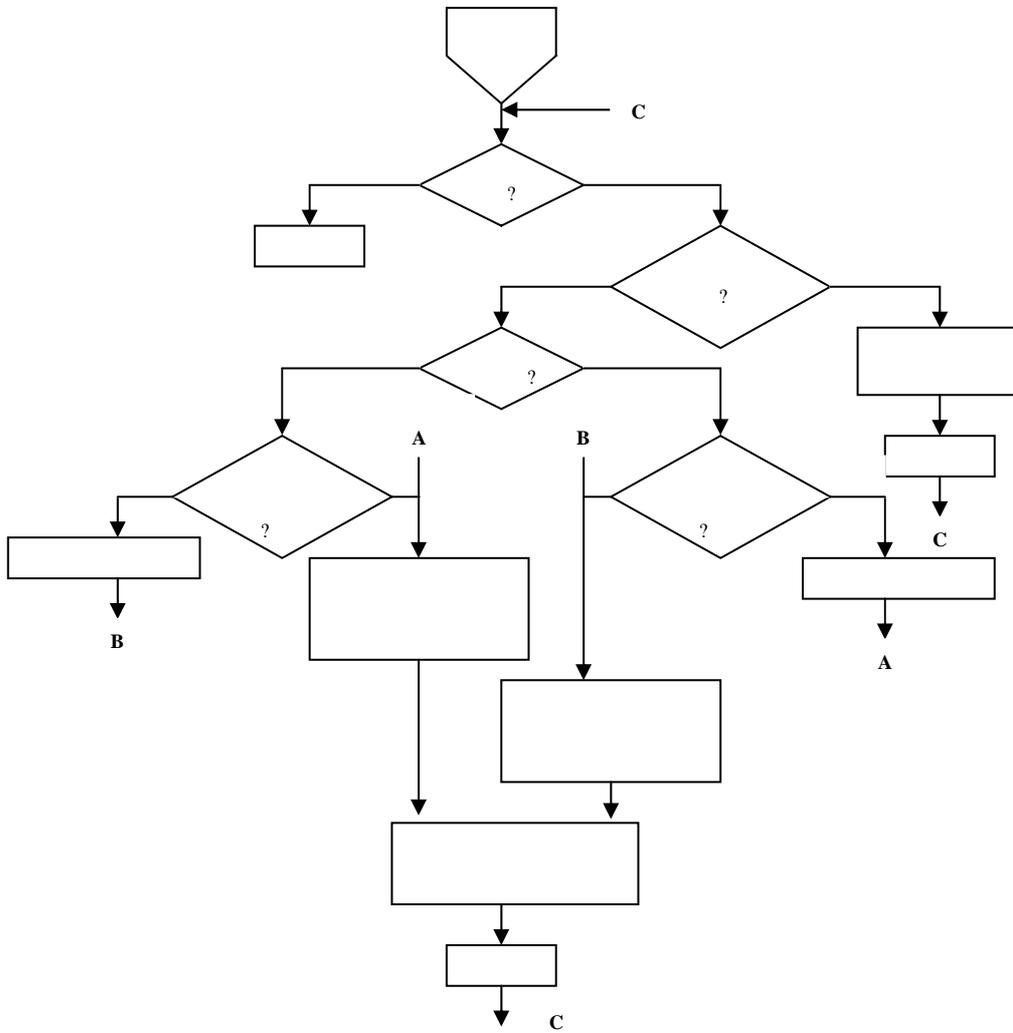
7 4 1
7 4 8
7 4 5
40 6 4
2 7 7

first-come first-served

0 7 0
7 40 2 7 7

7	4	1
7	4	5
7	4	8

7



1)"

"

5-10

5-10

I/O

I/O

"

"

I/O

I/O

2 " " **shortest seek time first algorithm** "

78 " " 40

3 " " **scan algorithm**

I/O

" " I/O " "

4 " " **N-steps scan algorithm** I/O

" " N " N "

5 " " **circular scan algorithm**

0 0

I/O

1 2 1 " "

" " " "

" " " "

5.5.6

RAID Redundant Array of Independent Disks
 Berkeley
 scheme

1987

RAID

CPU

4 RAID level 3

RAID3

RAID2

RAID3

X0 X3

X4

X4 i =X3 i X2 i X1 i X0 i
 X1 X4 i X1 i 2

X1 i =X4 i X3 i X2 i X0 i

RAID3 RAID4

RAID5

RAID

5 RAID level 4

RAID4

RAID5

I/O

I/O

RAID4 RAID5

bit-by-bit

X0 X3

X4

X1

X4 i =X3 i X2 i X1 i X0 i

X4' i =X3 i X2 i X1' i X0 i
 =X3 i X2 i X1' i X0 i X1 i X1 i
 =X4 i X1 i X1' i

6 RAID level 5

RAID5

RAID4

n

n

n

RAID4

7 RAID level 6 RAID level 7
RAID RAID6
RAID

RAID7 RAID6

“ ”

()

()

5.6.2

I/O I/O
I/O
/ () I/O / /
() (/)
I/O

5.7

5.7.1

4



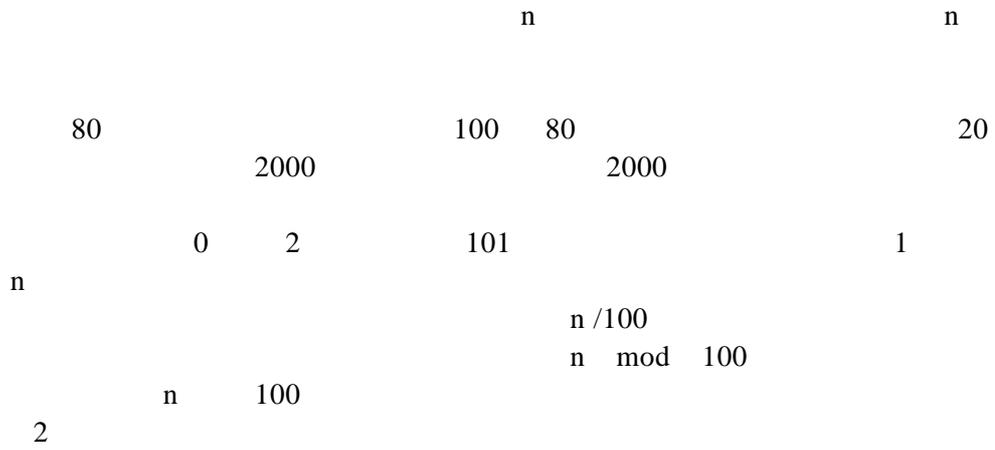
“ ”

“ ”

•

•

1



1

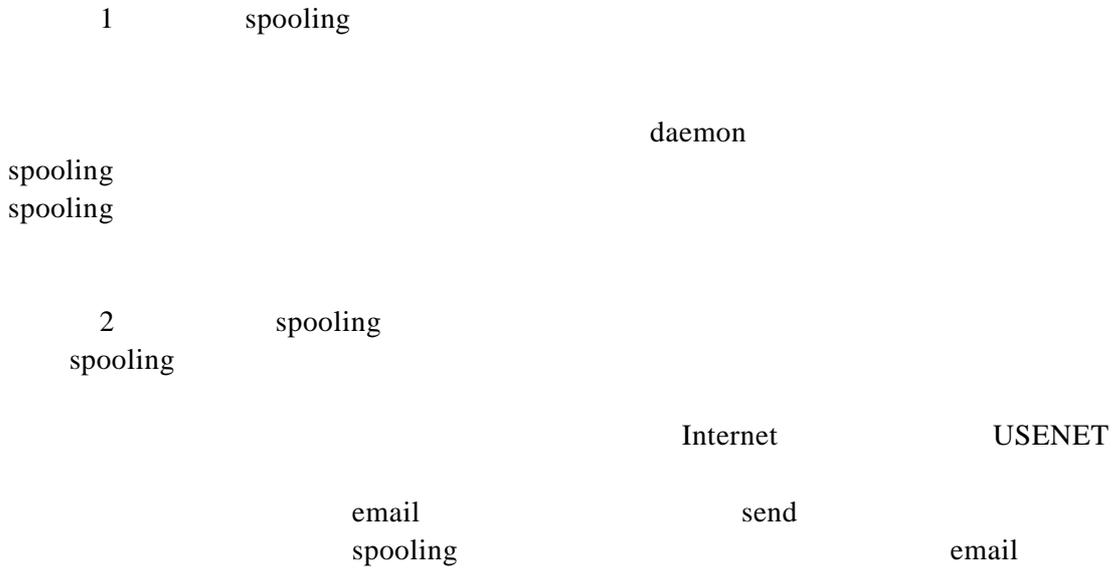
3

CPU

?

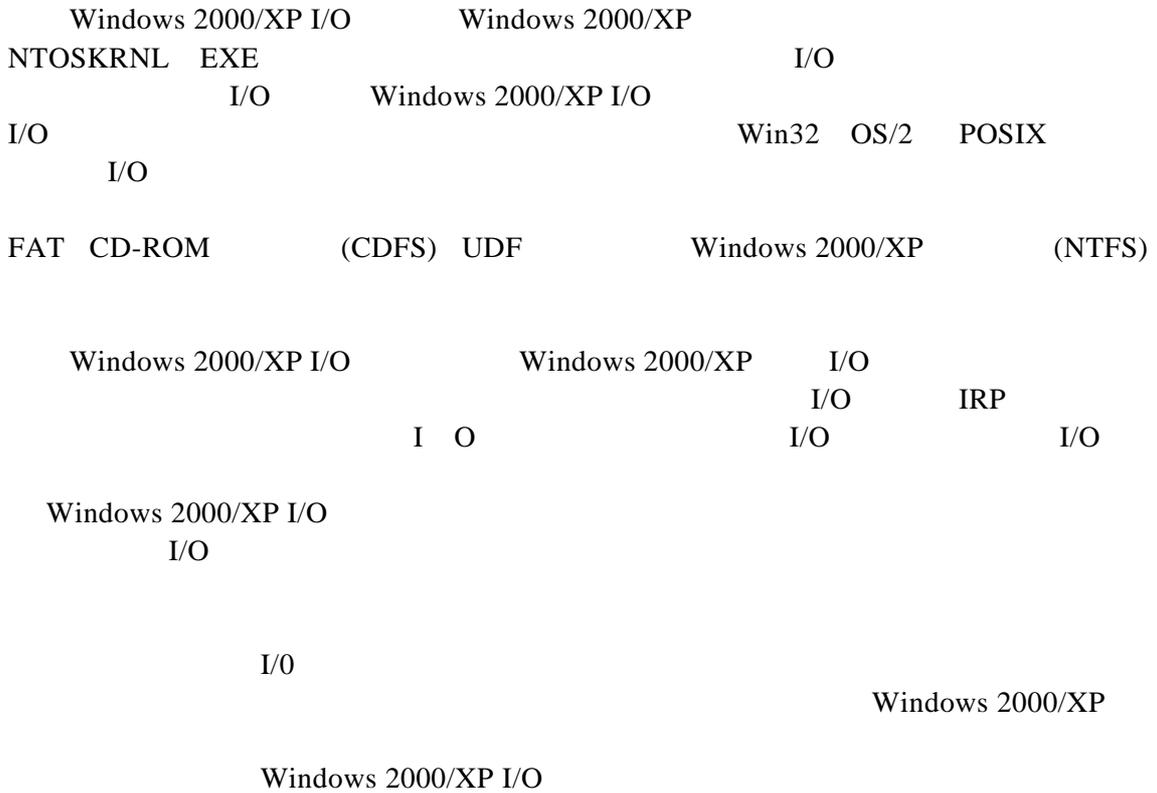
5.7.3 SPOOLING

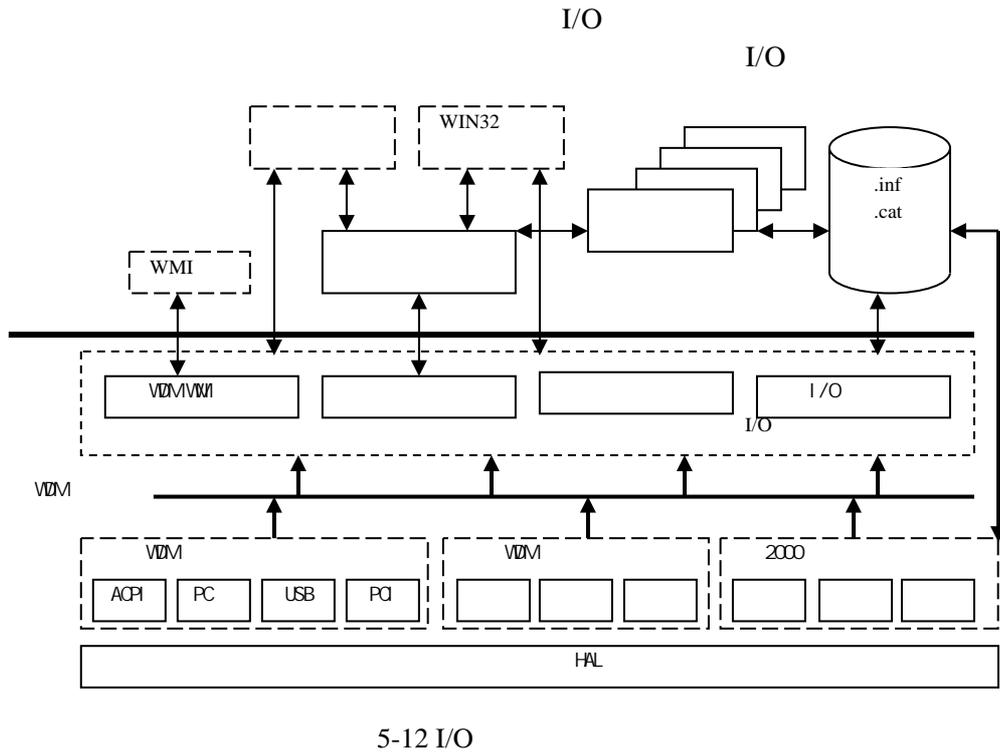
SPOOLING



5.8 Windows 2000/XP I/O

5.8.1 Windows 2000/XP I/O





Windows 2000/XP I/O 5-12

- API
- I O
- I/O IRP I/O
- I/O I/O
- I/O PnP(plug and play) I/O
- I/O
- WMI(Windows Management Instrumentation) WDM(Windows Driver Model)WMI WMI Windows
- WDM I/O WDM
- I/O I/O
- (HAL) I/O Windows 2000/XP

I/O
I/O

I/O

I/O

2 PnP

PnP (Plug and Play) I/O

PnP

PnP PnP I/O

PnP PnP I/O PnP

I/O PnP " " " "

PnP PnP PnP

PnP BIOS) PnP (PnP I/O

PnP Windows 2000/XP

PnP PnP USB USB

USB Windows 2000/XP PnP

- PnP
- PnP (I/O) (resource arbitrating)
PnP
- PnP I/O
- PnP Windows 2000/XP

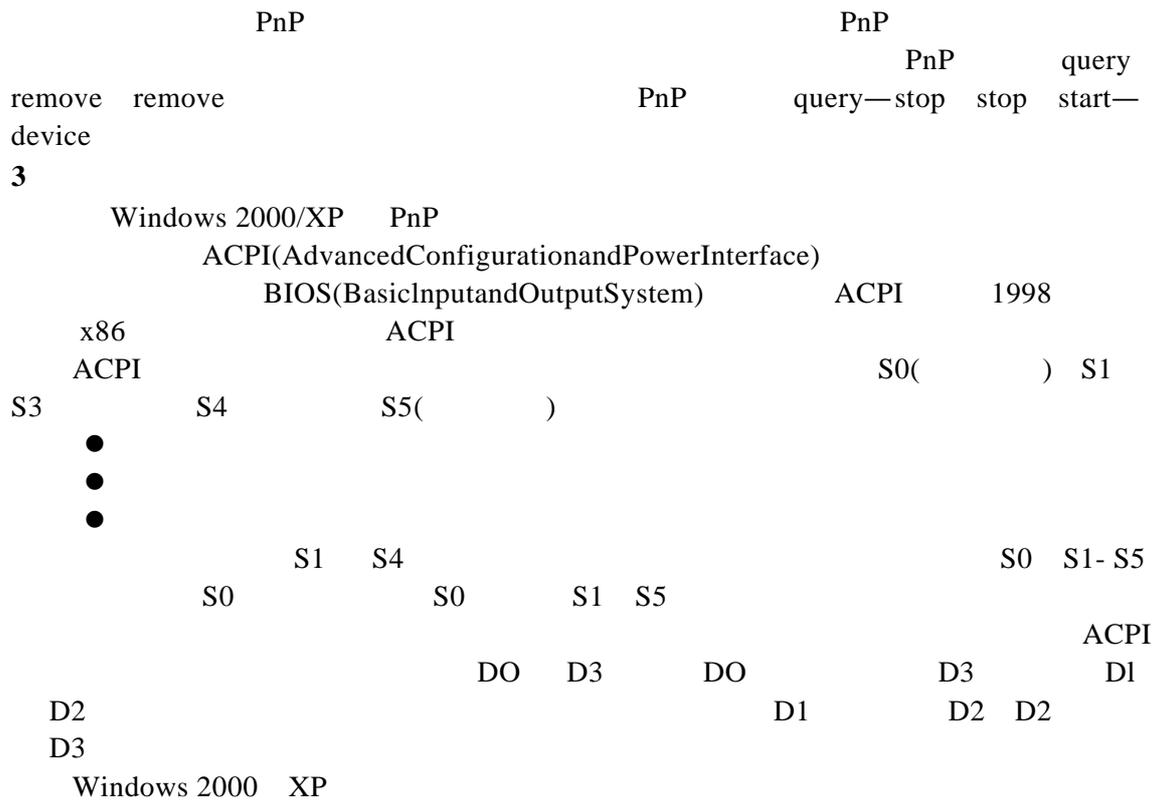
Windows 2000/XP PnP PnP PnP

PnP PnP PnP NT4 PnP

Windows 2000/XP PnP PnP PnP

PnP PnP PnP

(start-device) PnP PnP PnP

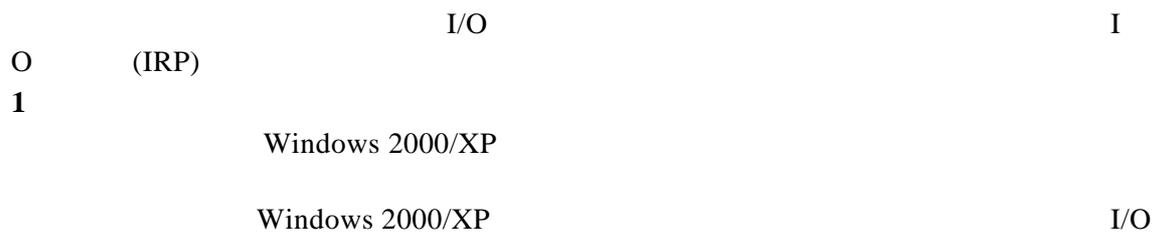


PoRegisterDeviceForIdleDetection

PoRegisterDeviceForIdleDetection

PoSetDeviceBusy

5.8.2 Windows 2000/XP I/O



I/O

Windows 2000/XP

(
) Windows 2000/XP I/O

	I/O
	I/O

I/O

(ACL)

I/O

ACL

Win32 LockFile

2

I/O

I/O

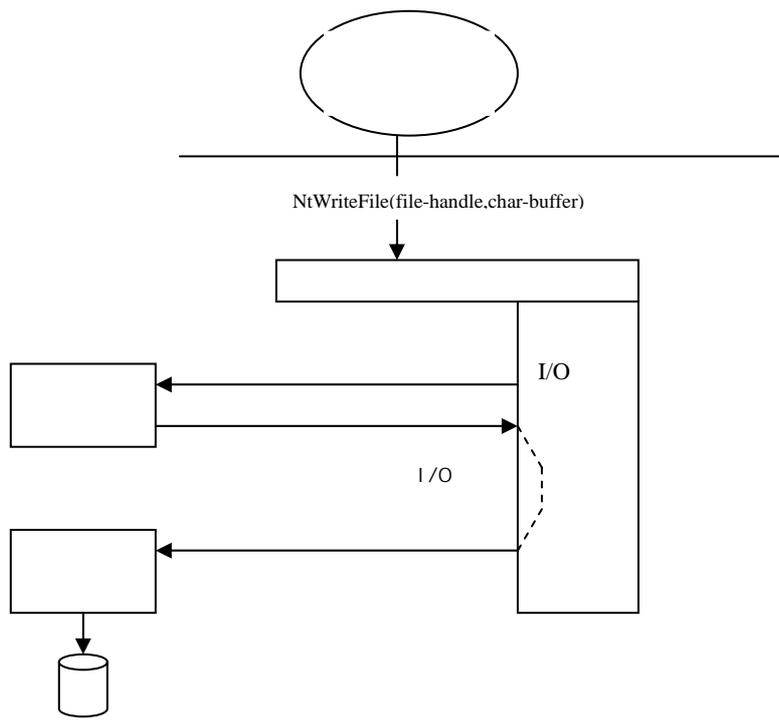
•

I/O

I/O

()

•



5-15

5-15

I/O

()

I/O

()

I/O

Windows 2000/XP

1

I/O

5-16

●

I/O

I/O

●

PnP

●

I/O

I/O

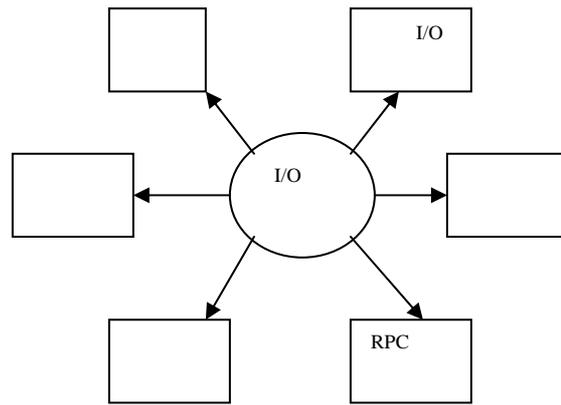
IRP

●

I/O

I/O

- (ISR)
Windows 2000/XP I/O ISR
- (IRQL) ISR IRQL (DPC) ISR
- (ISR) DPC DPC ISR
- DPC DPC ISR IRQ I/O I/O



5-16

- IRP
I/O
- I/O I/O I/O I/O
- I/O
- (I/O I/O)

2

- ()

- Windows 2000/XP

IRQL

I/O

ISR

ISR

ISR

Windows 2000/XP

ISR

(CPU)

ISR

Windows 2000/XP

ISR CPU IRQL

" "

CPU

ISR

()

ISR

5.8.4 Windows 2000/XP I/O

I/O I/O

I/O

I/O

1 I/O

I/O I/O

I/O I/O

1) I/O I/O

I O " "

ReadFile

I/O WriteFile

I/O

" I/O" I/O

I/O

FILE_FLAG_OVERLAPPED I/O Win32 CreateFile

I/O

(I/O)

I/O I/O IRP I/O

I/O I/O

Win32 HasOverlappedToCompleted

2) I/O I/O

I/O IRP

I/O

3) I/O I/O I/O I/O
 " I/O" I/O
 I/O

Win32 CreateFileMapping Map ViewOfFile I/O
 () I/O

I/O

Windows 2000/XP ()

4) / I/O
 Windows 2000 XP I/O " / "
 (scatter/gather) Win32 ReadFileScatter WriteFileScatter
 / I/O I/O
 I/O

2 I/O
 I/O
 (1)I/O DLL I/O
 (2) DLL I/O NtWriteFile
 (3)I/O IRP
 (4) I/O
 (5) CPU
 (6)I/O I/O
 I/O

1) I/O Windows 2000/XP
 I/O
 I/O
 ISR Windows 2000/XP ISR
 ISR IRQL
 DPC
 DPC I/O IRP
 I/O
 DPC IRQL
 Dispatch/DPC IRQL DPC

2) I/O
 DPC I/O

I/O " I/O " (I/O completion) I/O " I/O
 " (I/O status block) I/O I/O I/O
 I/O I/O
 I/O (APC)
 APC DPC DPC APC
 DPC
 I/O APC IRQL APC
 I/O APC I/O IRP ()
 I/O) I/O ()
)
 I/O I/O ReadFileEx
 APC I/O I/O
 APC I/O
 Platform SDK APC

5.8.5 Windows 2000 XP

1 Windows 2000 XP

Windows 2000 XP
 Windows 2000 XP ()
 Windows 2000 XP
 Win32 CreateFile
 Windows 2000 XP
 I O Windows 2000 XP
 I O
 Windows 2000 XP
 Windows 2000 XP
 (1)
 Windows 2000 XP
 CD-ROM
 (metadata)() Windows 2000 XP)
 (2)
 Windows 2000 XP

(section object)

I O (IRP)

(3)

(Win32 MapViewOfFile)
Windows 2000 XP

1 2
1 2
(4) (NetWare OpenVMS OS 2 UNIX
(logical block)
(virtual block caching) Windows 2000 XP
256KB

•

• I O (I O)

I O
(5)
Windows 2000 XP

NTFS
NTFS
Windows
2000 XP

(6)

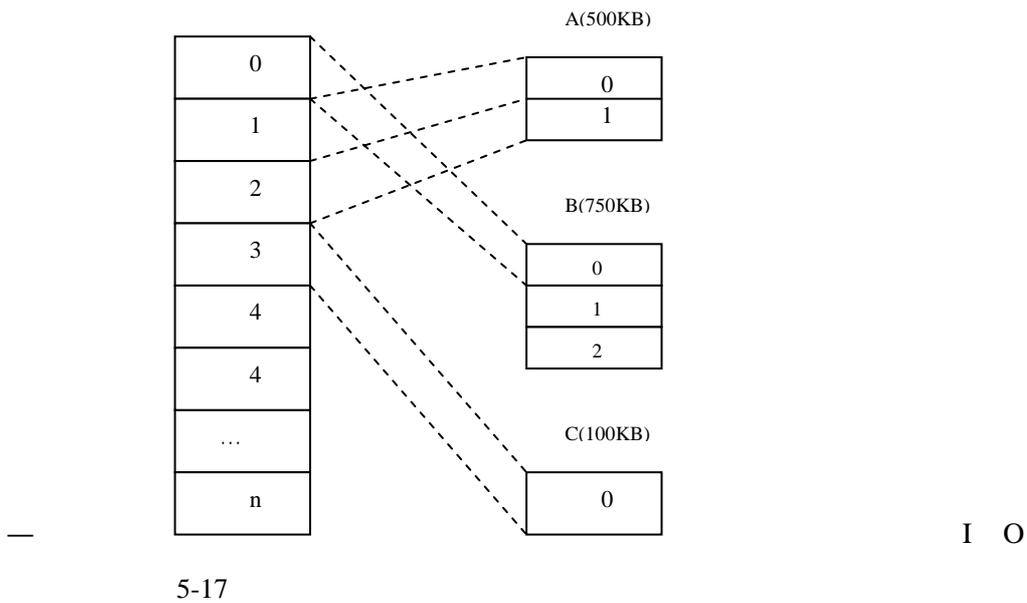
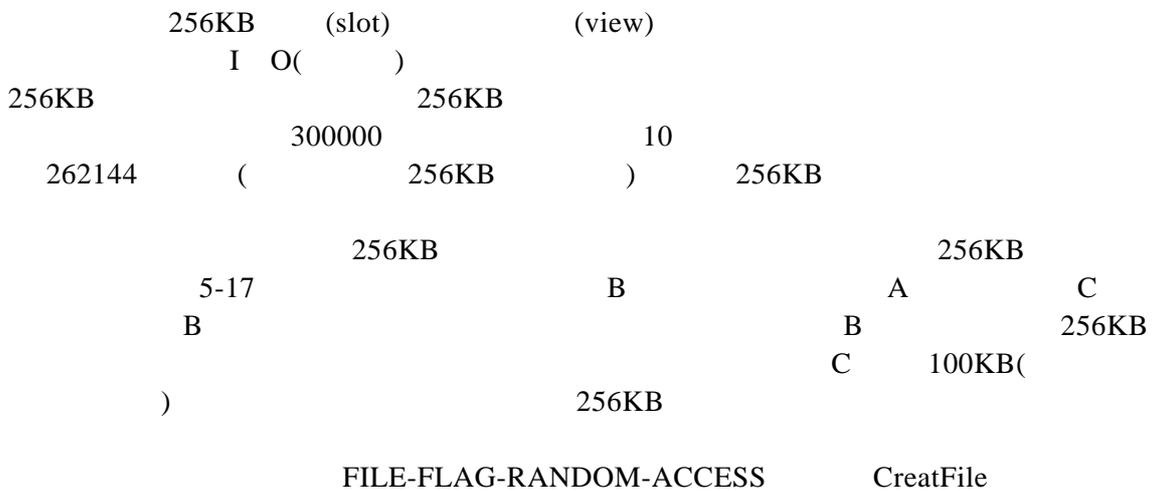
(recoverable file system) NTFS
I O
I O (log
file) ()

•

•

●
●
2

Windows 2000 XP



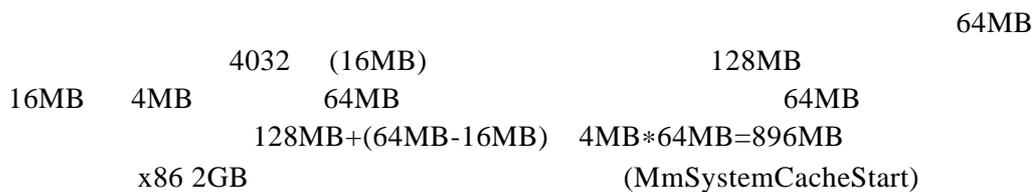
3

Windows 2000 XP

) ?

Windows 2000 XP

1)



(MmSystemCacheEnd)
0xC1000000-E0FFFFFF)

64/960MB(

512MB

2)

Windows 2000 XP

" "

)

4

(1)

256KB

VACB

(2)

(3)

1)

VACB(virtualaddresscontrolblock)

VACB

VACB

256KB

CcVacbs

VACB

VACB

5-18

VACB

(256KB)

VACB

O

VACB

I O

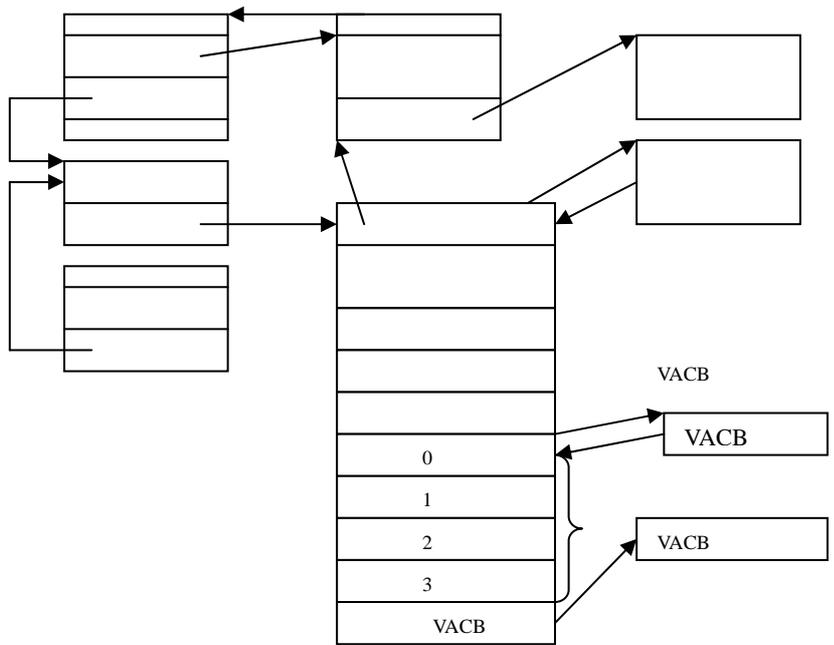
I



5-18 VACB

2)

) (VACB 5-19



5-19

? VACB ?
 VACB
 VACB " VACB " VACB
 256KB 256 KB 5-20

A
 VACB

VACB

VACB0	0
VACB1	1
VACB2	2
VACB3	3
VACB4	4
VACB5	5
...	
VACBn	n

337

VACB

(VACB) VACB

VACB

4 VACB
1 MB

VACB 256KB
1MB 256KB(1) VACB

32GB VACB
(sparse multilevel index array)

128 (-18) 7
7 18 VACB 256KB 256KB 2¹⁸B
128 2⁷ 128 2⁶³(
) 7 VACB

32GB 256KB
VACB 3 3

32GB 2³⁵ 3
VACB 3

128000 1000
5

I O (Win32 CreateFile)
FILE_FLAG_NO_BUFFERING

1) Windows 2000 XP (lazy writing)

(write-back)

I O

I O

) (

— 1 8 " " dirty () — —

2)

(threshold)

HKLM	SYSTEM	CurrentControlSet	Control	SeSSiOnManager
MemoryManagement	LargeSystemCashe	Windows 2000	XP Professional	
0	Windows 2000	XP Server	1	Windows 2000
				XP Server

Windows2000 XP Professional

4 8

4MB ()

2MB

3)

Win32 CreateFile FILE_ATTRIBUTE_TEMPORARY

—

4)

CreateFile

FILE_FLAG_WRITE_THROUGH

Win32FlushFileBuffers

5)

()

(Win32FlushFileBuffers)

" "

-
-

6)

Windows 2000 XP

(intelligent read-ahead)

ReadAheads Sec

CcReadAheadlos

Cache

File < 16 KB

7)

(virtual address read-ahead)

8)

30Tc<02c[(300)-5.200

I O

"

" (asynchronous read-ahead with history)

4000

3000

9)

I O

()
 ()
 •
 • ()

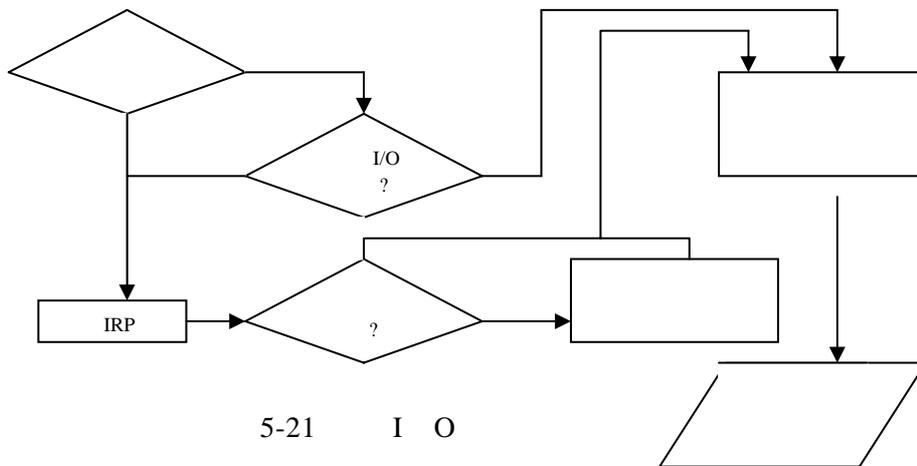
XP Professional	128	32 Windows 2000	64 XP Server	Windows 2000 256
-----------------	-----	--------------------	-----------------	---------------------

10) I O

I O (fast I O)
 I O (IRP) I
 I O
 IRP
 Windows 2000 XP

IRP
 I O
 ()
 I O
 I O
 I O
 I O
 (Win32 LockFile
 UnlockFile)

5-21



I O

1)
 2) I O I O
 I O
 I O
 3) I O I O
 (I/O
 I/O IRP
)
 4)
 5)
 6)
 ●
 ● " "
 ●
6

CcInitializeCacheMap

(1)" "
 (2)" "
 (3)" "

 " () IRP "

1)
 ()

CcCopyRead	
CcFastCopyRead	CcCopyRead 32 NTFS FAT
CcCopyWrite	
CcFastCopyWrite	CcCopyWrite 32

	NTFS	FAT
--	------	-----

2)

()

CcMapData	
CcPinRead	/
CcPreparePinWrite	
CcPinMappedData	
CcSetDirtyPinnedData	
CcUnpinData	

) (

3)

(direct memory access DMA)

DMA

DMA

(1KB 2KB)

DMA

(MDL) 4

DMA

DMA

--	--

CcMdlRead		MDL	
CcMdlReadComplete	MDL		
CcMdlWrite		MDL	0
CCMdlWriteComplete	MDL	" "	

7

Windows 2000 XP

9600

CcSetDirtyPageThreshold —

5.9 Linux

5.9.1 Linux

Linux

Linux

Linux

mknod

Linux

Linux

I/O

I/O

•

Linux

•

•

•

•

Linux

Linux

•

Linux polling

Linux Linux Linux

I/O I/O Linux Linux

I/O /proc/interrupts Linux

SCSI I/O Linux DMA I/O

7 DMA 16M DMA

16 8 DMA

DMA Linux dma_chan DMA dma_chan

Linux cat/proc/dma dma_chan

Linux Linux

Linux device_struct fs/devices.h chrdev device_struct

Linux chrdev blk_devs Linux blk_devs

chrdev blk_devs blk_dev_struct blk_dev_struct request request

5.9.2 Linux

Linux Linux DOS EXT2 Linux

Linux Linux gendisk

/include/linux/genhd

Linux IDE Inergrated Disk Electronic SCSI

Small Computer System Interface I/O Linux

IDE IDE IDE IDE

IDE Linux IDE IDE

ide_hwif_t ide_hwifs

ide_hwif_t ide_drive_t IDE IDE

Linux CMOS

Linux IDE

/dev/hda /dev/hdb /dev/hdc ... IDE Linux

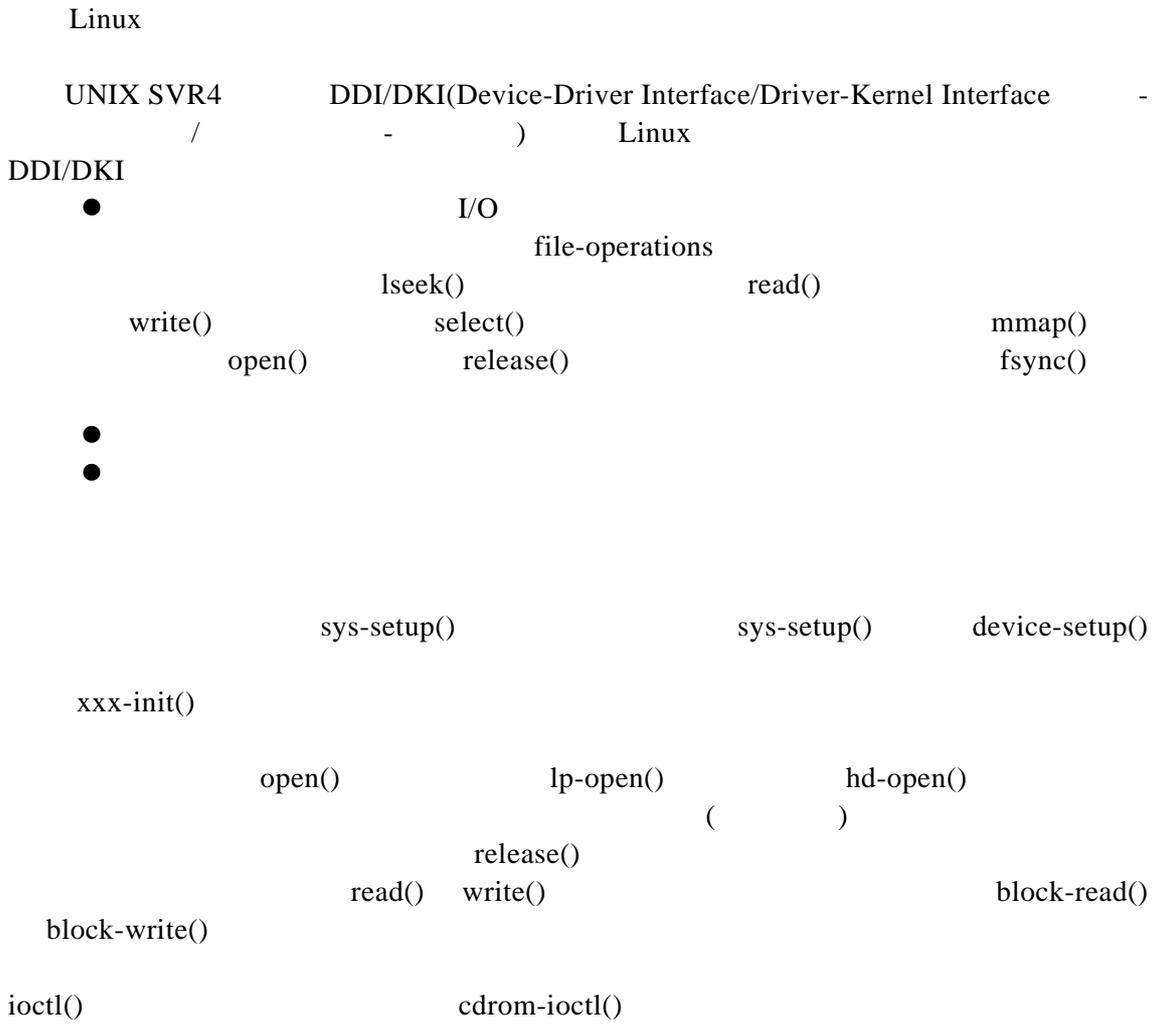
IDE IDE IDE 3 IDE

22 IDE blk_devs
 3 22 SCSI SCSI SCSI
 40MB/ SCSI 32
 8
 1 BUS FREE
 2 ARBITRATION SCSI SCSI
 SCSI SCSI
 3 SELECTION SCSI SCSI
 SCSI SCSI
 4 RESELECTION SCSI
 SCSI
 5 COMMAND 6B 10B 12B
 6 DATA IN DATA OUT
 7 STATUS
 8 MESSAGE IN MESSAGE OUT
 Linux SCSI host device
 Host SCSI SCSI SCSI SCSI
 SCSI host Device SCSI SCSI
 Device

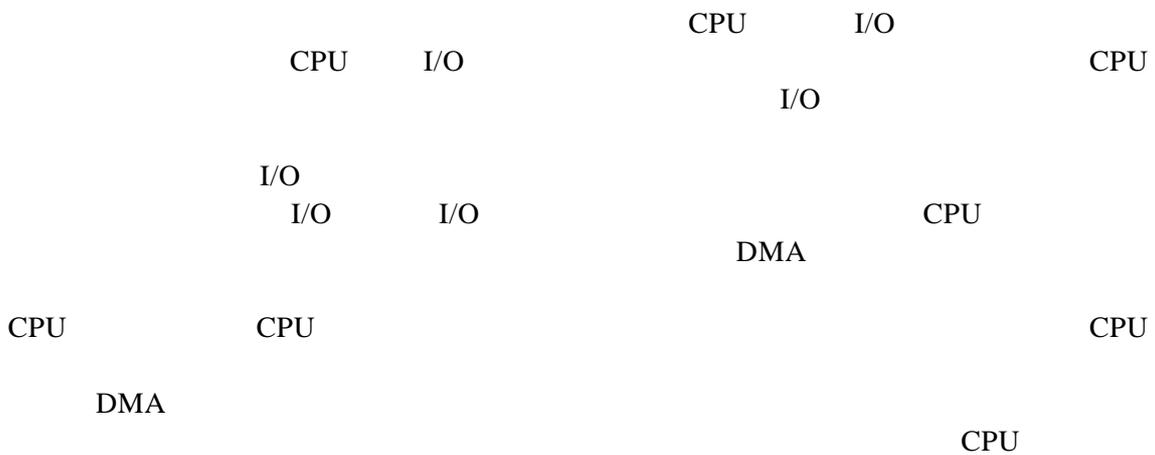
5.9.3 Linux

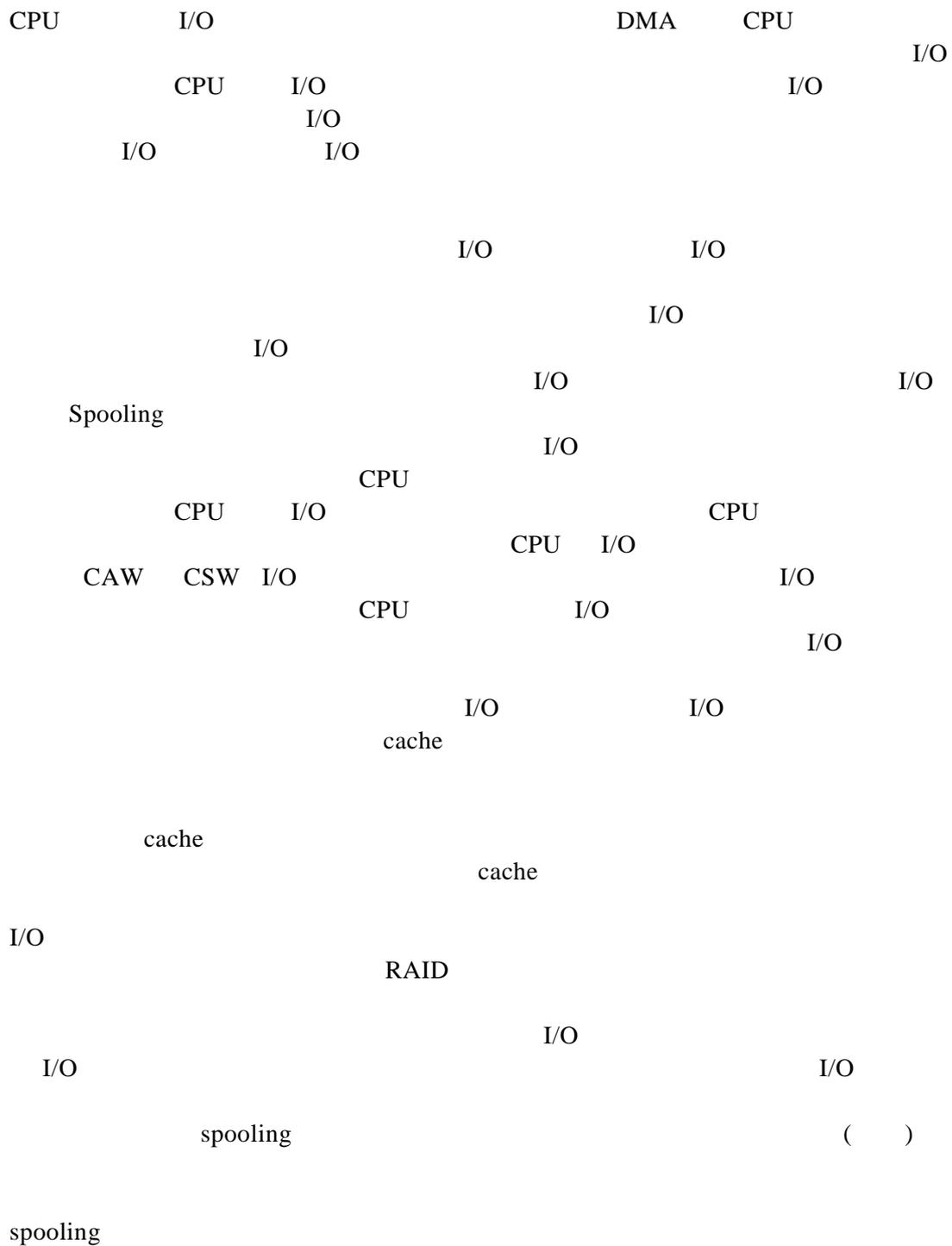
Linux device
 Device
 1
 0 /dev/ethN /dev/seN SLIP /dev/pppN PPP
 /dev/lo
 2 base address DMA DMA irq
 channel
 3 IP
 4 mtu AF_INET
 internet
 Linux
 X.25 SLIP PPP Apple Localtalk
 5 sk_buff
 6

5.9.4 Linux



5.10





- 1
- 2 I/O
- 3 DMA
- 4

5 I/O
 6 I/O
 7 I/O
 8 CPU ?
 9
 10
 11 I/O
 12
 13
 14
 15
 16 " " " "
 17 " "
 18
 19
 20
 21 Spooling
 22 Spooling
 23 Spooling ?
 24 Spooling ?
 25 Windows 2000/XP I/O
 26 Windows 2000/XP
 27 Windows 2000/XP
 28 T
 M C
 max C T +M
 29
 30 IBM370 I/O
 31
 32
 33 ?
 34
 35
 1 I/O 20
 2 1
 1 20 20 2 1 20

2 8 18 27 129 110 186 78 147 41 10 64 12

100

3

4

512

5

50 121 75 80 63

1569

?

5

1	7	2	8
2	7	2	5
3	7	1	2
4	30	5	3
5	3	6	6

1

6

40

0~39

11

1 36 16 34 9 12 1

FCFS 2

SSTF 3 SCAN

7

200

0~199

143

125

86 147 91 177 94

150 102 175 130

(1)

FCFS

(2)

SSTF

(3)

SCAN

(4)

8

FCFS

(1)

(2)

(3)

9

100

190 10 160 80 90 125 30 20 29 140 25

10

100

23 376 205 132 19 61 190 398 29 4 18 40

11

L

B

(1)

?(2)

K

I/O

?

12

200

20

8

1024B

606

CH6

" "

()

() (

)

-
-
-
-
-
-
-
-

6.1

6.1.1

" "

;

6.1.2

" "

" ?"

" *"

0 3 MS-DOS 1 8

Windows

\ / < > "

COM

LIB BAT EXE OBJ

MS-DOS

Windows-98 UNIX

" asm" C

255

Windows UNIX

" c"

6.1.3

UNIX



-
-
-

socket

UNIX/Linux

r w x

10

-rwxrwxrwx

- 1 (b/c) (-) (d) (l)
- 2-4
- 5-7
- 8-10

-rwxr-x--x

6.1.5

IBM

1

x i

i

2

	1024	4096
22	48	9

3

6.1.6

rmkdir JCL UNIX cat cd cp find mv rm mkdir

•

•

• /

/

•

•

•

6.2

6.2.1

" "

" "

FCB(File Control Block)

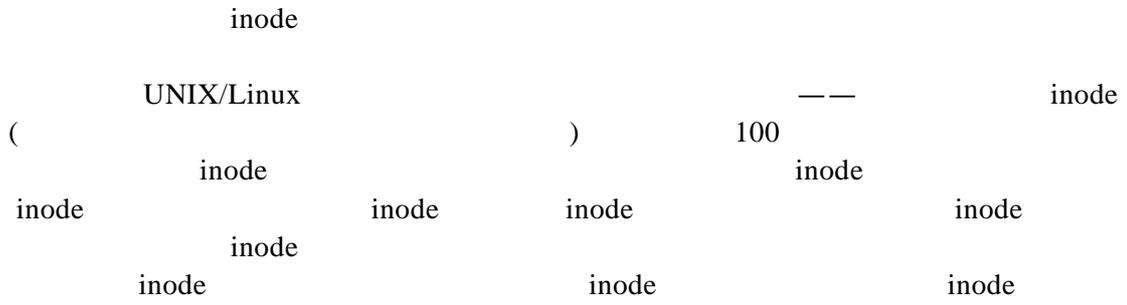
•

•

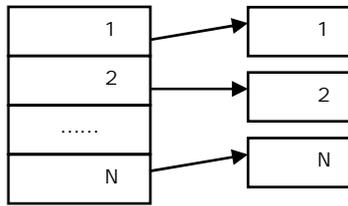
•

Z

- di-size
- di-add[8]
- di-atime
- di-mtime
- di-ctime



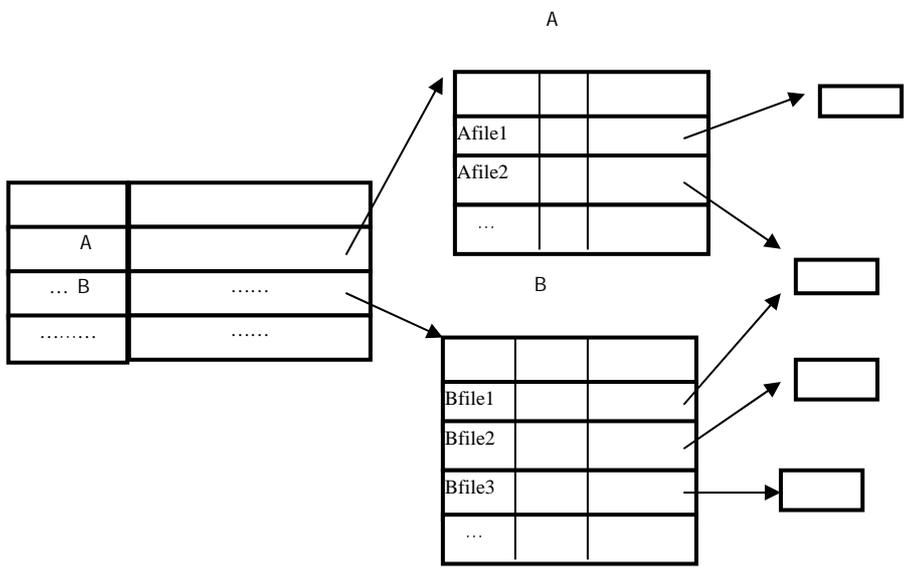
6.2.2



6-1

6.2.3

6-2



6-2

6.2.4

6-3

/ /

" ." inode " .." inode
../feil/myfile.c " .."
feil inode

6.3

6.3.1

6.3.2

1

COBOL

.....

COBOL

•

•

•

2

' M' ' F' " " 1
1 ' ' ' ' 2 1

•

•

FORTRAN

COBOL

2

80
80

80

800

10

10

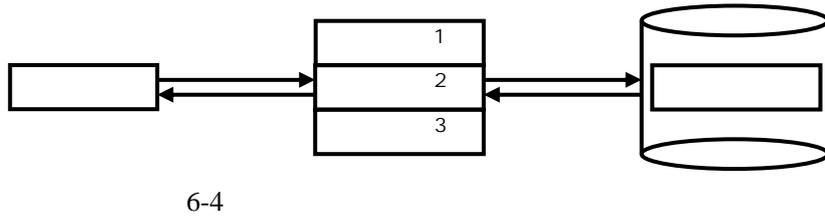
20

20

1600

6-4

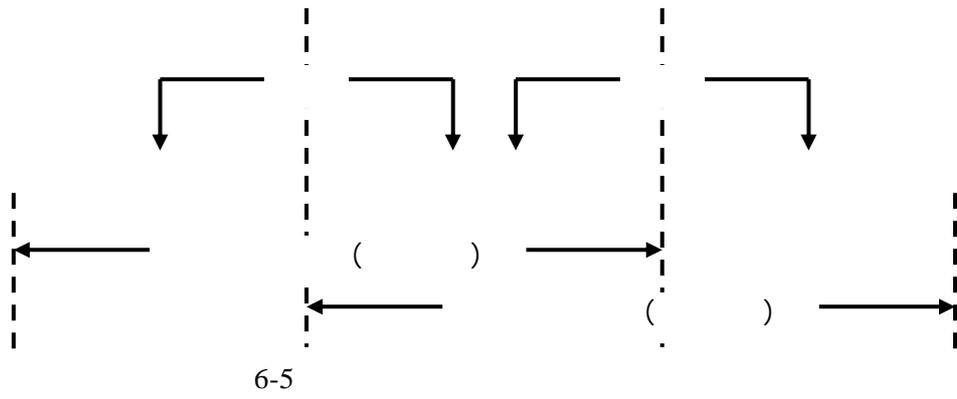
I/O



3

-
-
-

6-5



- F
- V
- S

F

V

-
-

RL

RL

BL

S

-
-
-
-

1000

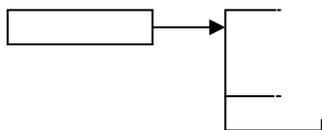
6.3.3

1

2

6-6

0



•

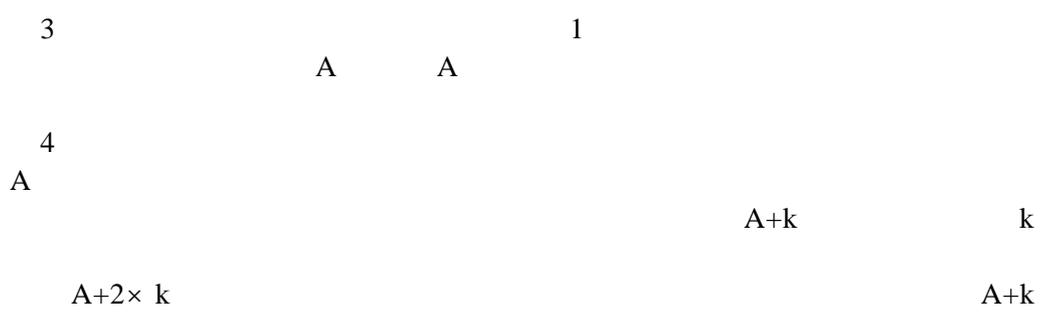
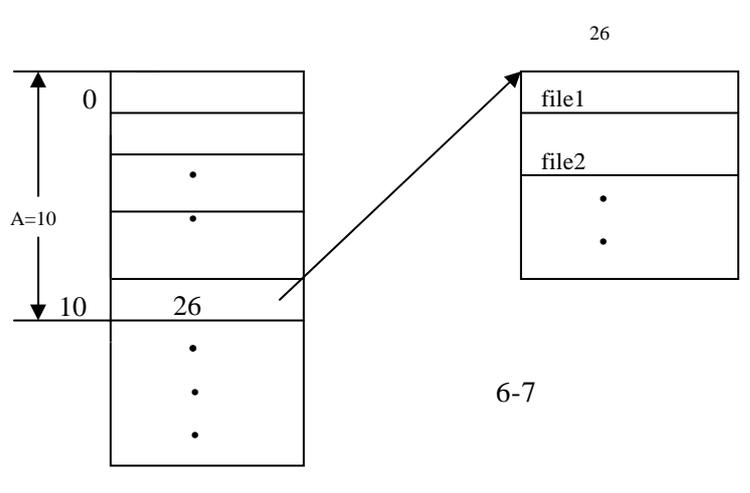
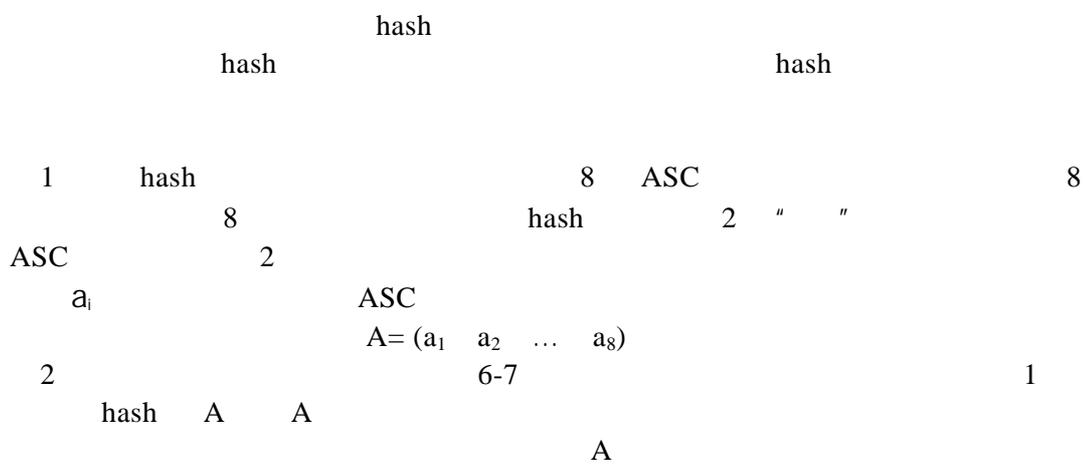
•

•

3

(hash)

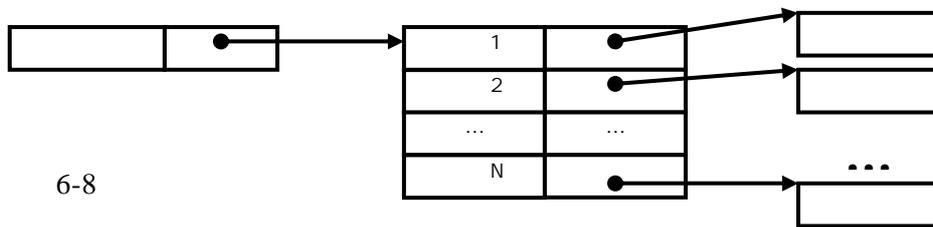
cache	hash	inode cache	directory cache	Linux	hash	cache
					buffer cache	Linux



4

(6-8)

()

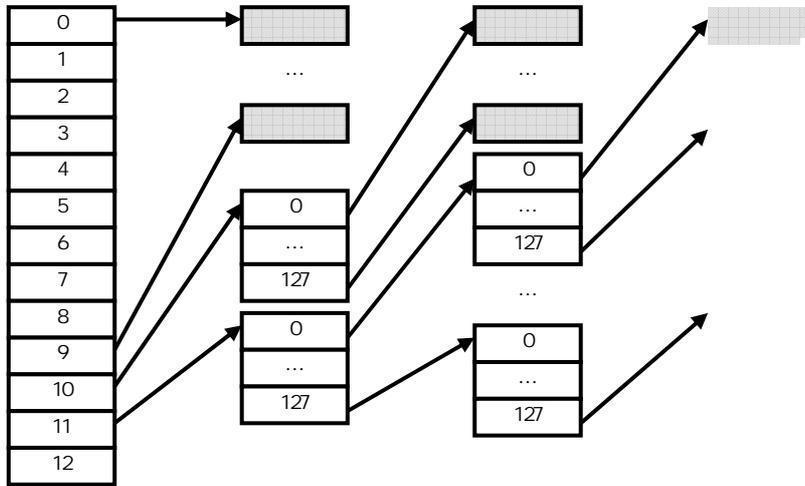


()

()

n n (n+1)

()



UNIX/Linux 13 4 6-9
 UNIX/Linux) 10 (0 9
 128 10 11 512
 UNIX/Linux 12 13 11
 80% 10 20%
 10

6.4

6.4.1

' ' ,
' ' ,

1

' ' ,

•

•

•

•

•

' ' ,

2

' ' ,

' ' ,

•

•

•

•

3

-
-
-
-

4 I/O

-
-
-

5

-
-
-

6.4.2 UNIX/Linux

UNIX/Linux

1

1

inode

C

```
int fd, mode;
char *filenamep;
fd = creat (filenamep, mode);
filenamep
```

mode

i_mode

fd

creat

" "

fd

/usr/lib/d2

C

creat

```
int fdlib;
```

```
fdlib = creat ("/usr/lib/d2", 0775);
```

```
                d2
d2                /usr/lib                i_mode    0775
                d2
i_nlink    " 1"
```

```
                f_flag    "  "                f_offset    " 0"
                d2
```

```
                "  "
2)
```

```
                i_link    " 1"
                unlink (filenamep)                creat
                "  "
```

2

```
                "  "
```

```
                "  "
```

1

```
                int fd, mode;
                char * filenamep;
                fd = open (filenamep, mode);
```

```
                mode                0                1
                2                creat                open
```

```
                mode
```

```
                "  "
```

```
                i_count    " 1"                i_count
```

(2)

```
int fd;
close (fd);
```

fd close

fd

f_count " 1" " 0"

i_count " 1" " 0"

f_count i_count

f_offset f_offset

3

" "

" "

f_offset

1)

```
int nr, fd, count;
char buf [ ]
nr = read (fd, buf, count);
```

fd

count

buf

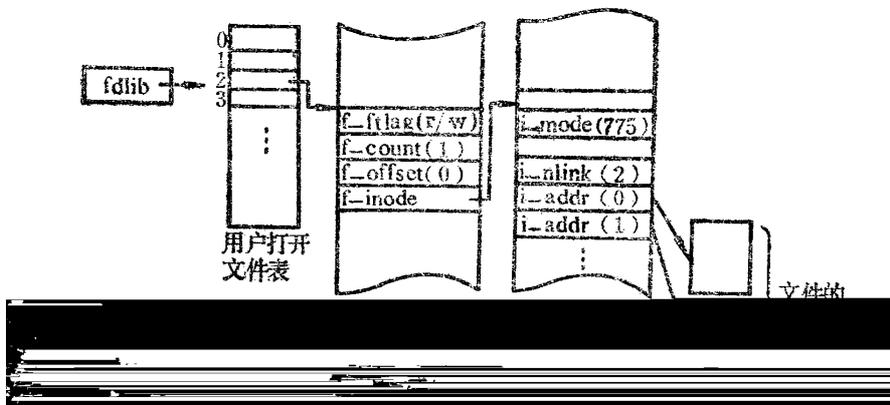
nr

count

nr

read

" 0"



6-10

/usr/lib/d2

6-10

d2 1500

bufp

number

```

    read
    number = read (fdlib, bufp, 1500);
    read
    f_flag
    f_offset
    i_addr
    bufp
    read

```

2)

```

    nw = write (fd, buf, count);
    fd, count  nw
    read
    buf
    nw
    buf
    count

```

4

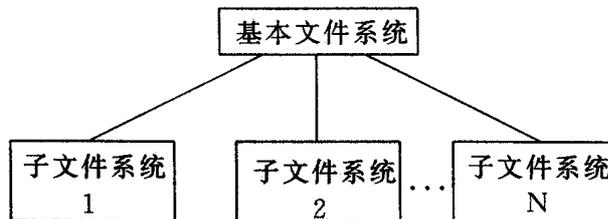
```

    " "
    offset
    lseek
    f_offset
    f_offset
    long lseek;
    long offset;
    int whence, fd;
    lseek (fd, offset, whence);
    fd
    whence " 0"
    f_offset
    offset
    offset
    whence " 1"
    f_offset
    offset

```

6.4.3

()
 () Windows Linux)



6-11 Unix/Linux

Windows/DOS
 Windows/DOS
 UNIX/Linux UNIX/Linux

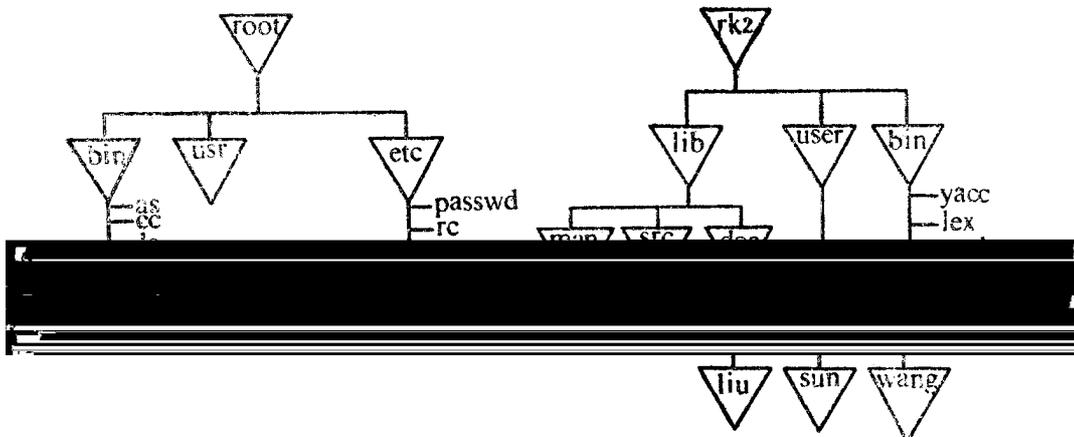
6-11

(1)

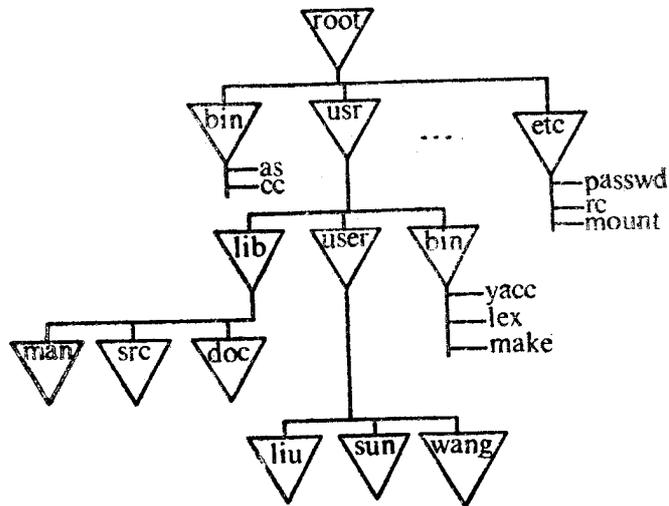
(2)

UNIX/Linux

C



(b) 待安装的文件卷 rk2



(c) 把文件卷 rk2 安装到 usr 节点之后

6-12

"

" "

" "

6-12

(a) rk2(rk2) (c) ()mount
 mount("/dev/rk2" "/usr" 0)
 rk2 usr
 (mount point) " " " "

inode (/usr) inode (/dev/rk2)

(1)

(2)

6.4.4

()

UNIX/Linux

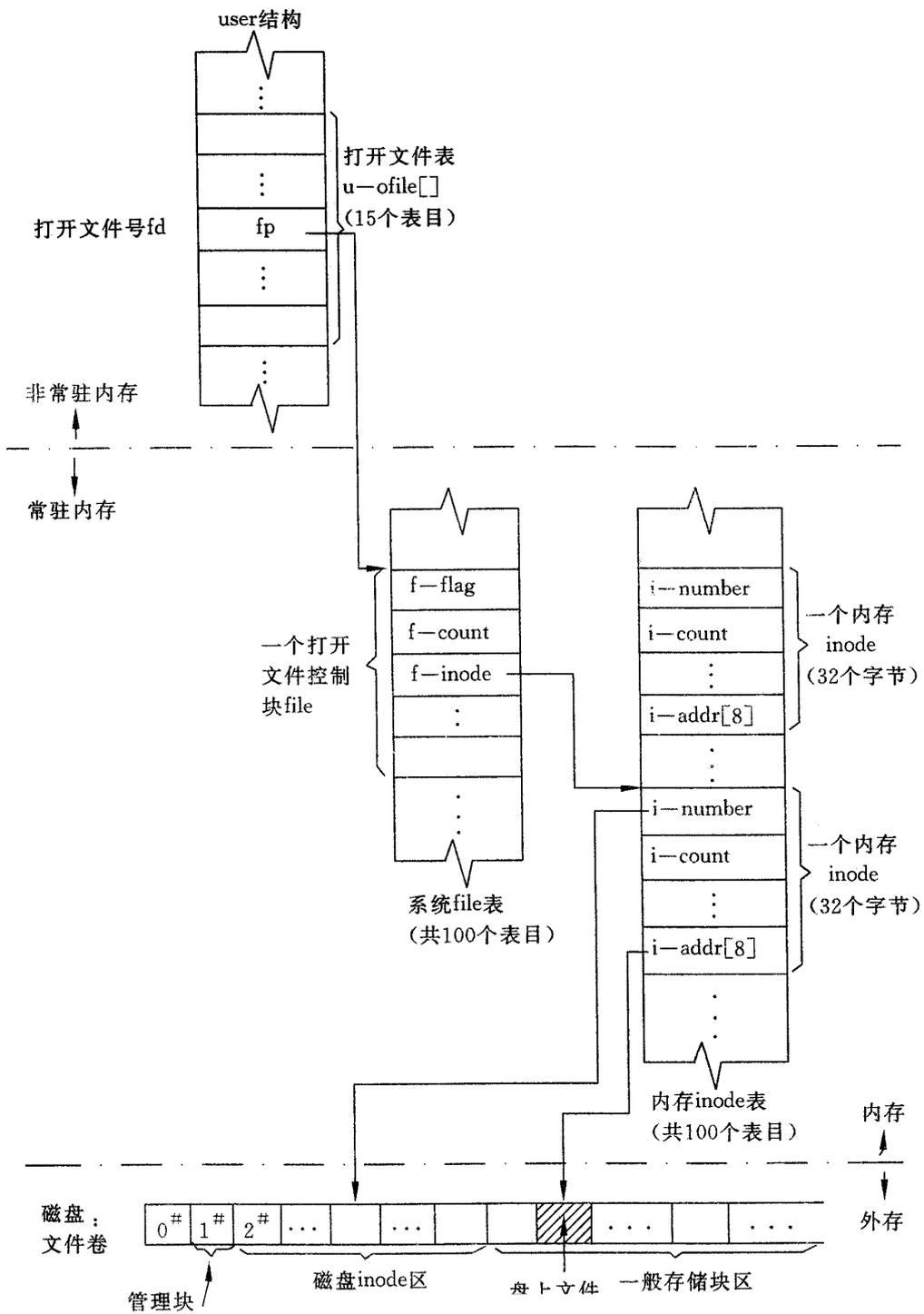
UNIX./Linux

(file link)

UNIX/Linux

myfile.c () ()
 include 6-3 fei1 fei2
 UNIX/Linux inode testfile.c

```
chat * oldnamep, * newnamep;  
link (oldnamep, newnamep);  
oldnamep newnamep
```



6-13 Unix/linux

UNIX/Linux

/

/

UNIX/Linux

user

user

/

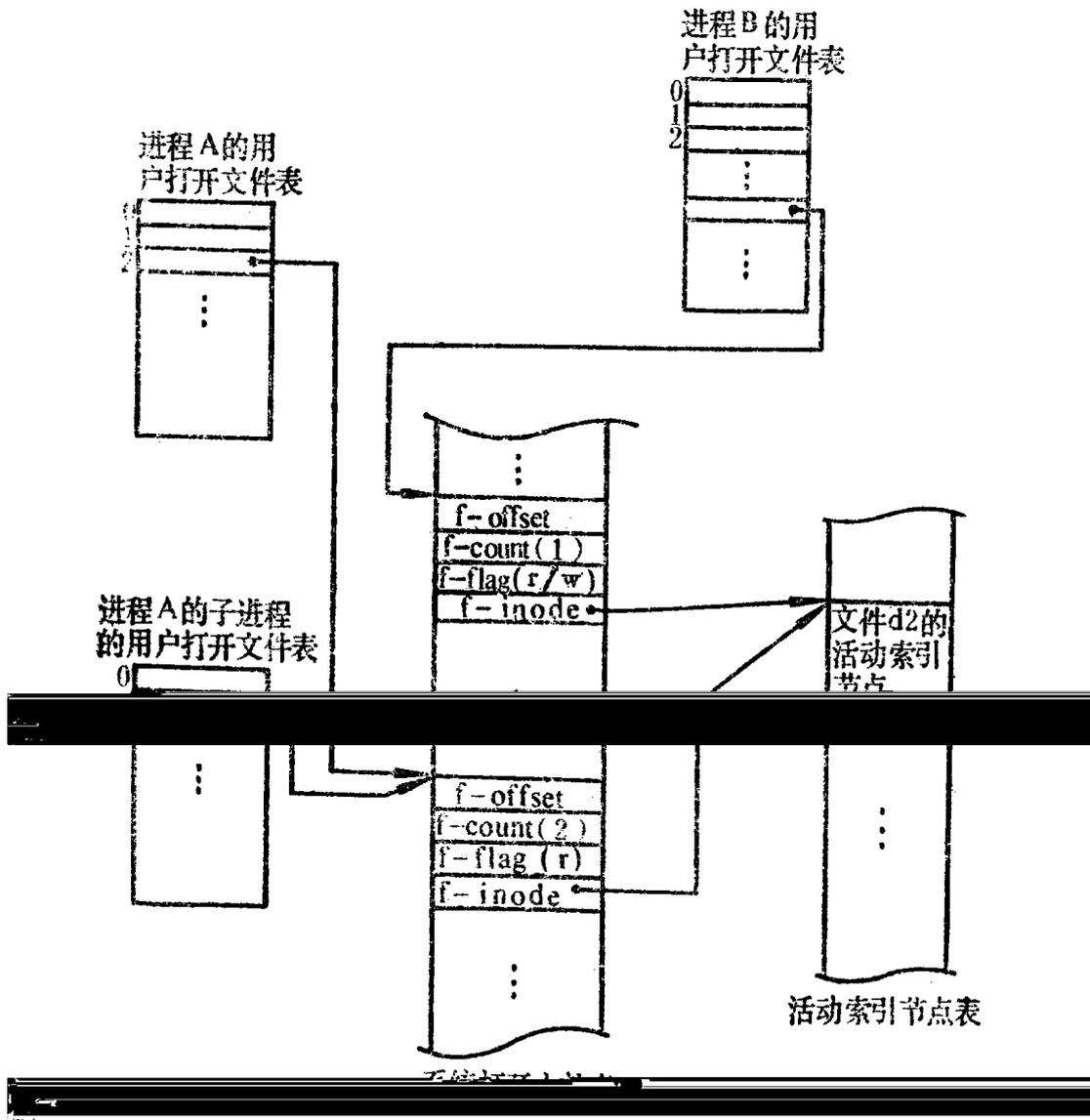
user

fork

6-14 UNIX/Linux

6-14 UNIX/Linux

Widely used, X[...]



6-15 UNIX

	100				
	f_offset		f_inode	f_count	f_flag
				6-14	6-15
6-15	A				d2
f_count	" 2"	B			
f_count	" 1"			A	d2

•

‘ ‘

1

CP/M VM/SP Windows Macintosh

'0'

2

3

DOS

4

6-17 UNIX/Linux 438 512 12 349

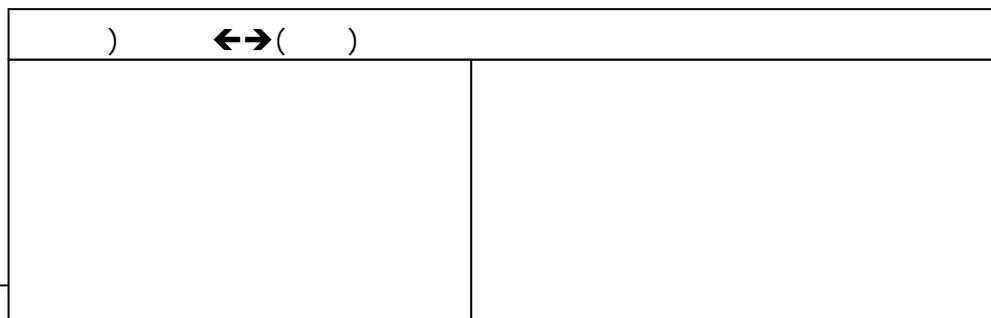
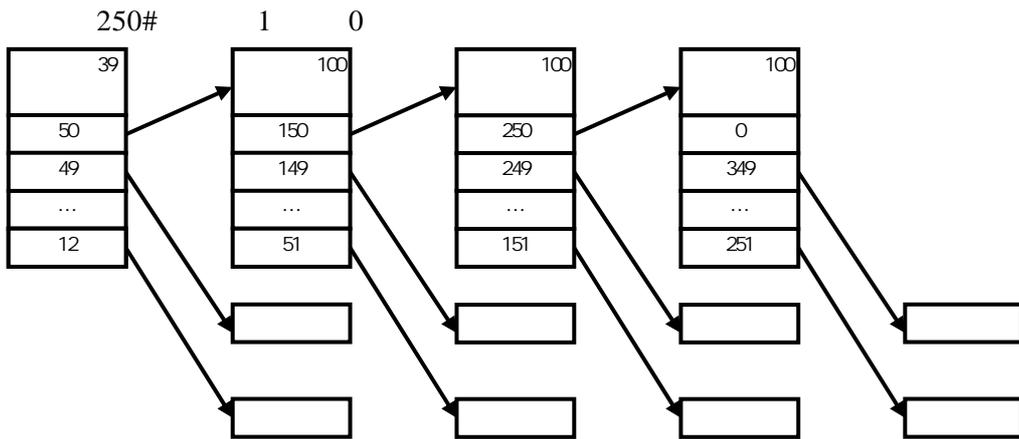
100

50#-12#

50# 150#

100 100

150#-51# 250#-151#



0

6.4.7

2000/XP
memory-mapped file
I/O

MULTICS
Windows
I/O

512KB
0

512KB+100

512KB
512KB

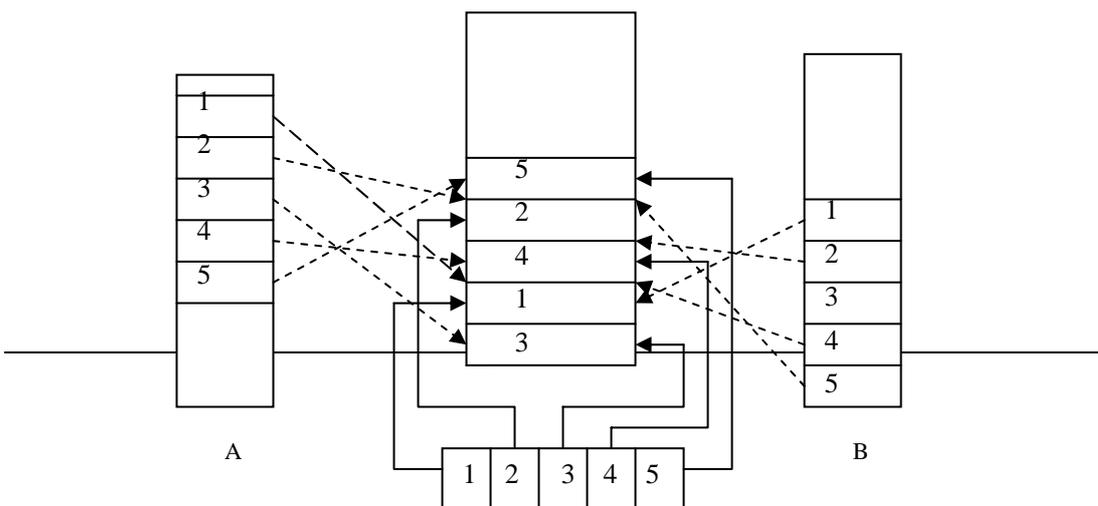
100

64KB

0

512KB+100

6-18



I/O

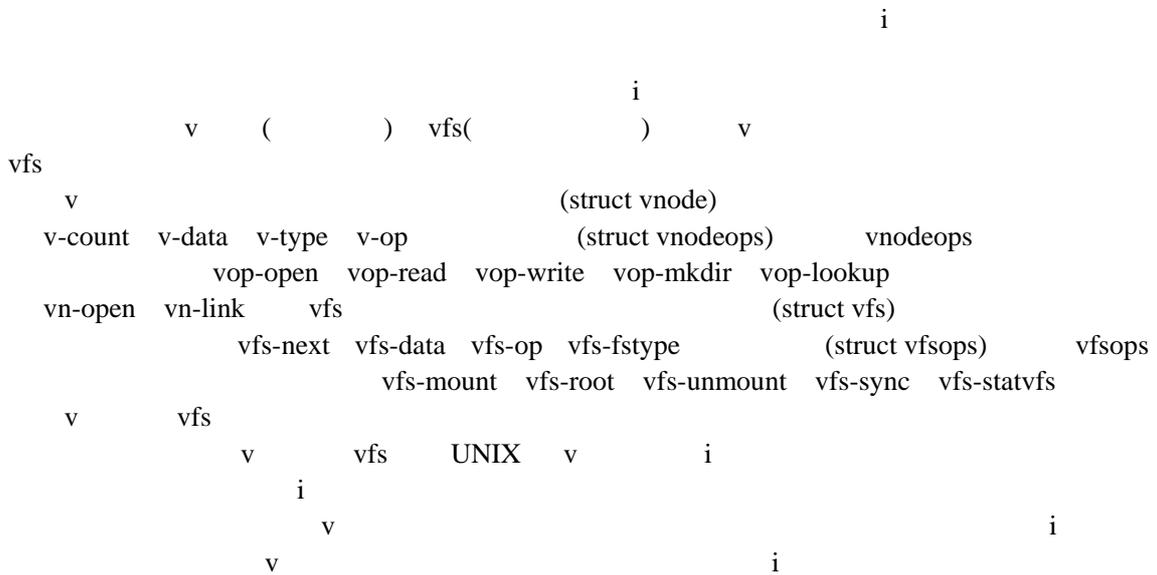
0

6.4.8



(1)

(2)



UNIX SVR4 v

```

Struct vnode {
    u-short v-flag; /* */
    u-short v-count; /* */
    struct vfs *vfsmountedhere; /* */
    struct vnodeops *v-op; /*v */
    struct vfs *vfs p; /* v */
    struct stdata *v-stream; /* */
}

```

```

struct page *v-page; /* */
enum vtype v-type; /* */
dev-t v-rdev; /* id*/
caddr-t v-data; /* */
}
UNIX SVR4 v
vops
( )
vfs vfs
struct vnode {
struct vfs *vfs-next; /*vfs */
struct vops *vops; /* */
struct vnode vfs-vnode covered; /* v */
int vfs-fs type; /* */
caddr-t vfs-data; /* */
dev-t vfs-dev; /* id*/
}
v v
v vfs
UNIX SVR4
( )
mount vfs mount
mount lookupn() v v
( )
vfs
vfs-mount v
mount

```

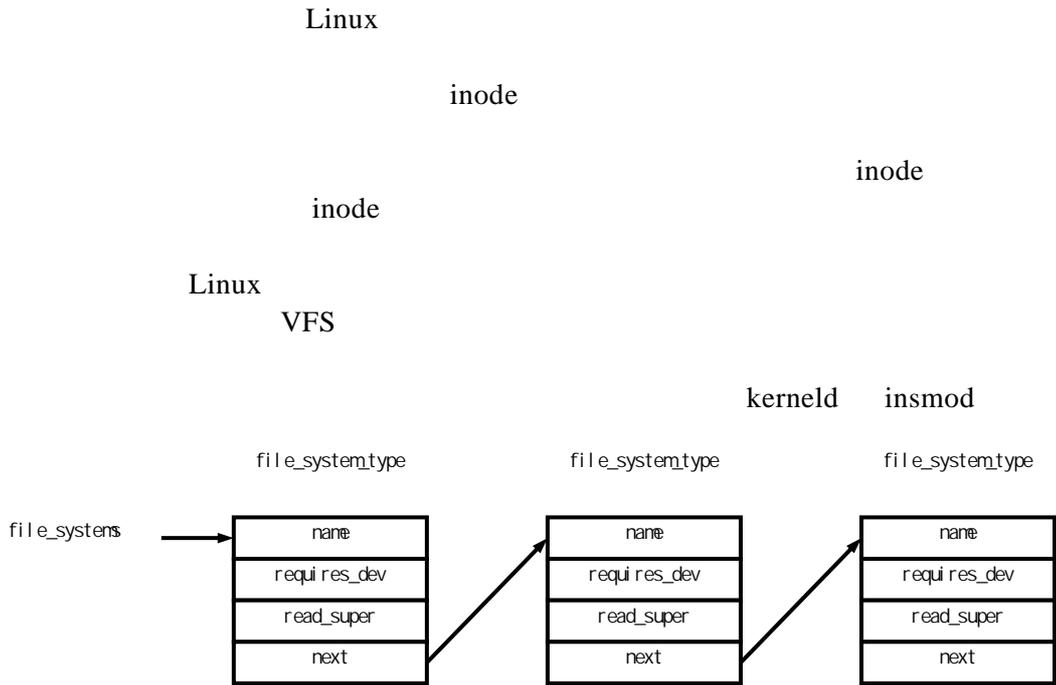
6.5 Linux

6.5.1 Linux

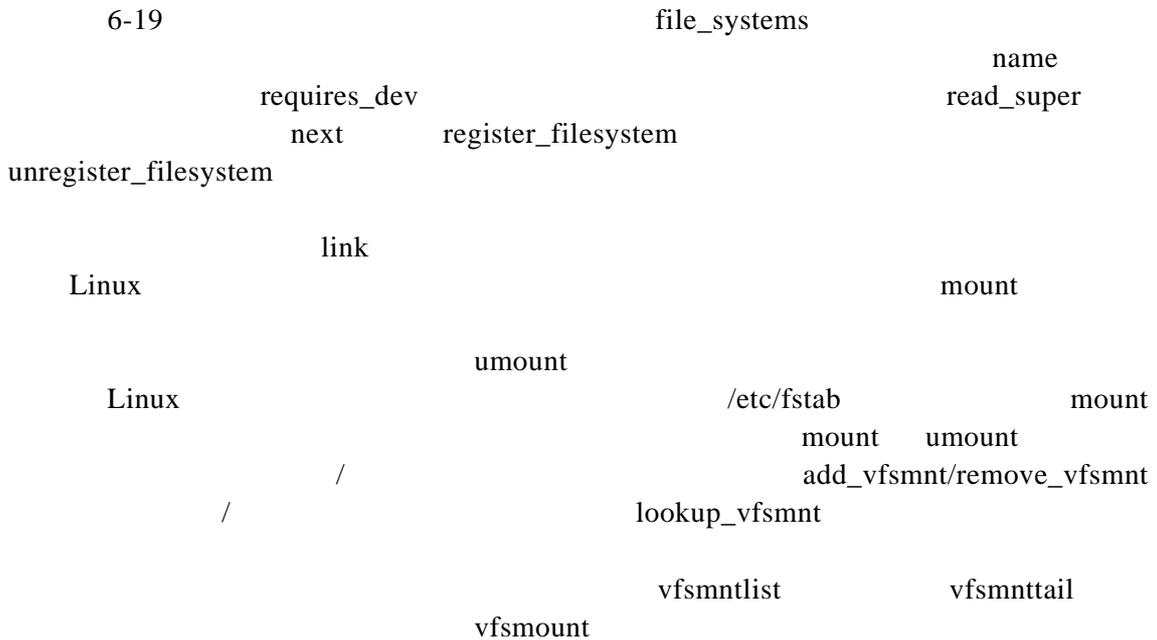
Linux EXT EXT2 MINIX UMSDOS NCP
 ISO9660 HPFS MSDOS NTFS XIA VFAT PROC NFS SMB SYSV AFFS
 UFS

Linux VFS(Virtual File System)
 VFS Linux
 VFS Linux
 Linux I/O EXT2 Linux
 inode hash
 inode

6.5.2 Linux



6-19 Linux



```

static struct vsmount vfsmntlist = (static struct vsmount )NULL; /* */
static struct vsmount *vfsmnttail = (static struct vsmount *)NULL; /* */
static struct vsmount *mru_vfsmnt = (static struct vsmount *)NULL; /* */

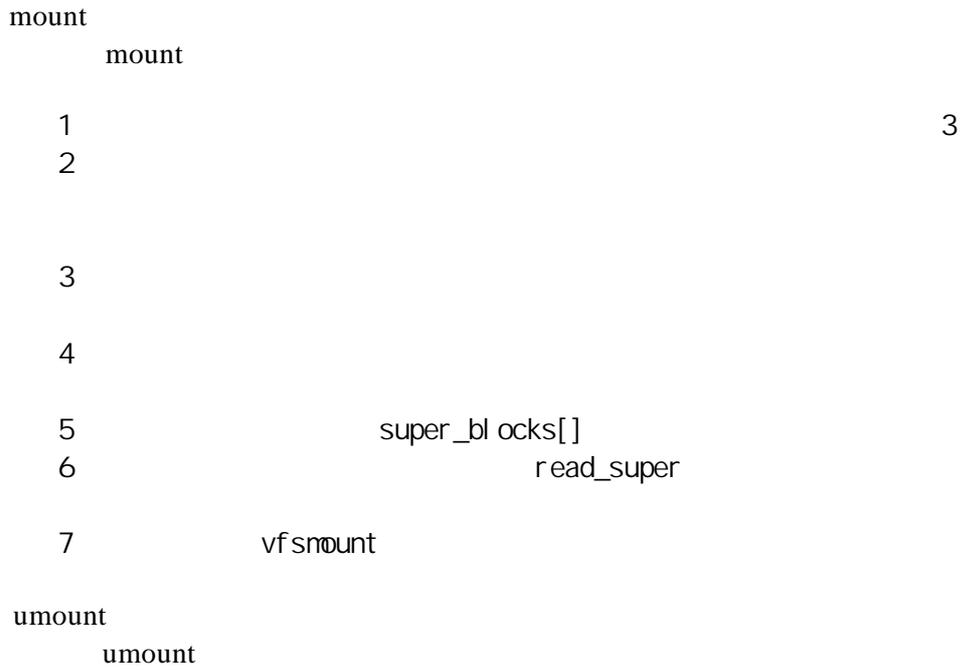
struct vsmount {
    kdev_t mnt_dev; /* */
    char *mnt_devname; /* /dev/hda1 */
}

```

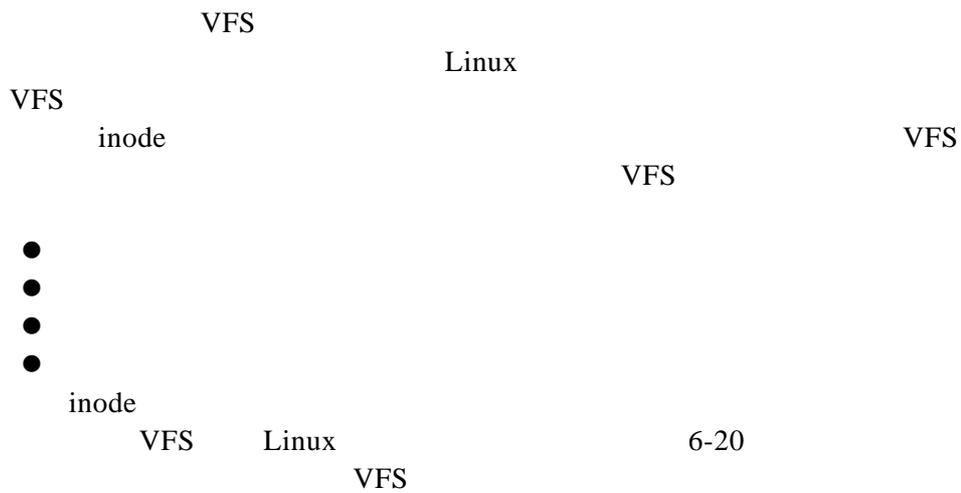
```

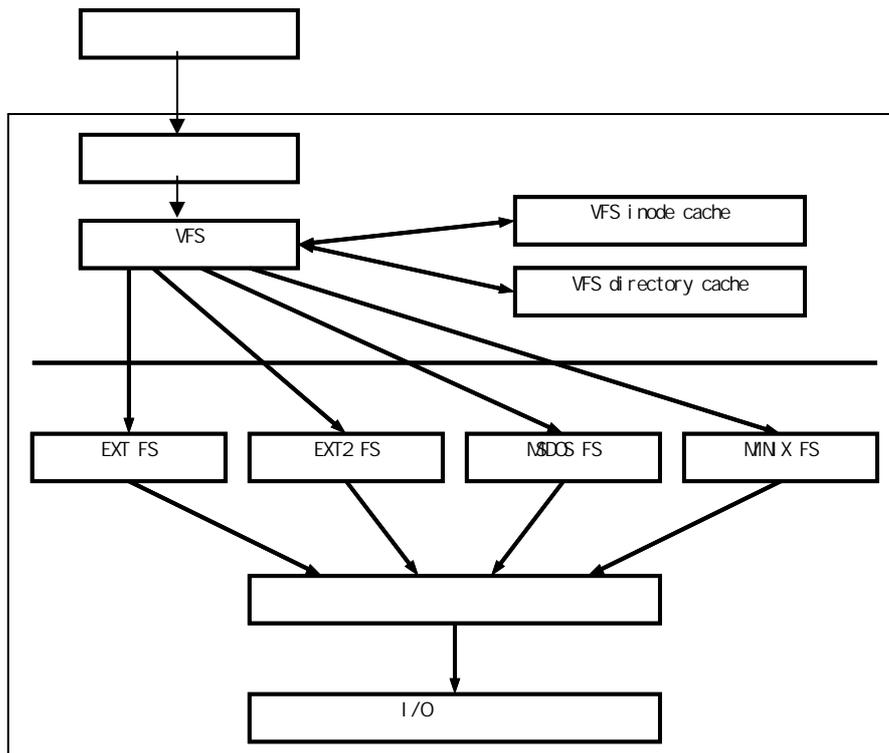
char *mnt_dirname;          /* */
unsigned int mnt_flags;     /* ro */
struct semaphore mnt_sem;   /* */
struct super_block *mnt_sb; /* */
struct file *mnt_quotas[MAXQUOTAS]; /* */
time_t mnt_iexp[MAXQUOTAS]; /* expiretime for inodes */
time_t mnt_bexp[MAXQUOTAS]; /* expiretime for blocks */
struct vfsmount *mnt_next;  /* */
};

```



6.5.3 VFS





6-20 Linux

VFS

VFS

inode

VFS

read_super

```

struct super_block {
    struct list-head s-list;
    kdev_t   s_dev; /* */
    unsigned long s_blocksize; /* */
    unsigned char s_blocksize_bits; /* 2 */
    unsigned char s_lock; /* */
    unsigned char s_rd_only; /* */
    unsigned char s_dirt; /* */
    struct file_system_type *s_type; /* */
    struct super_operations *s_op; /* */
    struct dquot_operations dq_op;
    unsigned long s_flags,s_magic,s_time;
    struct dentry *s-root;
    struct wait_queue *s_wait; /* */
    struct inode *s-ibasket;
    short int s-ibasket-max,s-ibasket-count;
    struct list-head s-dirty;
    union { /* */
        struct minix-sb-info minix-sb;
        struct ext2-sb-info ext2-sb;
        struct hpfs-sb-info hpfs-sb;
    };
};

```

```

    struct ntfs-sb-info ntfs-sb;
    struct msdos-sb-info msdos-sb;
    struct isofs-sb-info isofs-sb;
    struct nfs-sb-info nfs-sb;
    struct sysv-sb-info sysv-sb;
    struct ufs-sb-info ufs-sb;
    struct romfs-sb-info romfs-sb;
    struct smb-sb-info smb-sb;
    struct hfs-sb-info hfs-sb;
    struct adfs-sb-info adfs-sb;
    struct qnx4-sb-info qnx4-sb;
    void *generic-sbp;
} u;
};

```

	super_block.u	VFS		inode
	inode	s-mounted		Linux
				VFS inode
inode_cache				
	inode			

```

struct inode {
    struct list-head i-hash,i-list,i-dentry;
    unsigned long i-ino;
    unsigned int i-count;
    kdev_t i_dev; /* */
    umode_t i_mode; /* */
    nlink_t i_nlink; /* link */
    uid_t i_uid;
    gid_t i_gid;
    kdev_t i_rdev; /* */
    off_t i_size; /* */
    time_t i_atime i_mtime i_ctime;
    unsigned long i_blksize; /* (1KB) */
    unsigned long i_blocks; /* */
    unsigned long i_version;
    unsigned long i_nrpages; /* */
    struct semaphore i_sem i-atomic-write;
    struct inode_operations *i_op; /* */
    struct super_block *i_sb; /* VFS */
    struct wait_queue *i_wait; /* */
    struct file_lock *i_flock; /* */
    struct vm_area_struct *i_mmap;
    struct page *i_pages; /* */
    struct dquot *i_dquot[MAXQUOTAS];
    struct inode i_next, i_prev; /*inode */
    struct inode i_hash_next, i_hash_prev; /*inode cache */
    struct inode *i_bound_to, *i_bound_by;
    struct inode *i_mount; /* inode */
}

```

```

unsigned long i_count; /* 0 */
unsigned short i_flags, i_writecount;
unsigned char i_lock; /* inode */
unsigned char i_dirt; /* inode */
unsigned char i_pipe i_sock i_seek i_update i_condemned;
union { /* inode */
    struct pipe-inode-info pipe-i;
    struct minix-inode-info minix-i;
    struct ext2-inode-info ext2-i;
    struct hpfs-inode-info hpfs-i;
    struct ntfs-inode-info ntfs-i;
    struct msdos-inode-info msdos-i;
    struct umsdos-inode-info umsdos-i;
    struct isofs-inode-info isofs-i;
    struct nfs-inode-info nfs-i;
    struct sysv-inode-info sysv-i;
    struct affs-inode-info affs-i;
    struct ufs-inode-info ufs-i;
    struct romfs-inode-info romfs-i;
    struct coda-inode-info coda-i;
    struct smb-inode-info smb-i;
    struct adfs-inode-info adfs-i;
    struct qnx4-inode-info qnx4-i;
    struct socket socket-i;
    void *generic-ip;
} u;
};

```

```

inode i_prev i_next first_inode inode i_dev
i_ino i_count 0 inode
inode first_inode inode VFS
grow_inodes
inode first_inode first_inode VFS
insert_inode_free() remove_inode_free() put_last_free() insert_inode_hash()
remove_inode_hash() clear_inode() get_empty_inode() lock_inode() unlock_inode()
write_inode()

```

6.5.4

Linux VFS

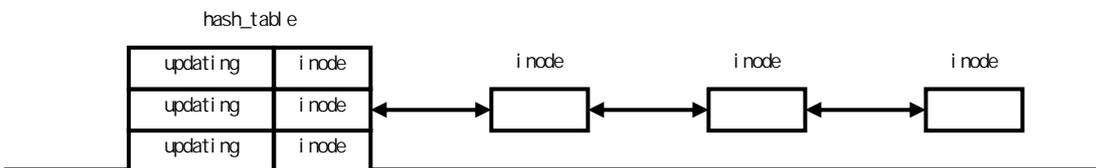
cache

1 VFS inode cache

```

VFS inode VFS inode cache hash inode
VFS inode hash h 6-21
hash_table[h] i_hash_next i_hash_prev cache
inode inode i_count 1 inode

```



2 VFS directory cache

```

inode
Linux      inode      VFS directory cache      Linux
      directory cache
namei()
struct hash_list { struct dir_cache_entry *next *prev; };
struct dir_cache_entry {
    struct hash_list h;
    kdev_t dc_dev;
    unsigned long dir; /* inode */
    unsigned long version;
    unsigned long ino; /* inode */
    unsigned char name_len;
    char name[DCACHE_NAME_LEN];
    struct dir_cache_entry *lru_head;
    struct dir_cache_entry *next_lru *prev_lru;
}
VFS directory cache level1_cache level2_cache cache
level1_head level2_head level1_cache level2_cache
dir_cache_entry.next_lru dir_cache_entry.prev_lru 128
      LRU level1_cache
level2_cache level1_cache
level2_cache Linux level1_cache level2_cache
hash

```

3 buffer cache

```

Linux      buffer cache buffer
cache
      buffer cache

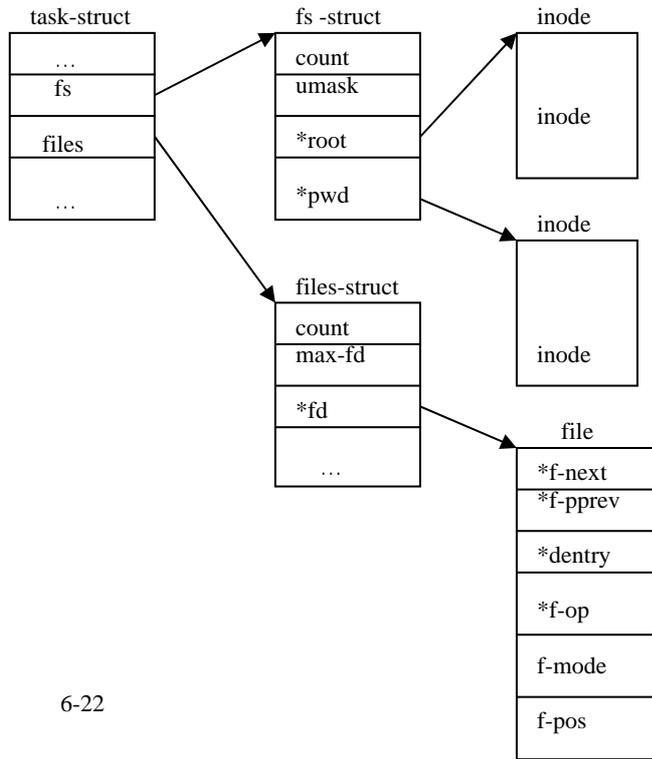
```

```

Linux struct buffer_head
buffer_head
      buffer_head
      buffer cache hash buffer
● hash
      static struct buffer_head ** hash_table;
●
      4 lru_list[0]
lru_list[1] lru_list[2]
inode lru_list[3]
      buffer_head
      static struct buffer_head *lru_list[NR_LIST];
●
      buffer 512 1024 2048 4096 8192 5
      buffer_head
      static struct buffer_head *free_list[NR_SIZES];
●
      buffer_head

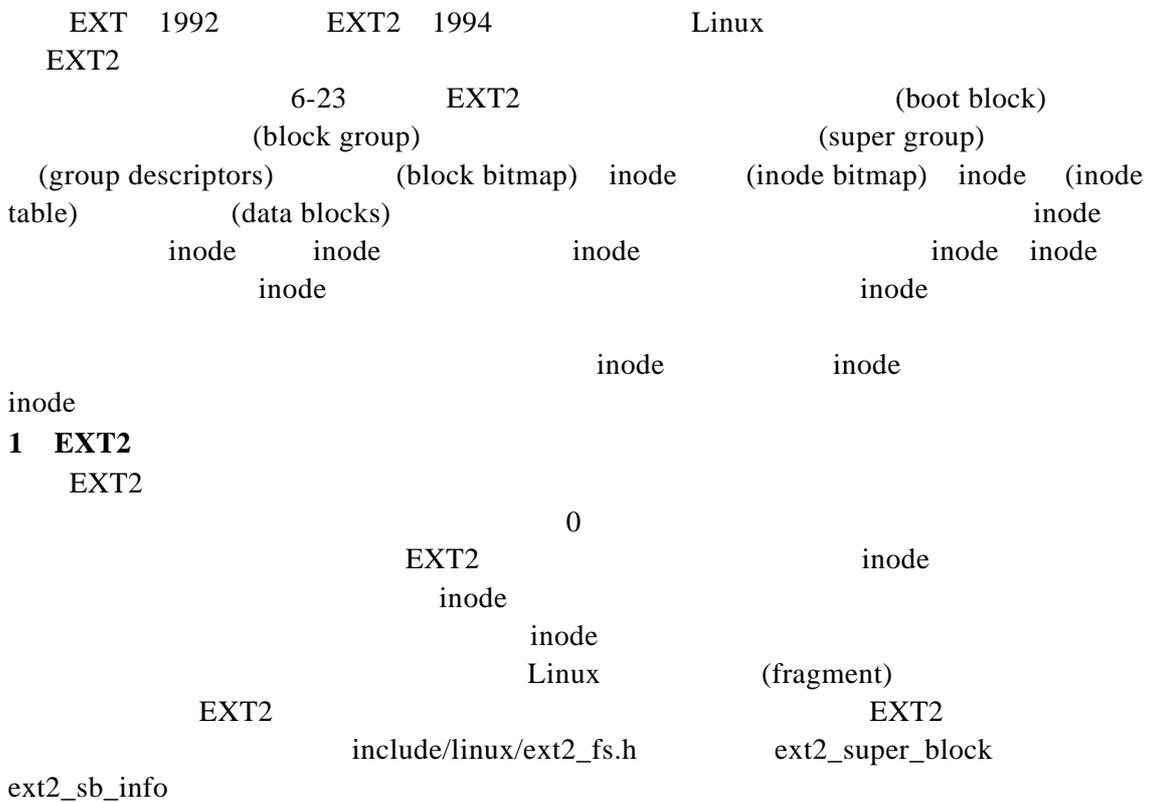
```

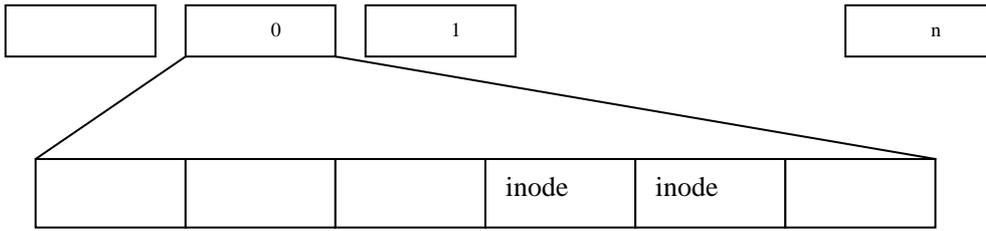

Linux
include/linux/fs.h



6-22

6.5.6 EXT2





6-23 EXT2

2 EXT2

```

                                     inode      inode
include/linux/ext2_fs.h      inode      ext2_group_desc

```

3 EXT2 inode

```

inode      inode
inode      inode
/ / /      link
                                     include/linux/ext2_fs.h      ext2_inode

```

4 EXT2

```

                                     " ."      " .."

struct ext2_dir_entry {
    _u32  inode;      /*      inode */
    _u16  rec_len;    /*      */
    _u16  name_len;  /*      */
    char  name[EXT2_NAME_LEN]; /*      */
}

```

5

```

                                     "      "
Linux      defrag(defragmentation program)
                                     EXT2
●
64
●      EXT2      8      ( EXT2_PREALLOCATE )

                                     EXT2      inode      ext2_inode_info
prealloc_block      prealloc_count

```

6.6 Windows 2000/XP

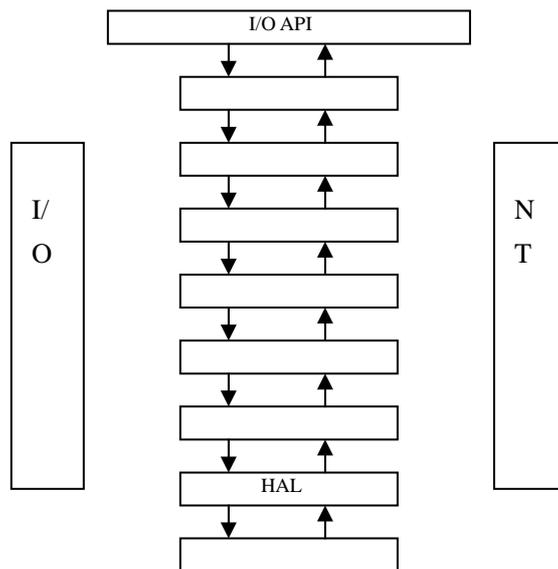
6.6.1 Windows 2000/XP

Windows 2000/XP	FAT	FAT	DOS
Windows3.x	Windows95		
2GB	FAT	2^{12} 2^{16}	FAT
Windows9x	Windows Me	FAT	FAT32
FAT16		32	4GB
FAT	Windows NT	FAT32	Windows98
FAT32		CDFS	NT
2000/XP	FAT32		Windows
HPFS			UDF
Microsoft	Windows NT	FAT	
NTFS(New Technology File System)	NTFS		
NT4	NTFS	/	Windows 2000/XP
	NTFS5	NT4	NTFS4
●	NTFS		
●	NTFS	(ACL)
●	Windows2000		EFS(Encrpyting File System)
●	NTFS		RAID
●	NTFS	64	
●	NTFS		NTFS
●	Unicode	NTFS	16
	Unicode		Unicode
●	Unicode	NTFS	255
		ID	(
			Windows 2000
			OLE
			NTFS
●	Windows2000		
●	NTFS		
●	NTFS		

- Windows2000 NTFS
- Windows2000
- Windows2000 50% 40%
- Windows2000 NTFS ID ID (OLE)
- POSIX POSIX
- Windows2000 " (HSM)"
- Windows 2000/XP (DFS) Windows NT4 DFS
- 2000/XP Windows 2000/XP DFS
- " " DFS
- DFS

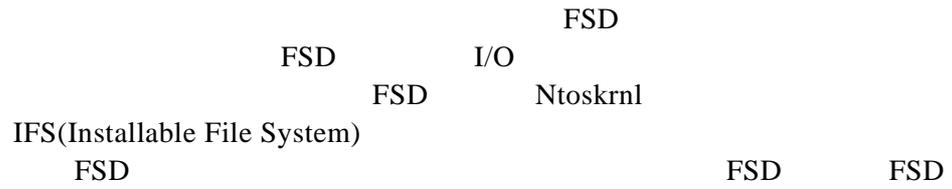
6.6.2 Windows2000/XP FSD

Windowsx2000/XP I/O I/O
6-24

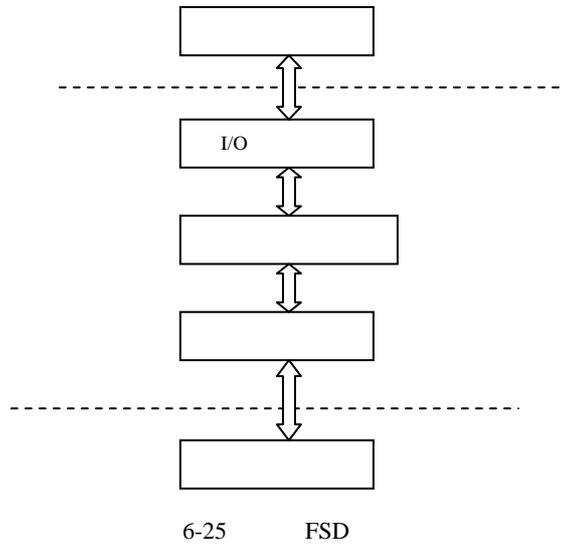


6-24 Windows

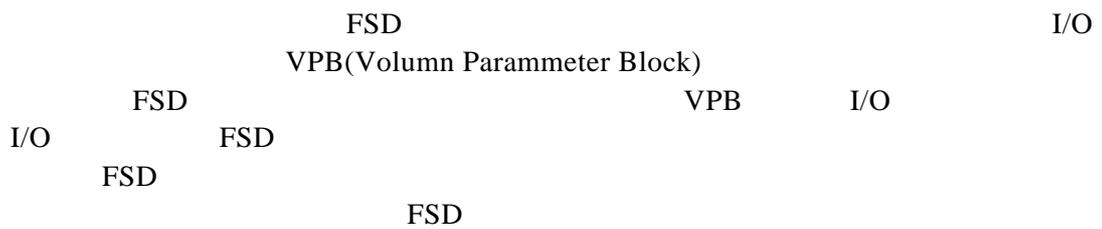
- I/O
 - I/O
 - FSD(File System Driver)
(NTFS)
 - I/O
- API



1. FSD
- | | | | | | | |
|------|-----|----------|-------------|----------|----------|---------|
| | FSD | Ntfs.sys | Fastfat.sys | Udfs.sys | CDfs.sys | Raw FSD |
| 6-25 | FSD | I/O | | | | I/O |
| FSD | | | | | | |



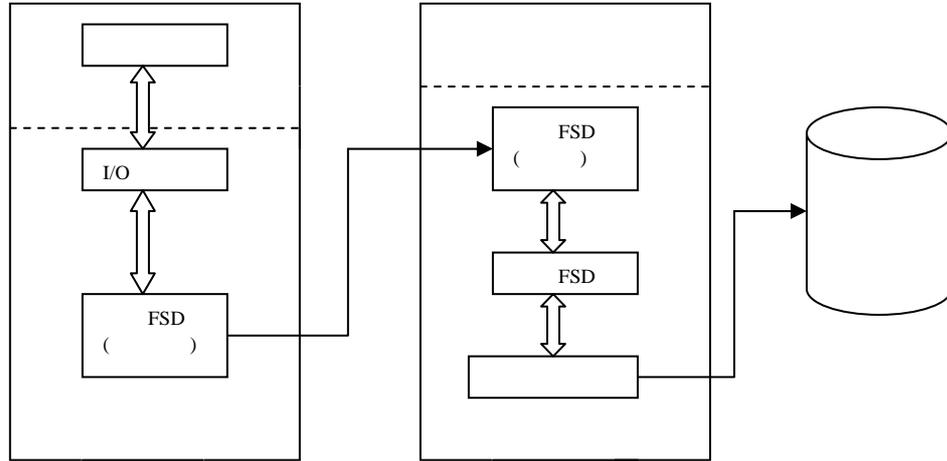
Windows2000/XP



2. FSD
FSD

FSD

FSD FSD I/O
FSD FSD 6-26 FSD



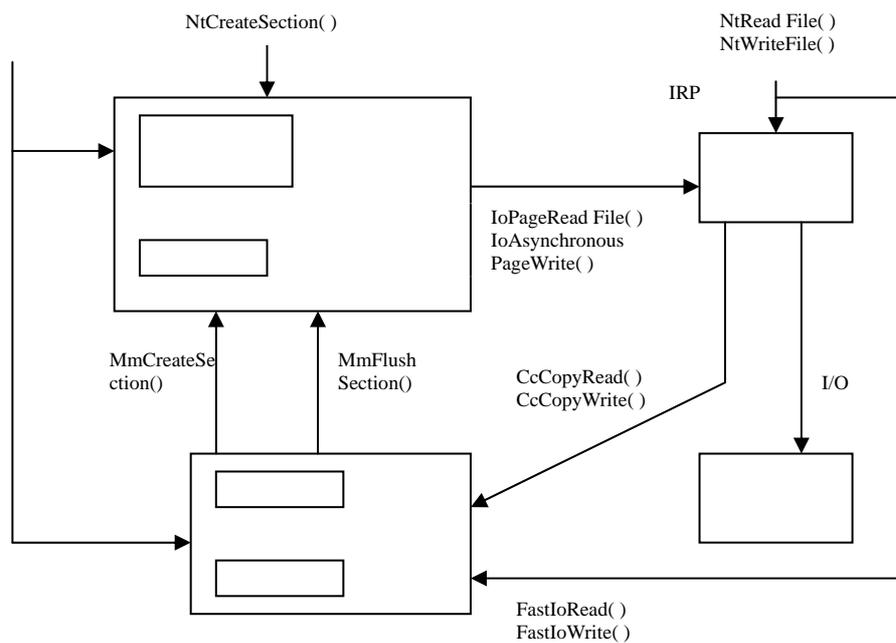
6-26 FSD

Windows2000/XP (LANMan Redirector) FSD LANMan FSD LANMan (LANMan Server) (LANMan)

/

CIFS(Common Internet File System)
3. FSD
Windows

FSD 6-27



6-27 FSD

I/O Win32 I/O CreateFile ReadFile WriteFile

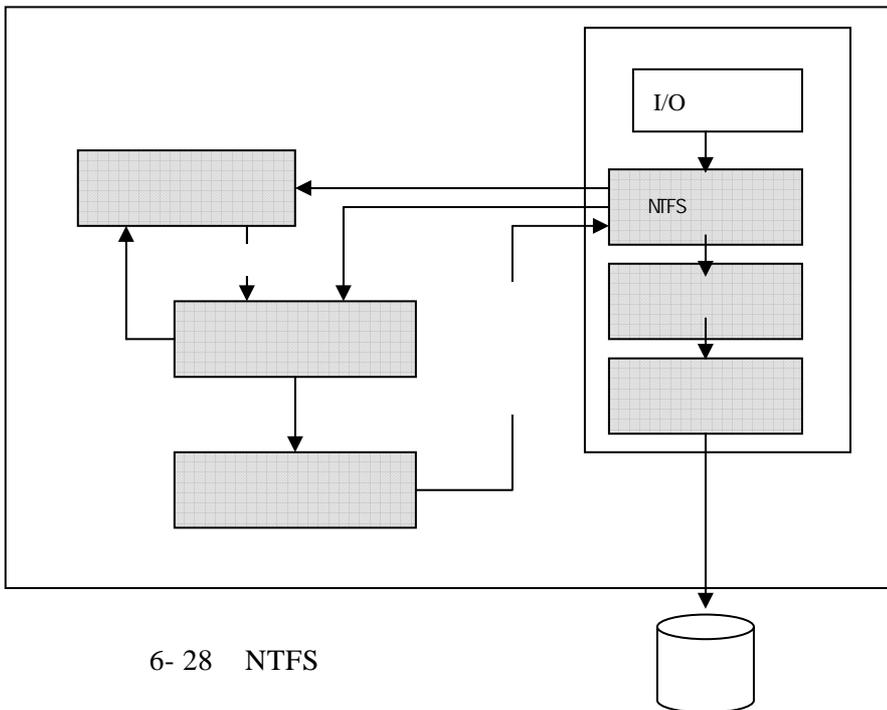
MmFlushSection IoAsynchronousPageWrite MmFlushSection FSD

IoAsynchronousPageWrite IRP IRP
 FSD
 I/O
 MmAccessFault IoPageRead IRP

6.6.3 NTFS

Windows2000/XP NTFS FAT HPFS POSIX
 I/O

I/O
 Window2000/XP
 6-28 Windows2000/XP I/O
 NTFS
 I/O I/O I/O
 NTFS



6-28 NTFS
 LFS(log file server)

NTFS
 NTFS
 (Windows2000)
 Windows2000

Windows2000
 NTFS
 " " (lazy writer) I/O
)

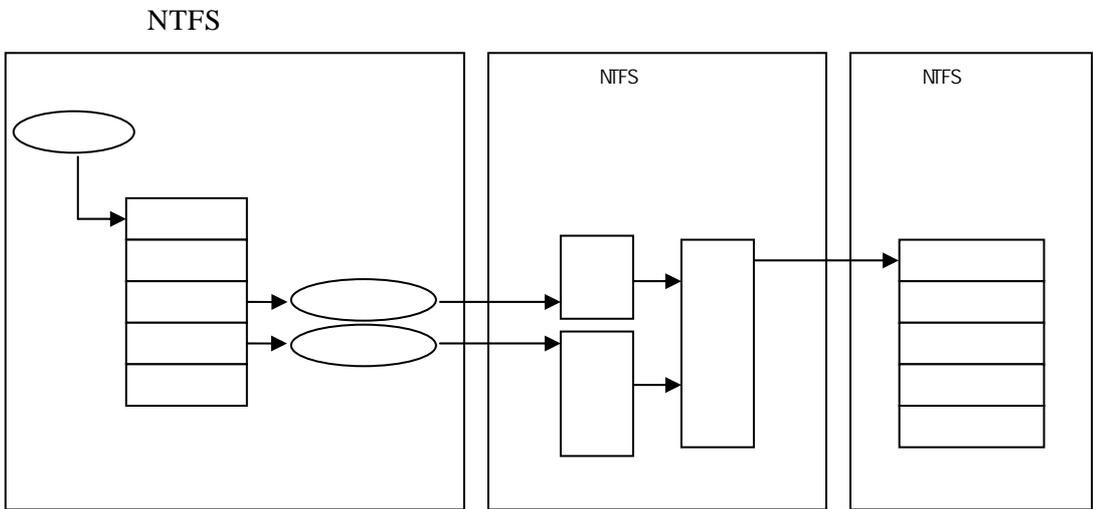
NTFS
 Windows2000/XP
 Windows 2000/XP
 I/O NTFS

I/O NTFS

6-29 NTFS I/O

Block SCB Stream Control SCB

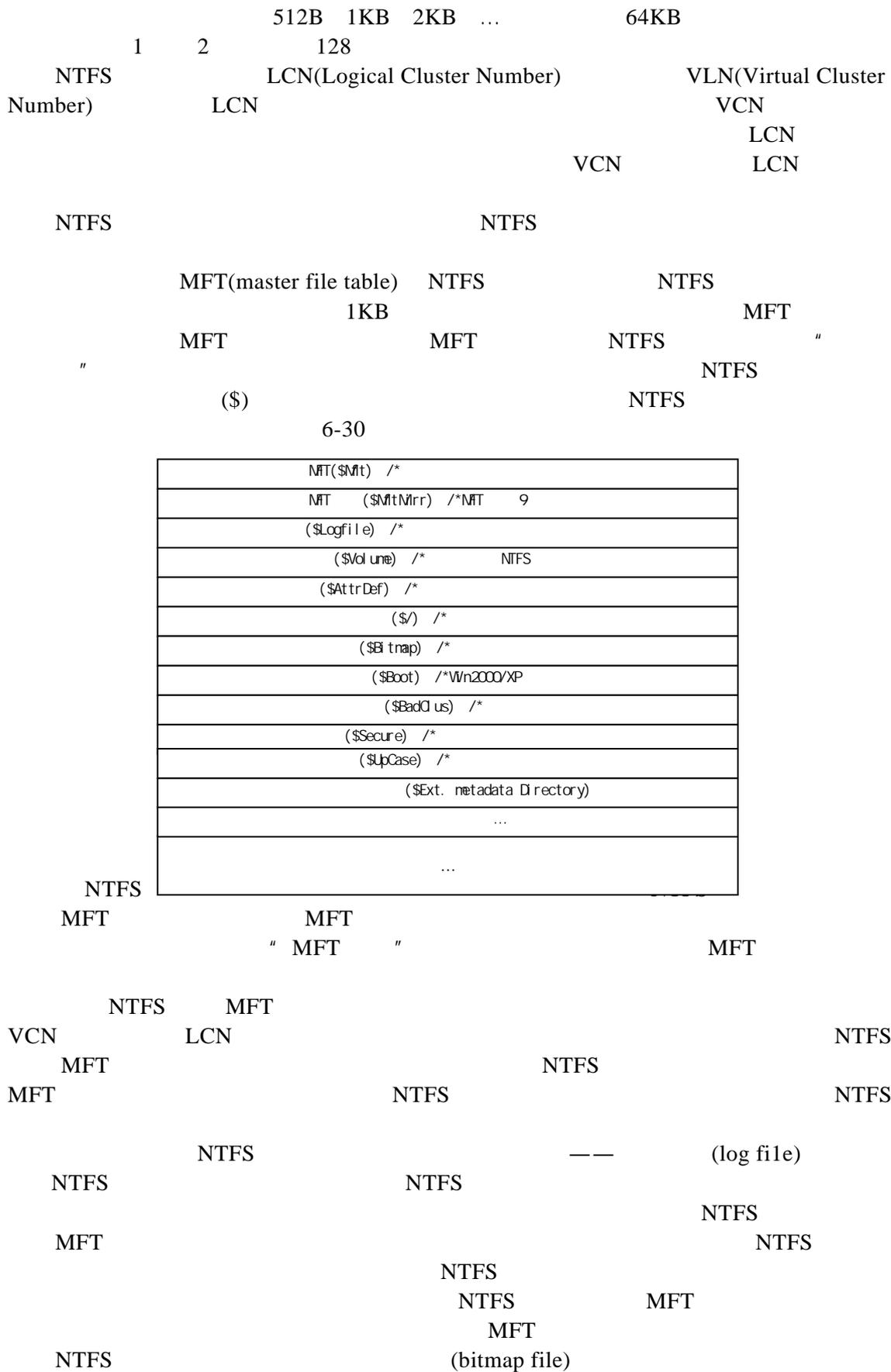
SCB File Control Block FCB FCB MFT



6-29 NTFS

6.6.4 NTFS

NTFS
 NTFS FAT NTFS
 NTFS



(bootfile) Windows2000 Format

NTFS NTFS

Windows 2000 "

" I/O

NTFS " " (bad-cluster)

" " (volume file) NTFS

NTFS Chkdsk (attribute definition table)

NTFS " " 64

MFT 1 (

MFT 1)

NTFS NTFS

NTFS / (

)

NTFS () () NTFS

NTFS FAT 255 Unicode

5IX-6.4(e)]TJ/TT33

" " (index buffers)

6.6.5 NTFS

NTFS

(logging)

NTFS

"

(Write-ahead logging)"

●

NTFS
area)

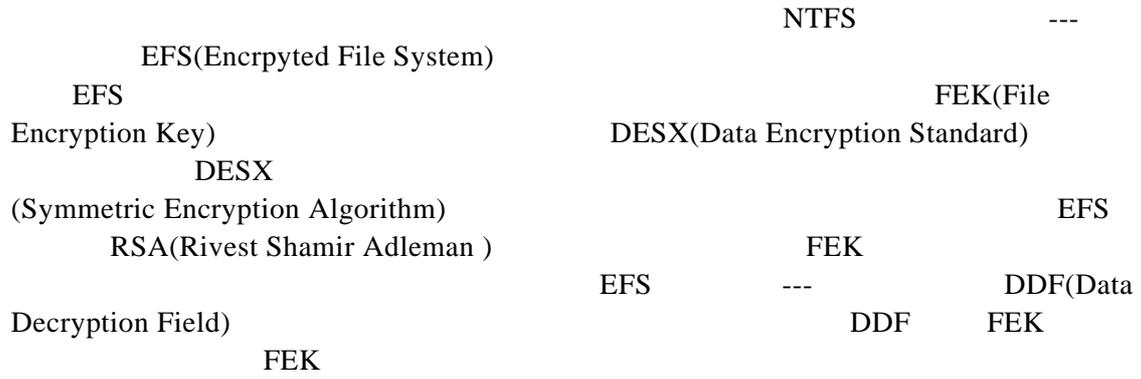
LFS(Log File Service)
LFS

" NTFS

NTFS
LFS

(restart

-
-



6.7

IBM

i-node

()

13. Windows 2000 NTFS

14.

(1)

(2)

(3)

15.

3

?

?

16.

17.

?

?

18.

hash

?

19.

hash

?

20.

hash

21.

?

22.

?

23.

'

'

?

?

24.

?

?

25.

UNIX

26. UNIX/Linux

27. UNIX/Linux

i

28. UNIX/Linux

0

29. UNIX/Linux

30.

OPEN CLOSE

1

2

OPEN COLSE

31.

?

32.

rename

33.

36

UNIX/Linux

2

1

1

2

1

i

j

2

B=

R=

F=

(

)

F

3

500

32

(5)

(51)

(6)

CH7

7.1

-
-
-

(Authorization) (Encryption) (Authentication)
1985 " (Audit) "

- D
- C1
- C2
- B1
- B2
- B3
- A1

C2 OSF/1 B1 UNIX Ware 2.1 B2 DOS D Windows NT Solaris

7.2

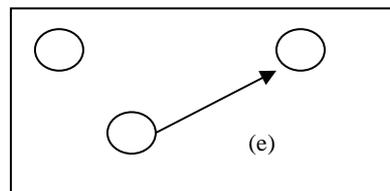
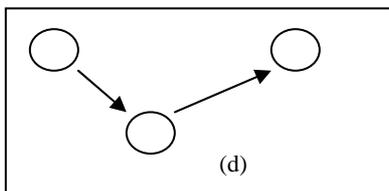
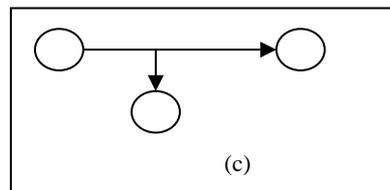
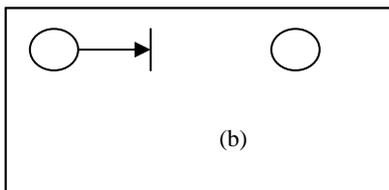
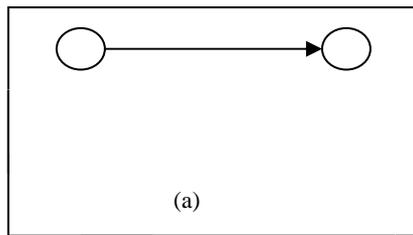
()

- confidentiality

- integrity
- availability
- authenticity

7-1a

-
-
-
-
-



7-1

			/

1

2

3

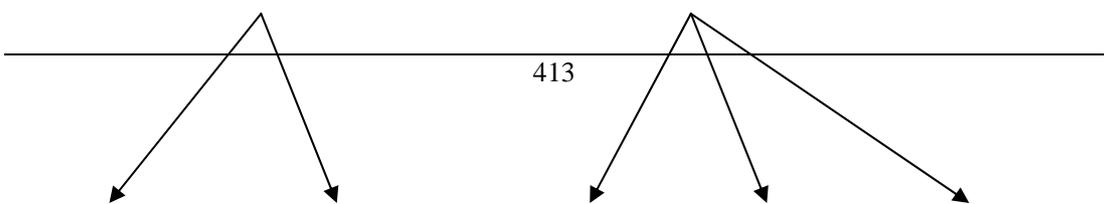
A B C D

E

A B C D E

4

7-2



7.3

7.3.1

I/O

-
-
-
-
-
-

7.3.2

7.3.3

ID

ID

ID

ID/

ID/

ID

ID

7.3.4

(access matrix)

- (subject)
- (object)
- (access authority)

7.4

7.4.1

cracker

Anderson

hacker

- masquerader
- misfeasor
- clandestine user

•

•

1

2

3

4

5

6

7

8

6

7

8

7.4.2

ID

ID ID

•

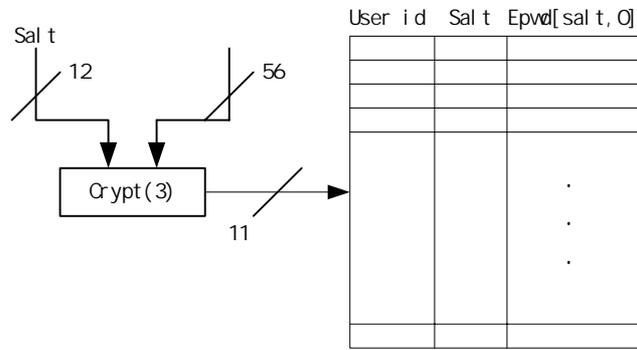
ID

ID

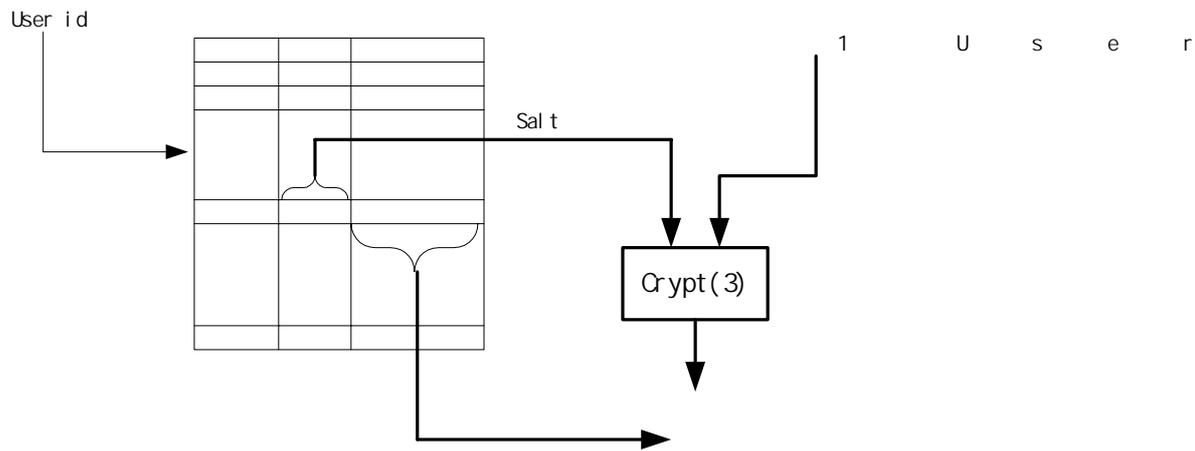
• ID

• ID

ID



(a)



UNIX
 salt
 25
 DES
 25
 UNIX
 (Purdue)
 7000
 3%
 54
 6
 8

1	55	0.004
2	87	0.006
3	212	0.02
4	449	0.03
5	1260	0.09
6	3035	0.22
7	2917	0.21
8	5772	0.42
total	13787	1.0

2

Spafford,E.

●

UNIX

●

●

3

8

4

•

•

•

•

8

" "

" "

7.4.3

•

•

•

Porras,P.Stat

•

1

2

•

Dorothy Denning

Smith

GAME

<library>

COPY GAME.EXE TO <LIBRARY> GAME.EXE

Smith		<library>copy.exe	0	CPU=00002	11058721678
-------	--	-------------------	---	-----------	-------------

Smith		<smith>GAME.exe	0	RECORDS=0	11058721679
-------	--	-----------------	---	-----------	-------------

Smith		<library>copy.exe	write-viol	RECORDS=0	11058721680
-------	--	-------------------	------------	-----------	-------------

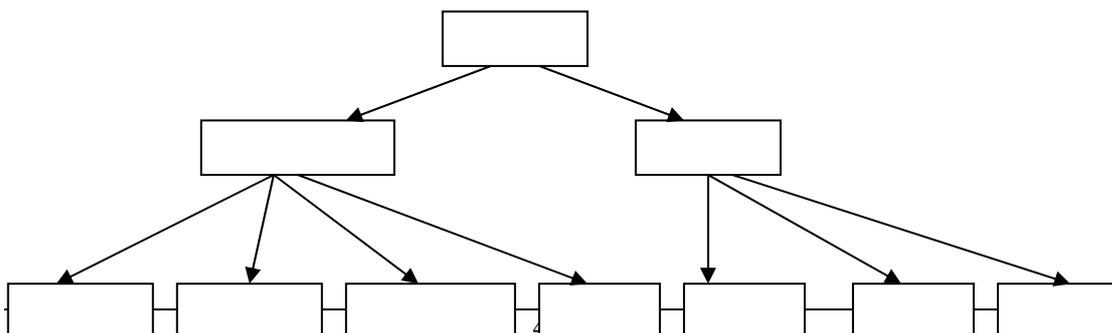
Smith <Library>

7.5

7.5.1

malware

malware



7-4

1 trap door

7.5.2

-
-
-
-

4

7.5.3

-
-
-
-
-
-

word

Microsoft Excel

7.5.4

-
-
-

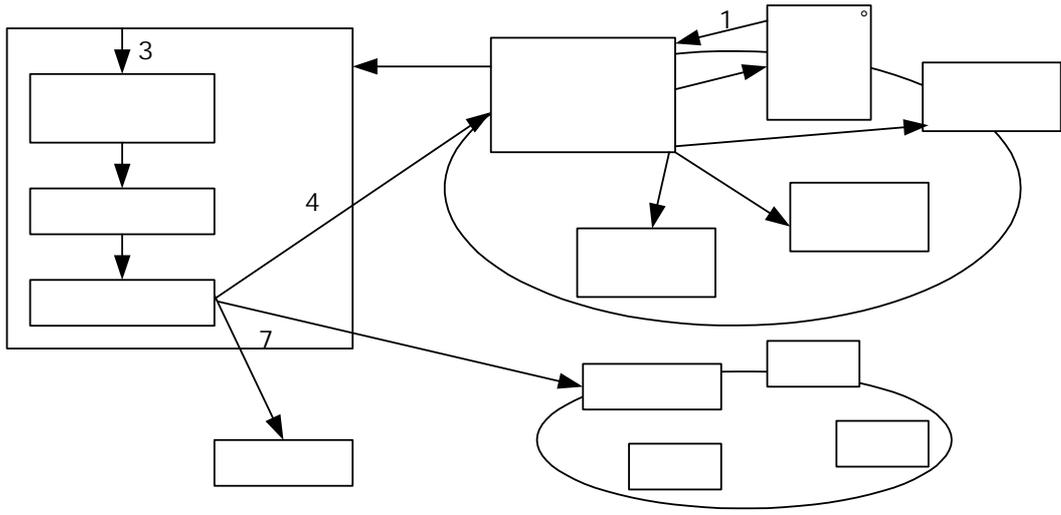
1

GD(General Decrypt)

● CPU GD GD

●
●

GD



7.6

7.6.1

policy

mechanism

?

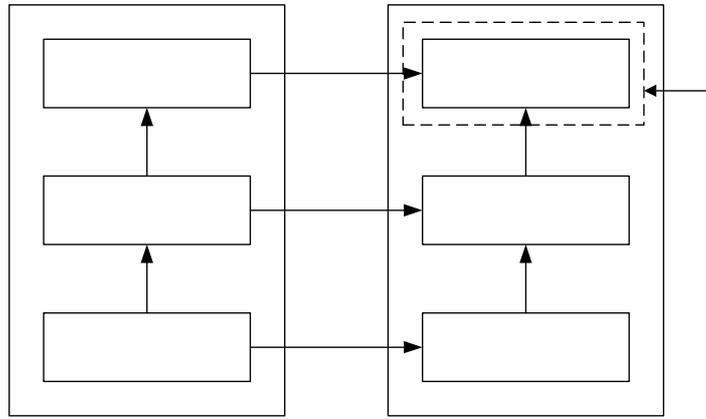
" "

—

-
-
-
-

1

7-6



7-6

7-7

7-7

a

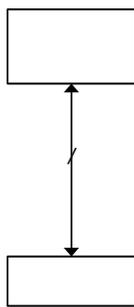
b

a

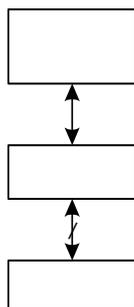
b

a

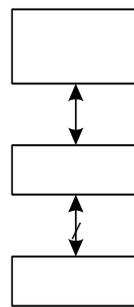
b



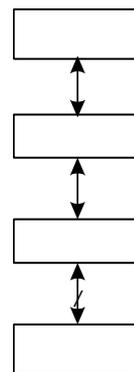
(a)



(b)



(c)



(d)

7-7

7-7

c

d

c

b

c

b

c

c

b
d

2

Bell&Lapadula
3

()

(information flow control)

D.E.Denning

Honeywell

Secure Ada Target

" "

4

(1)

(2)

(3)

(4)

(5)

(6)

()

(multilevel security)

•

•

• a

• b

• a

• b

S O A

S				
O				
sclass(a)=	s			
oclass(o)=	o			
A(s,o)=				
r =	s	o		
w =	s	o		
r,w =	s	o	o	
=	s	o	o	
content o =	o			
subj=				

a b

contents(o)

subj

subj

state={S,O,sclass,oclass,A,contents,subj}

a b

		a	S	o	O
i f r	A(s, o),	then	scl ass(s)	ocl ass(o)	
i f w	A(s, o),	then	ocl ass(s)	scl ass(o)	

A

Creat_object(o,c)	o
Set_access(s,o,modes)	s o
Creat/Change_object(o,c)	o c o
Write_object(o,d)	contents(o)
Copy_object(from,to)	form to
Append_data(o,d)	d o

" "

```

FUNCTION 1:Create_object(o,c)
  If o O
  Then O = O {o}
  And
  ocl ass(o) = c
  and
  for all s S, A(s, o)=
FUNCTION 2:Set_access(s,o,modes)
  If o O and s S
  And if {[r modes and scl ass(s) ocl ass(o)] or r modes}
  And
  {[w modes and ocl ass(s) scl ass(o)] or w modes}
  then A(s,o)=modes

```

-
-
-
-

if conditions then
else conditions

Creat_object o
c
Set_access A
A(a,x) x o

	<pre> Create_object for all s S, A(s, o) = S_i </pre>	<pre> mode r Create_object </pre>
mode	<pre> sclass S_i oclass(o) </pre>	
	<pre> sclass(s) </pre>	
	<pre> Create_object Set_access </pre>	
	<pre> Create_object Set_access </pre>	

$S_0, O_0, Sclass_0, Oclass_0, Contents_0, Subj_0$

1 $S_0 = \dots, O_0 = \dots$

S

2 For all s $S_0, o \in Q_0$
 $Sclass_0(s) = c_0$
 $Oclass_0(o) = c_0$
 $A_0(s, o) = \{r, w\}$

A_0

-
-
-

1
Create_object

FUNCTION 3 Create/Change_object(o,c)
 oclass(o)=c; and
 if $o \in O$ then $o = O \{o\}$; and
 for all $s \in S, A(s,o) =$

- **c**
Create/change_object

1 For all $o \in Q$ oclass(o) = oclass(o)

2

" "

- **d**
subj
2 For all $o \in O$
 if $r = A(\text{subj}, o)$
 then for all $s \in S, A(s, o) = A(s, r)$
 Set_access

Set_access

FUNCTION 4 Write_object(o,d)
 If $o \in O$ and $w \in A(\text{subj}, o)$
 Then $\text{contents}(o)=d$

Write_object

$w \in A(\text{subj}, o)$
 b

/

● e

3 For all $o \in Q$
 if $w \in A(\text{subj}, o)$
 then $\text{contents}(o)=\text{contents}(o)$

/

5 Bell&LaPadula

Bell&LaPadula

David Bell Leonard LaPadula

Case Western Reserve University

BLP

BLP

BLP

" "

" " " " " " " "
" >" " >" " "
" " " "
" "

-
-
-

-
-

SYSTEM HIGH

SYSTEM LOW

SYSTEM HIGH
SYSTEM LOW

BLP

BLP 20

-
-
-

BLP

/

/

Bell
Multics

LaPadula

Multics

BLP

6

" " " "

1

-
-
-
-
-
-
-

I/O

—

CPU

Saltzer Schroeder

2

"

() ()

•

•

•

•

•

3

5:00

8:00

4

5

7.6.2

UNIX

UNIX

UNIX

" "

UNIX

"

"

" "

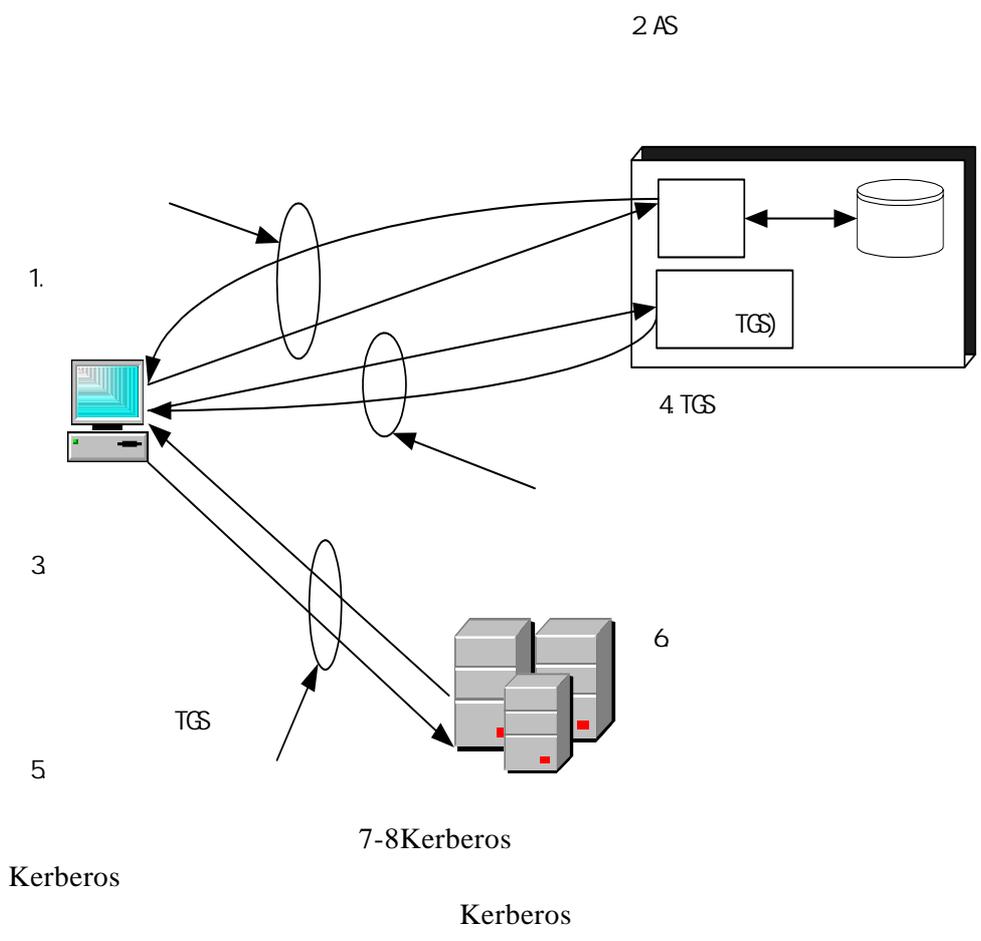
1

2

" Morris
 11.c 99
 name@host finger UNIX UNIX " finger
 finger Finger finger .plan
 .plan
 finger
 Morris finger finger
 finger finger
 finger shell

3 Kerberos

Kerberos MIT 1980
 " Athena" " — "
 Kerberos Kerberos UNIX TCP/IP
 Kerberos
 ● Kerberos
 ● Kerberos Kerberos Kerberos
 ● Kerberos
 ●



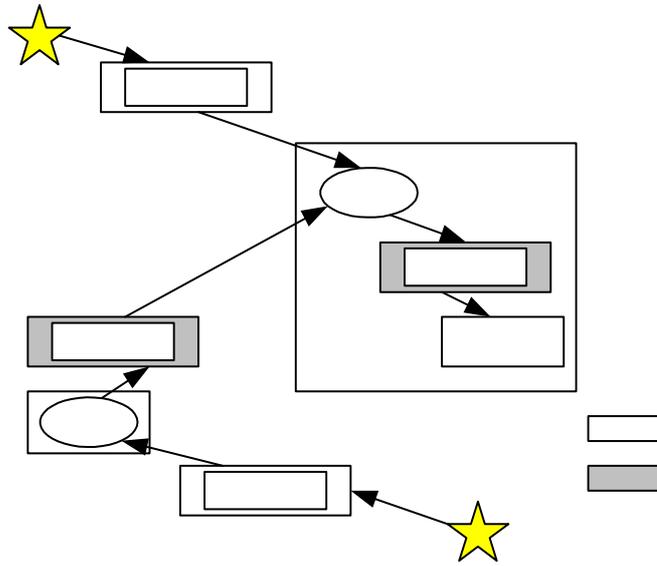
- 7-8
- 1
 - 2
 - 3
 - 4
 - 5
 - 6

" " _ _

7.6.3

1

7-9



7-9

shell

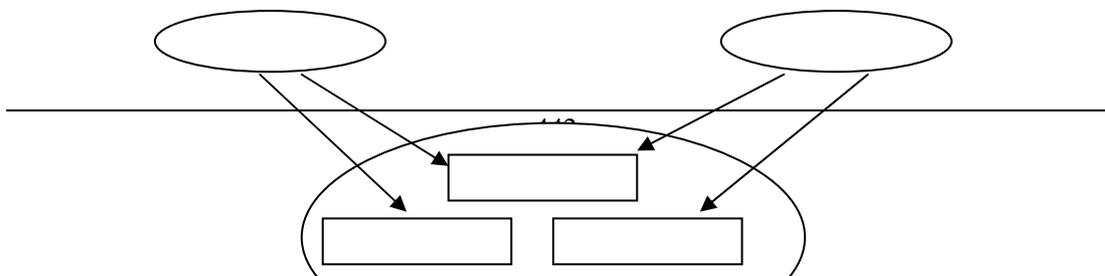
7-9

A W X Y Z 7-10

B W C W

B Y C Y C X B X A

Z



7-10



subjects = {S₁ S₂ S₃}
 objects = subjects {F1 F2 D1 D2}
 F1 F2 D1 D2 7-13
 S₁ S₂ S₃
 F1 S₁ * * S₂ F1 S₃ F1
 S₂ F2
 S₂ update F2 F2
 A[S₂ F2] A[S₂ F2]
 S₂ F2
 S₂ execute F2 F2
 A[S₂ F2] A[S₂ F2]

S ₁	S ₂	S ₃	F1	F2	D1	D2
			*			
			*			
						*

7-13

3

Graham Denning 7-14
 S₀ S₃ D2
 A[S₀ O] A[S₃, D2] A[S, O] 7-13

S ₀
1 transfer{ *} to [S, O] * A[S ₀ , O] A[S, O] = A[S, O] { *}
2 grant{ *} to [S, O] owner A[S ₀ , O] A[S, O] = A[S, O] { *}
3 delete from [S, O] control A[S ₀ , S] or A[S, O] = A[S, O] - { }
owner A[S ₀ , S]

7.14

* S₀ O S 7-14
 S₁ F1 O S₂ S₃ S₁
 2 S₀ O S₀

O

Graham Denning[1972] 7-14

●

X

S
X

●

●

" "

" "

" "

" "

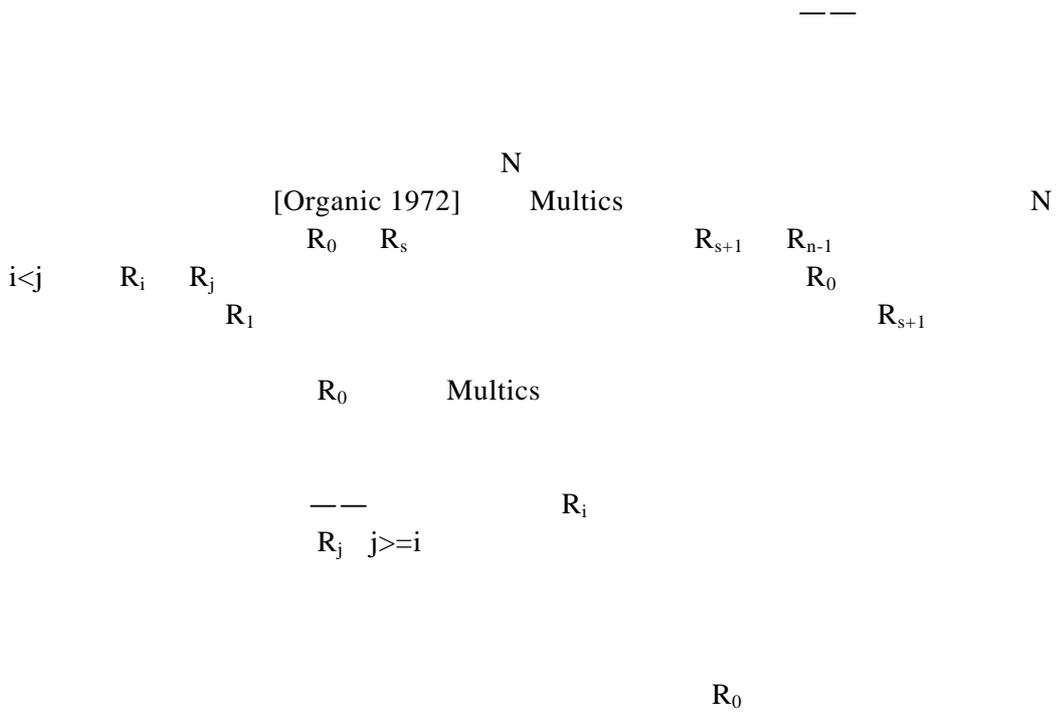
" "

●

4

-
-
-
-
-

(1)



1

0

2

Intel 80486

0

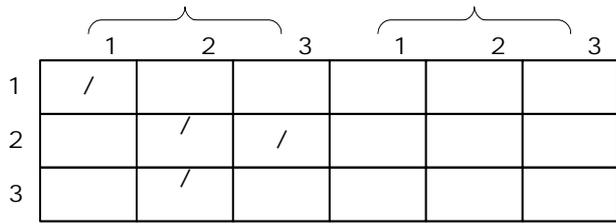
2

VAX/VMS

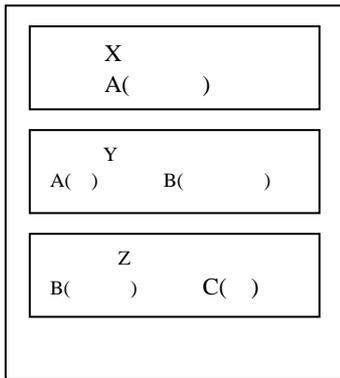
7-15VAX/VMS

7-15

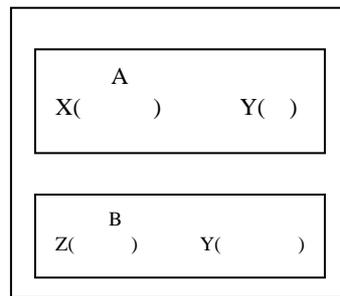
- kernel VMS
- I/O
- executive
- supervisor
- user



7-16



7-17a



7-17b

1

" " ACL Access Control List

X 2 Y 1 3 A 2 B 3 C 1
Z 7-17a

ACL

ACL

UNIX
UID

ACL

UNIX

UID GID

GID

" d"

UNIX

" -"

9

" "

" r"

" w"

" -"

" x"

Windows setUID UID ACL [Solomon 1998]

ACL

2

" " CL(Capability List)

Kerberos

" " k i " " j

7-17b

(unforgeable)

•

" "

•

•

•

•

•

```
struct capability{
    type tag
    long addr
}
```

c tag —c.tag c.addr address
—c.addr

c.tag

tag

tag

Burroghs

capability other

tag

tag

Mach Rochester Intelligent Gateway RIG

Accent

Mach

IPC

[Accetta et al. 1986]

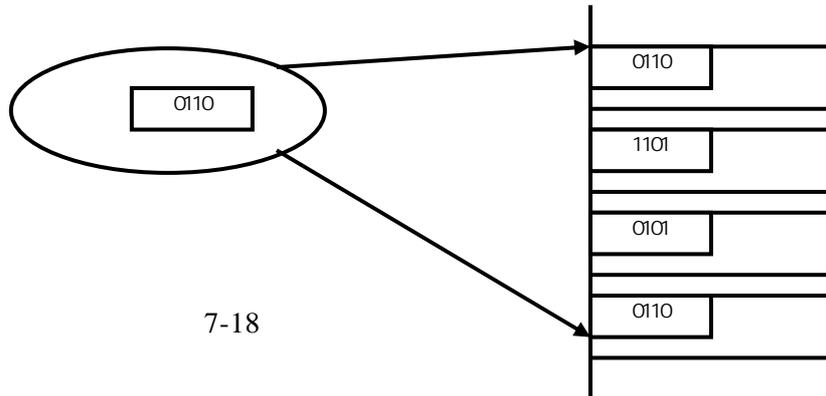
Mach IPC

(4)

ACL

20 70

ACL



20 60
k

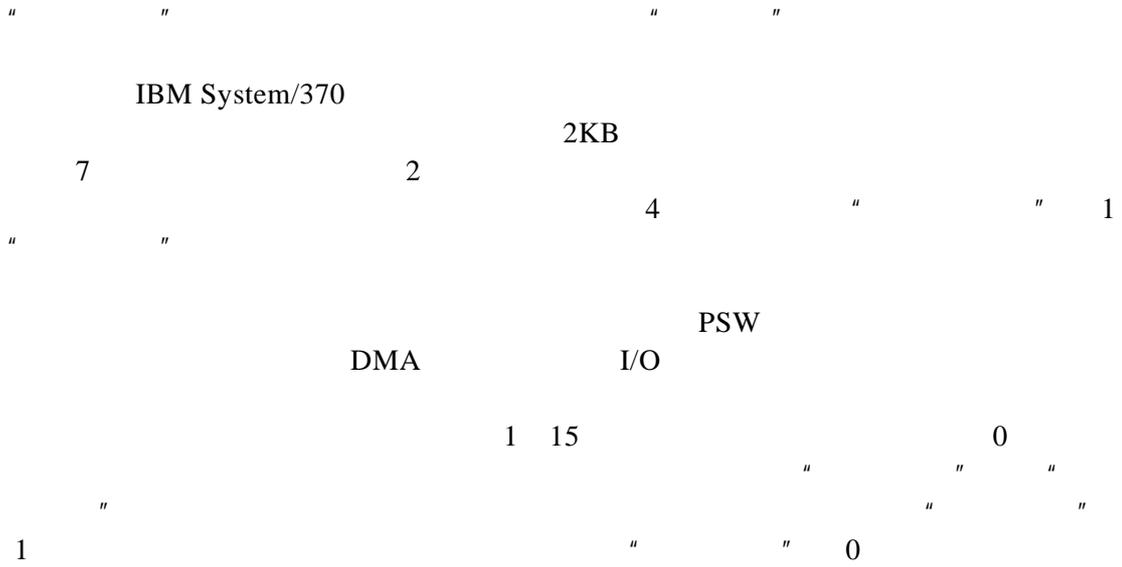
h

k

7-18

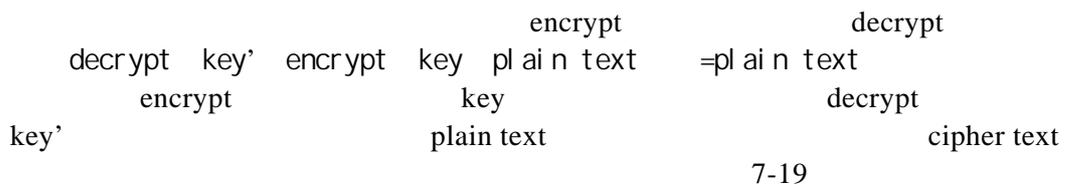
CPU

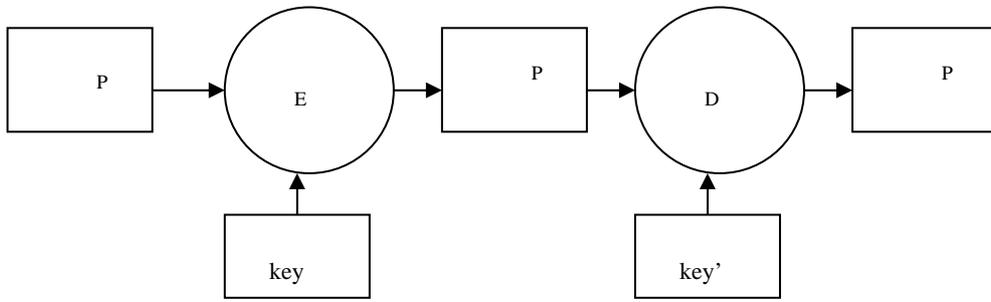
20 60 h=16 k=4
16



7.6.4

1





7-19

-
-
-
-

Kerckhoff

/

2

(1)

M	E	G	A	B	U	C	K
7	4	5	1	2	8	3	6
P	l	e	a	s	e	t	r
a	n	s	f	e	r	o	n
e	m	i	l	l	i	o	n
d	o	l	l	a	r	s	t
o	m	y	S	w	i	s	s
B	a	n	k	a	c	c	o
u	n	t	s	i	x	t	w
o	t	w	o	a	b	c	d

Please transfer one million dollars to my Swiss
Bank account six two two ...

AFLLSKSOSELAWAIATOOSSCTCLNMAN
TESILYNTWRNNTSOWDPAEDOBNO ...

7-20

7-20

A

1 U

8

(2)

Julius Caesar

house

krxvh

26 k
house 26 7-21
qifsb

a	b	c	d	e	f	g	h	i	j	k	l	m	n	o	p	q	r	s	t	u	v	w	x	y	z
↓	↓	↓																					↓	↓	↓
z	x	c	v	b	n	m	q	w	e	r	t	y	u	i	o	p	a	s	d	f	g	h	j	k	l

7-21 26

3

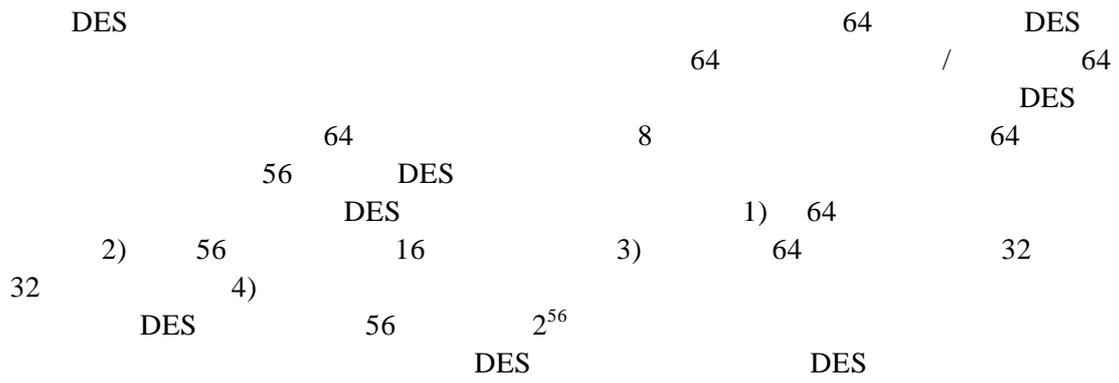
(1)

DES(Data Encryption Standard)
IBM
DES
DES

70

20

DES

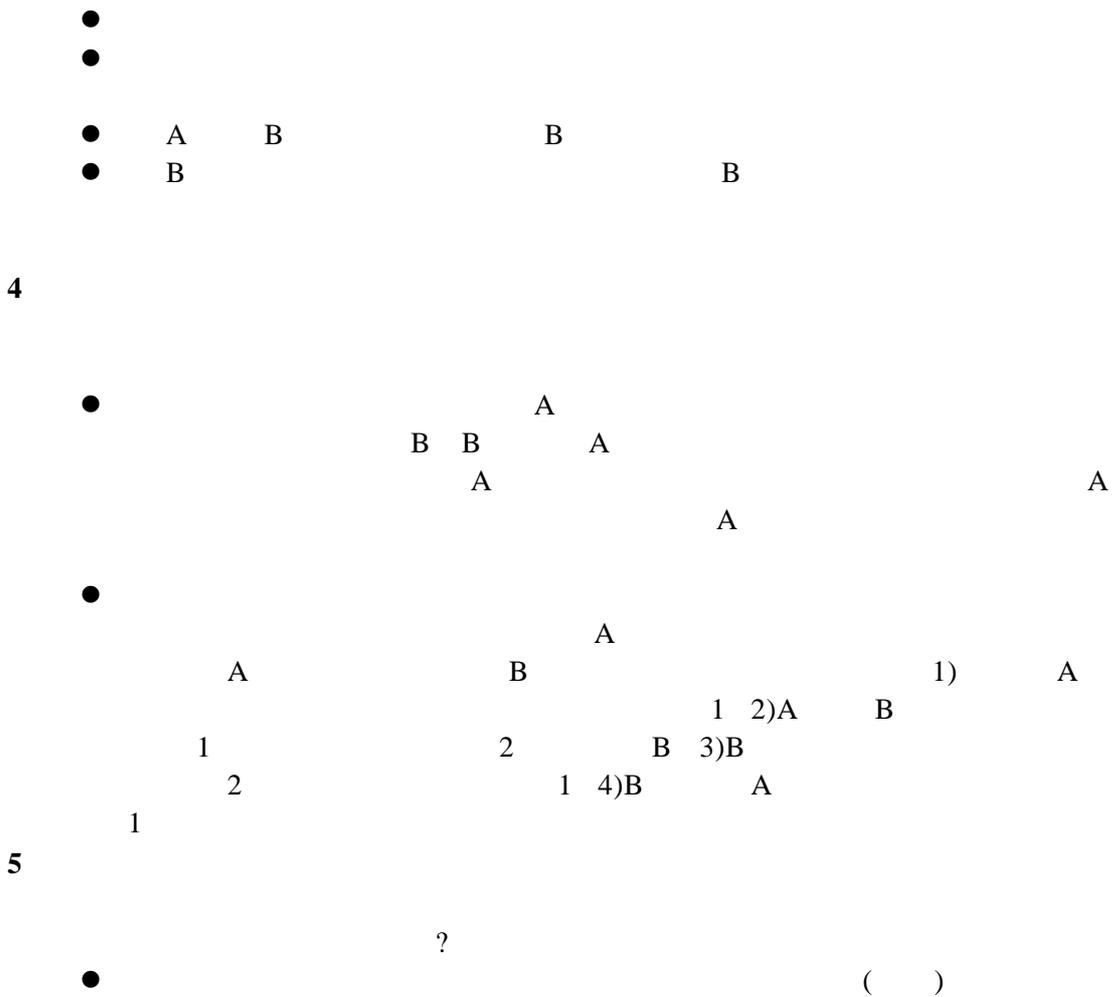


(2)

DES

1976 Diffie

Hallman



/ ?

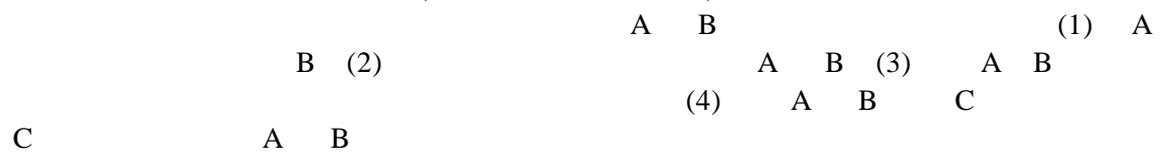
/



?

KDC(Key Distribution Center)

ACC(Access Control Center)



7.6.5

(auditing)

) () () () ()

7.7 Windows 2000/XP

7.7.1. Windows 2000/XP

Windows2000/XP

	C2	1995			Windows NT
Server	Windows NT Workstation 3.5				NCSC United
	States National Computer Security Center	C2			
	http://www.radium.ncsc.mil	1996	Windows NT Server	Windows NT Workstation	
	3.51			ITSEC	UK Information
	Technology Security Evaluation and Certification			F-C2/E3	
	C2		http://www.itsec.gov.uk		Windows
NT 4.0	Windows2000		NCSC	ITSEC	

-
-
-
-

Windows2000/XP

Windows2000/XP

ID

Windows NT4

Windows2000/XP

-
- Kerberos 5
Internet
- Secure Sockets Layer 3.0
- CryptoAPI 2.0

7.7.2. Windows2000/XP

Windows2000/XP

- SRM NTOSKRNL.EXE
- LSA LSASS EXE
- " " LSA

HKEY-LOCAL-MACHINE\security

- SAM
- LSASS
- SAM
- HKEY-LOCAL-MACHINE\SAM
 - MSV1_0 DLL
 - Windows LSASS DLL
 - SAM
- WINLOGON.EXE
- LSA
- SERVICES.EXE
- LSASS

7.7.3. Windows2000/XP

Windows2000/XP

Windows2000/XP

Windows2000/XP

7.7.4.

Windows 2000/XP

Windows 2000

/

ID SID

-
-

SID

Windows 2000/XP
Windows 2000/XP

Windows 2000/XP

RPC

/

7.7.5.

7.22 a

● ID

● SID

● SID SID

●

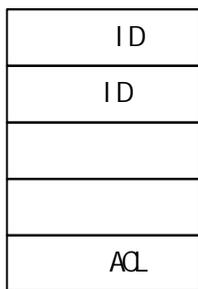
●

● ACL

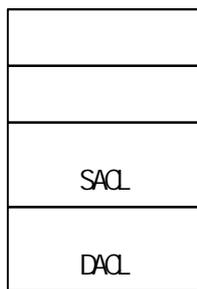
ID
SID

SID

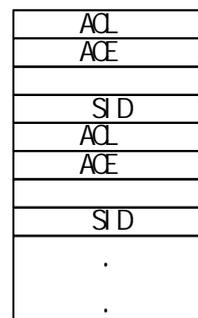
ACL



(a)



(b)



(c)

7.22 WINDOWS 2000/XP

7.7.6.

7.22 b

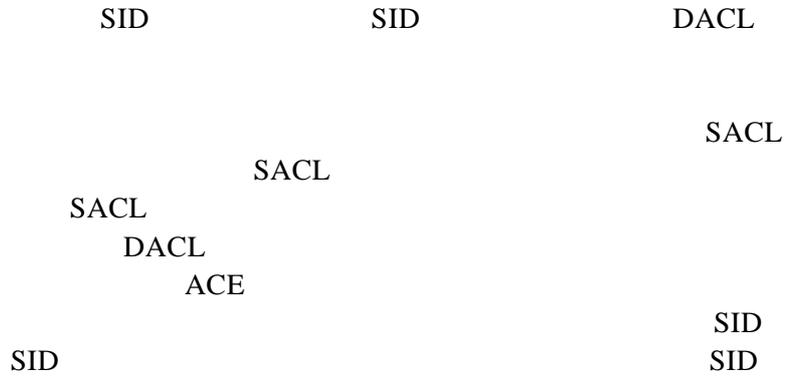
●

SACL

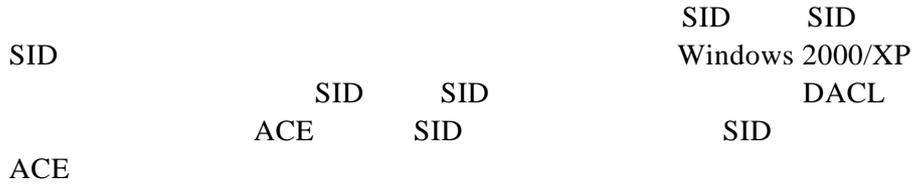
DAACL

RPC

-
-
-

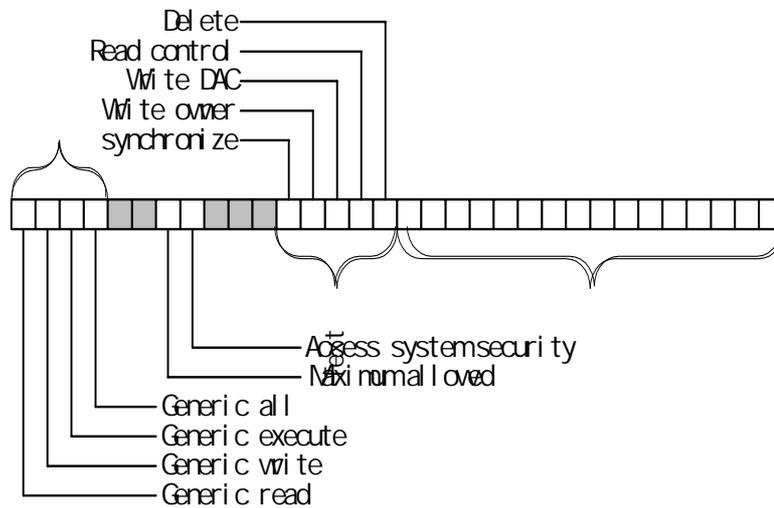


Windows 2000/XP



7.23

0 File_Read_Data 16 0 Event_Query_Status



16

7.23

5

- Synchronize
- Write_Owner

DAC

- Write_DAC DACL
- Read_Control DACL
- Delete

4

ACE ACE ACE ACE

- Generic_all
- Generic_execute
- Generic_write
- Generic_read

Generic_Read

Read_control Synchronize
 File_Read_Data File_Read_Attribute File_Read_EA ACE
 SID SID Generic_Read SID

5

Access_System_Security
 SID ACE SID

Maximum_Allowed
 DACL Windows 2000/XP Windows 2000/XP DACL SID
 ACE ACE
 DACL DACL

Maximum_Allowed

SID

Windows 2000/XP
 Windows 2000/XP

/

/

7.8

I/O

DES

ISO

(

RSA

)

1

2

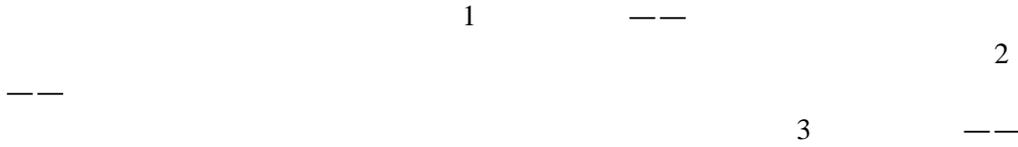
1				UNIX/Linux	rwX		
		?(1)Rick	Jennifer			(2)Helen	Anna
		(3)Cathy				UNIX/Linux	
2							
3	26		4				
						?	
4		5000				4990	
	(1)		(2)				
5		A1	A2	A3	B1	B2	
			S				
6							
						WJ	

CH8

8.1

8.1.1

1



2

-
-
-

é z

5
4 5 1 2 3
6
1
2
3

8.1.2

1

2

3
●
●
●
4
●
●
●
5
●
●

6



7



8



packet
X.25



ATM

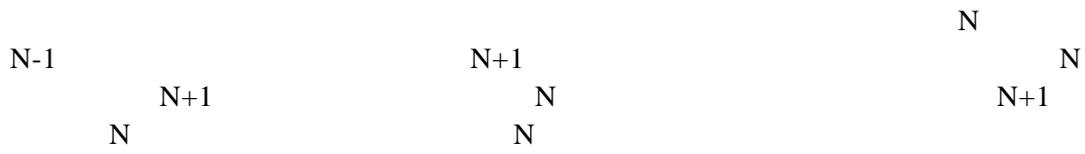
9.

10

/ Kb/s " " / Mb/s " " / b/s " "

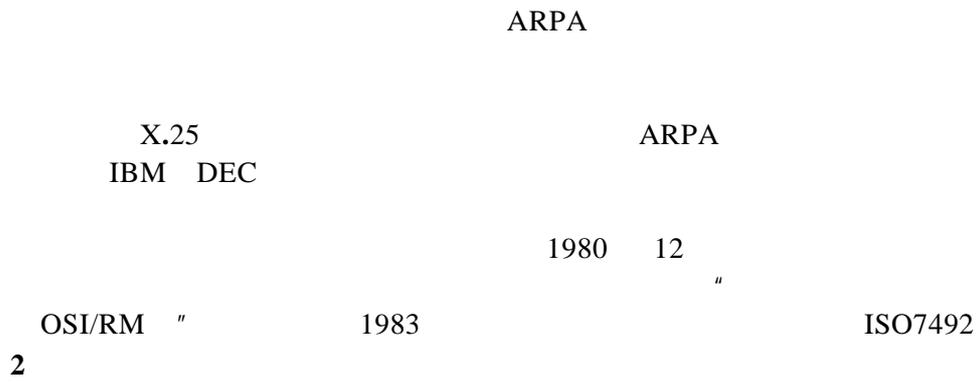
8.1.3

1

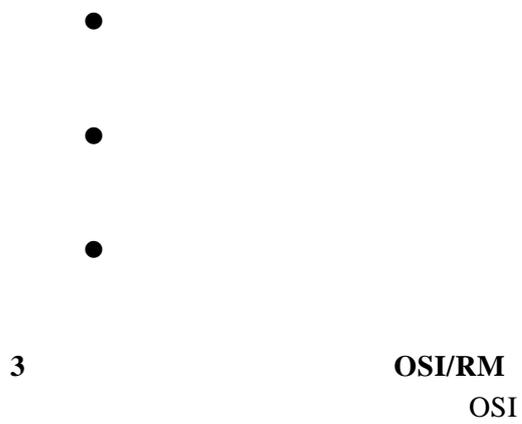


Architecture

Network



Network Protocol



3

7

8-1 OSI/RM

1

bit

2

RS-232



X. 25

IP

IP

IP

IP

X. 25

4



X. 25 IP

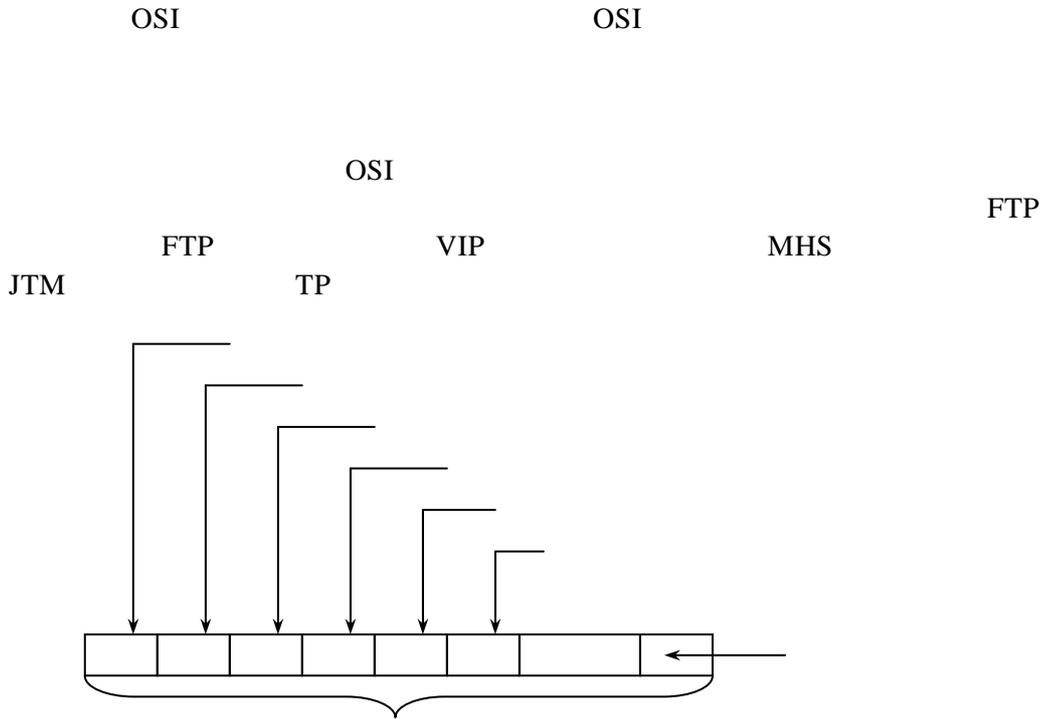
5

6

OSI



•
•
7

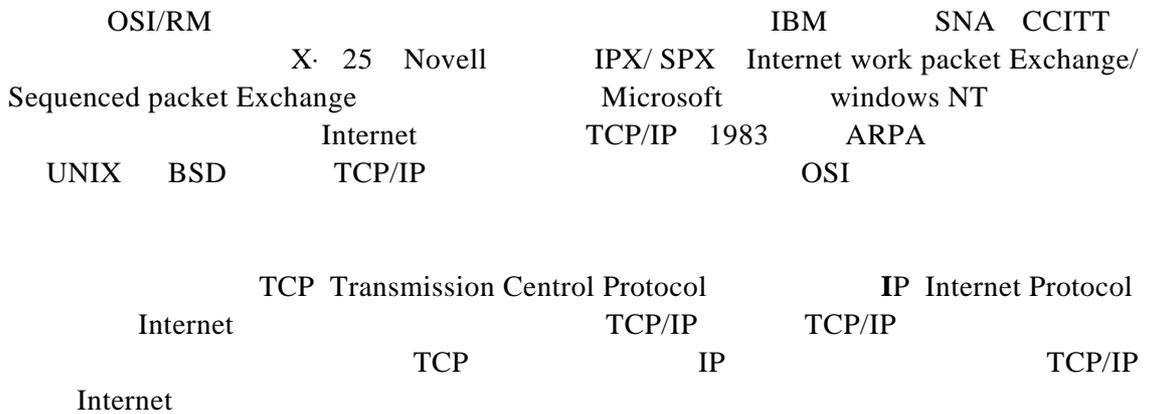


8-2

OSI/RM

8-2

4 TCP/IP



8.2

8.2.1

windows NT
for wrokgroup UNIX

MS-DOS OS/2 Windows 98 Windows

WWW

1

2 /

UNIX

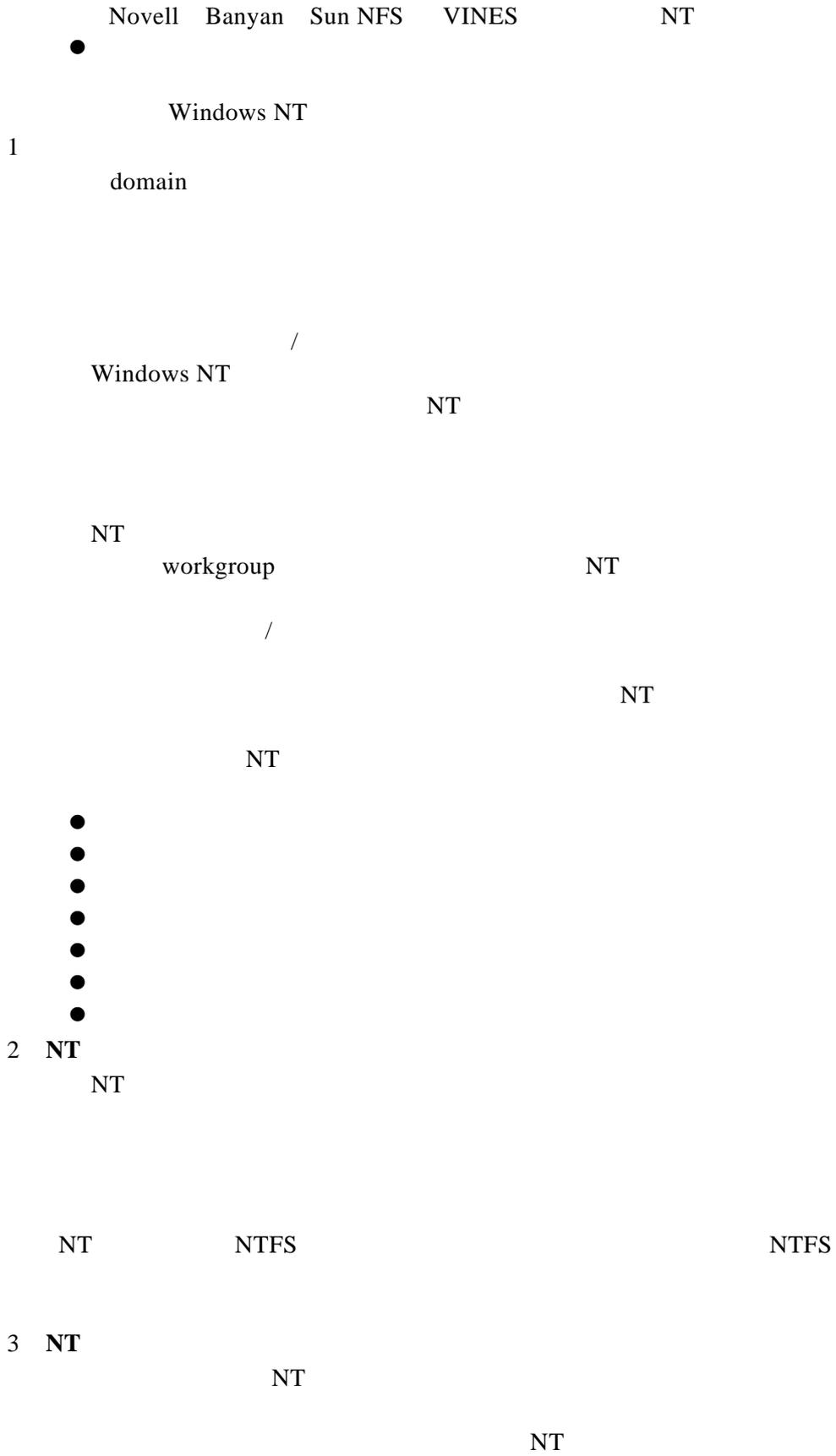
Netware Windows NT

C/S

C/S

3

Netware Lite Windows for workgroup



4 NT

NT

Internet/Intranet

● DNS

IP

DNS

Windows Internet

● Internet

Internet

IIS
Internet

NT

●

IIS
NT

NT

IP

IPX

LAN

WAN

●

IP

TCP/IP

● NT

NT

●

(PPTP)

NT

Internet

Windows NT RAS

RAS

WAN

Windows NT

●

DCOM NT
Internet

COM

5 NT

8.3

8.3.1

" Single Computer System Image "

●

●

●

●

"

"

"

"

-
-
-
-
-

Distributed Operating System

-
-
-
-

8.3.2

Procedure
1

socket

RPC Remote

/

send

receive

1

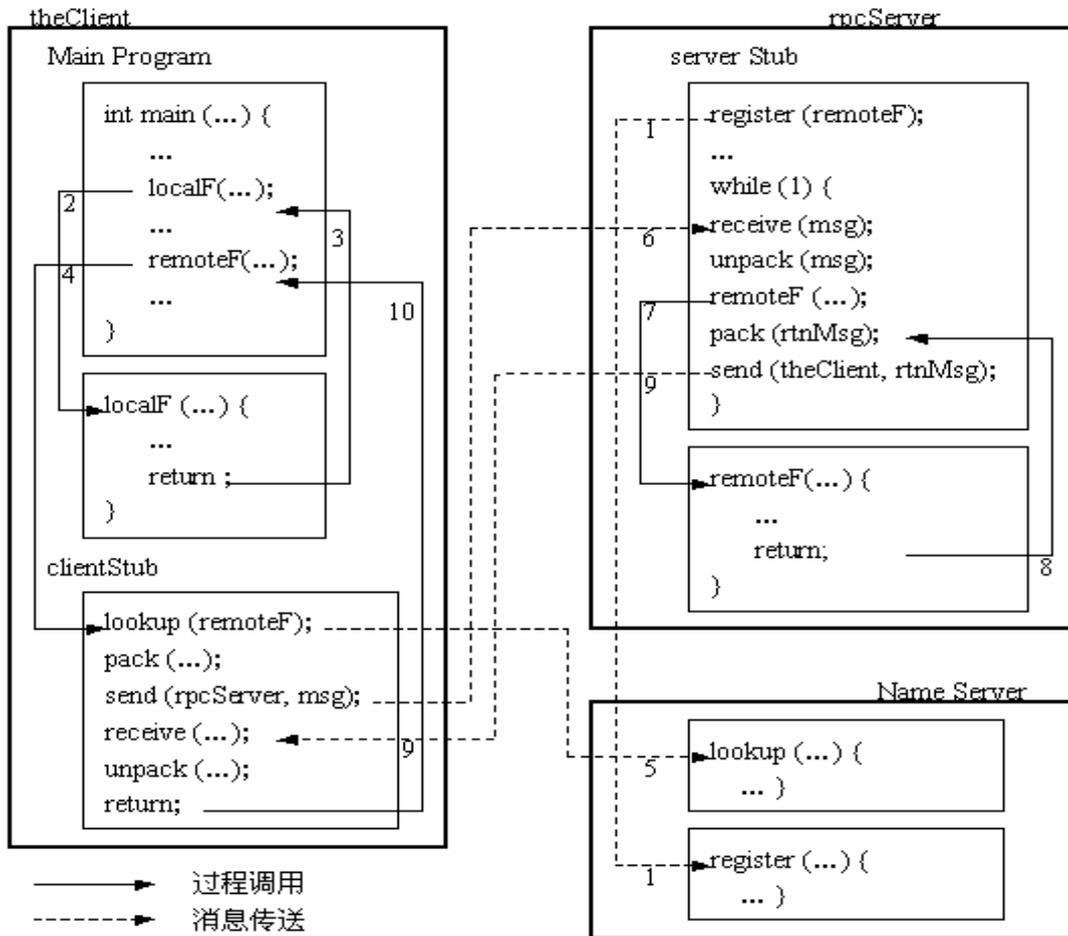
stub

/

RPC

Q4AAÄDWYÄqCQ4]QäbÄ
RPC

server stub



8-5

RPC

RPC

ASCII

0 1

3

UNIX BSD

socket

socket

socket

socket

socket

socket

socket

Socket

socket

C/S

socket

/

socket

socket

socket

C/S

socket

socket

socket

Socket

bind() send() sendmsg() recvmsg()

socket

socket

connect() write() writev()

listen()

socket()

accept() read() readv() recvfrom()

socket TCP/IP

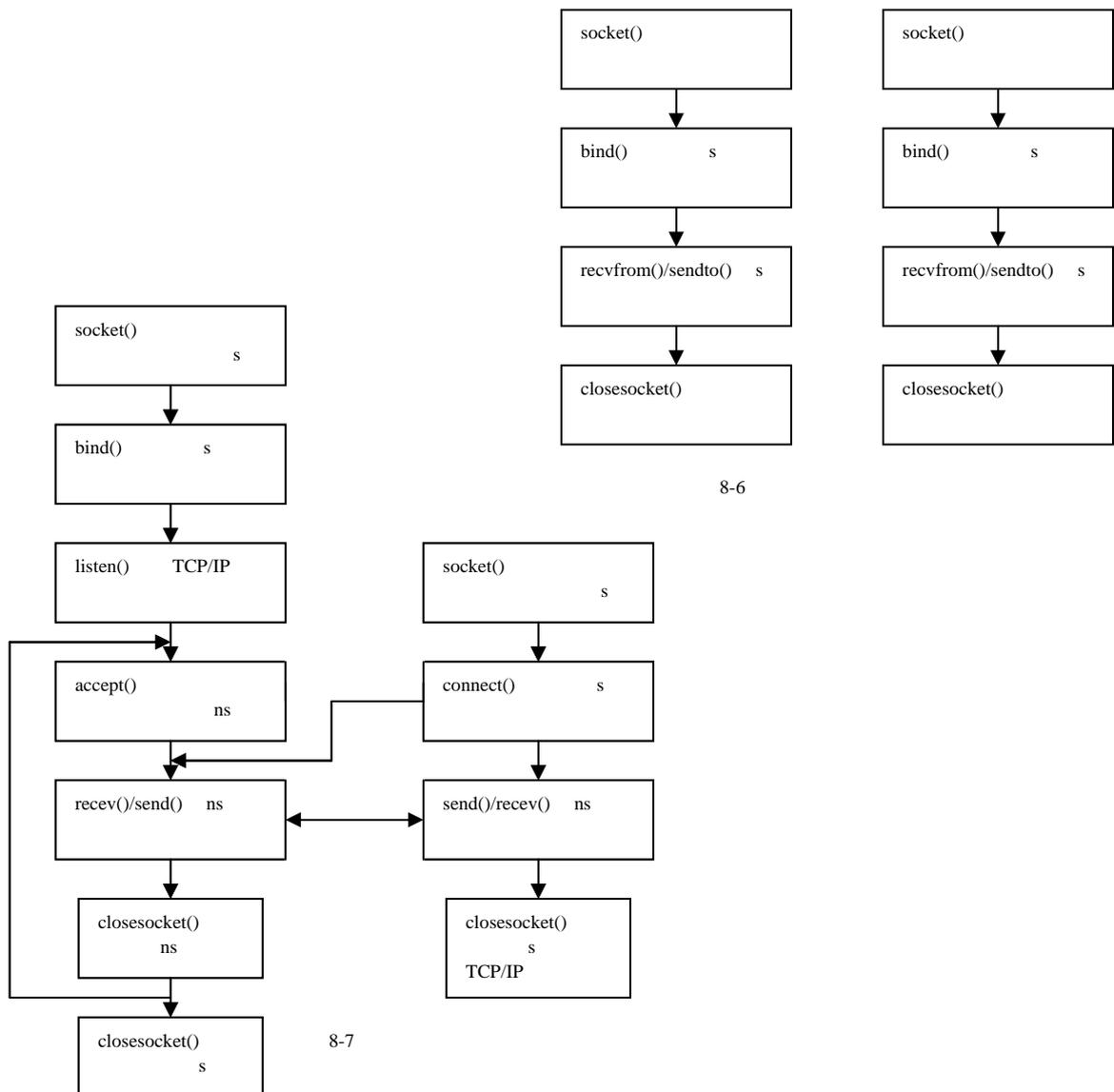
TCP/IP

C/S

8-6 8-7

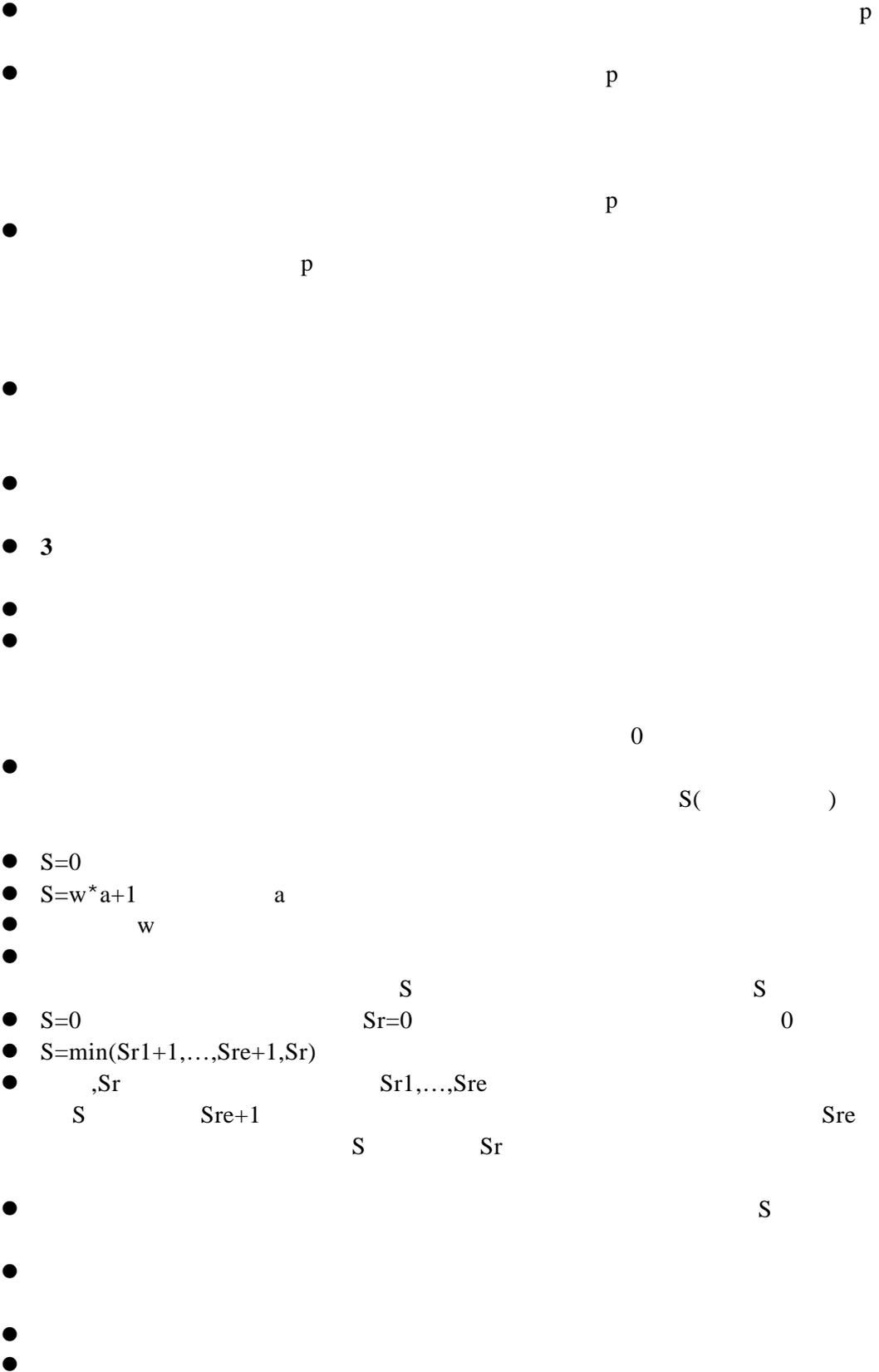
fork

Windows2000/XP



8.3.3

é



8.3.4

a f x y y x x y

timestamping

b C(a) C(b) C(a) C(b) C a b a
 a b a
 b

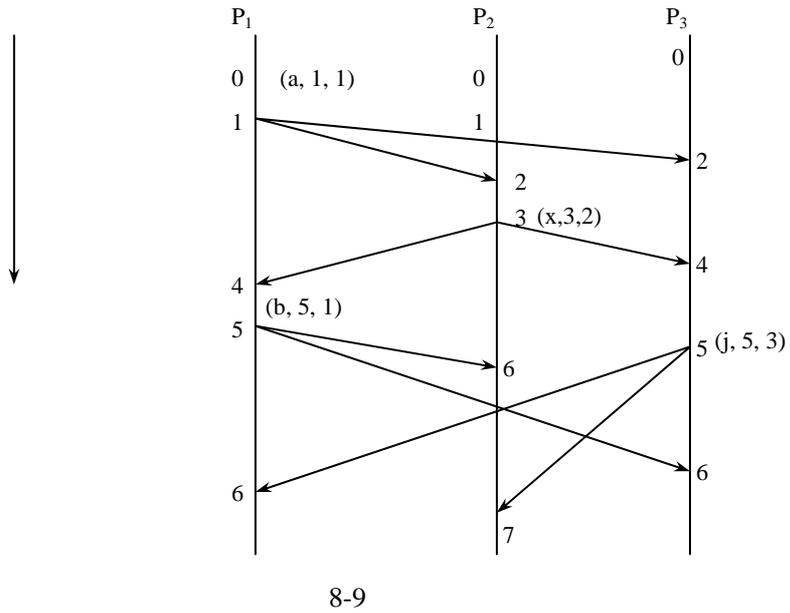
C
 (1) P e_j e_j P 1
 C(e_j)=1
 e_j P j j-1 e_{j-1}
 C(e_j)=C(e_{j-1})+1
 (2) P e_r e_r P 1
 C(e_r)=1+C(e_{s'})
 ● e_{s'} P' e_r P r r-1
 ● C(e_r)=1+max[c(e_{r-1}) C(e_{s'})]

n i C_i (m, T_i, i) m
 T_i i
 j (2)
 1
 i x j
 y
 (1) T_i < T_j
 (2) T_i = T_j i < j

Lamport

8-8 P₀
 6 P₁ 8 P₂
 10
 6 P₀ A P₁

P_1 b P_3 j a x b j
 8-10 P_1 P_4 3 P_1 P_4 2
 P_4 a q
 =1 3 a q ? P_1 P_4 $(T_1$
 $T_4=1)$ $i < j$ a q P_4



2

(1) Lamport

Lamport
FCFS

N

Lamport
Lamport

1~N

applicationstack = (release 0 i) i=1,...,N

- (request T_i i) P_i
- (reply T_j j) P_j
- (release T_k k) P_k

Lamport

(1) P_i
(request, T_i, i) applicationstack[i]

(2) P_j (request, T_i, i)
applicationstack[i] P_j (request, T_i, i)
 P_j

(3) P_i
1) P_i applicationstack request

(T_i,i)

```
(4)          Pi          applicationstack [i]          (request,Ti,i)
              (release, Ti,i)
(5)          Pj          Pi          release          applicationstack[i]
              Pi          (request,Ti,i)
              3(N-1)          (N-1) request          (N-1) reply
(N-1) release

/*          */
type message=record
  class: (application,reply,release); /*          */
  source: 1..n; /*          */
  timestamp:integer; /*          */
  clock:integer; /*          */
  valid:Boolean; /*          */
end;
/*          */
var
  T:integer; /*          1*/
  applicationstack:array [1..n] of message; /*          */
  replycount:0..n-1; /*          */
procedure Apply; /*          */
var M:message;
    i:integer;
begin with M do
  begin class:=application; /*          */
        source:=me; /*me          */
        timestamp:=T /*          */
        valid:=true;
  end;
  T:=T+1;
  applicationstack[me]:=M; /*          */
  replycount:=n-1; /*          n-1          */
  For i:= 1 to n do
  begin M.clock:=T;T:=T+1;
        If i = n then send(M,i) /*          i          */
  end;
  waitFor(Replycount=0); /*          n-1          */
  for i:= 1 to n do /*          */
  waitFor (not Applicationstack[i] Applicationstack[me]);
/*          "          "          */
```

```

    get resource;          /*          */
    end;
procedure Receive(var M:message);
/*          waitfor          */
var R:message;
begin T:=1+max(M.clock,T);
    with M do
        case class of
application:
            begin applicationstack[source]:=M;
                with R do
                    begin
                        class:=reply;
                        source:=me;
                        clock:=T;
                    end;
                    T:=T+1;
                    send(R,source);          /*          */
                end;
            reply:
                replycount:=replycount-1;
            release:
                applicationstack[source].valid:=false;
        end;
procedure Release;          /*          */
var R:message;
    i:integer;
begin with R do
    begin class:=release;
        source:=me;
        clock:=T;
    end;
    T:=T+1;
    applicationstack[me].valid:=false;
    for i:= 1 to n do
        if i = n then send(Mi);          /*          */
    end;
end;

```

(2) G.Ricart

Ricart		release
2(N-1)		
(1)	P_i	(request, T_i, i)

```

(2)          Pj          (request,Ti,i)
      1)      Pj
              (reply,Tj,j)    Pi
      2)      Pj
              Ti < Tj      Ti = Tj    i < j
                              (reply,Tj,j)    Pi          reply
(3)          Pi          reply
(4)          Pi
              reply

```

```

type message=record
    class: (application,reply );          /*          */
    source: 1..n;                          /*          */
    timestamp:integer;                     /*          */
end;
var
    T:integer;                              /*          1*/
    Applicationtime:integer;               /*          */
    Replydeferred:array [1..n] of Boolean;
/*    i                                     Replydeferred[i]          */
    Replycount:0..n-1;                    /*          */
    Requesting:Boolean;                   /*          */

procedure Apply;                          /*          */
Var M:message;
begin with M do
    begin class:=application;
        Timestep:=T;
        Source:=me;
    end;
    applicationtime:=T;
    T:=T+1;
    requesting:=true;
    replycount:=n-1;
    broadcast(M);                          /*          */
    waitfor(replycount=0);                 /*          replycount=0*/
end;
procedure Receive(var M:message);
/*          waitfor          */
begin with M do
    case class of

```

application:

```

replydeferred[source]:=
requesting and ((timestap>applicationtime) or
(timestap=applicationtime and source>me));
/* source */
T:=1+max(T,timestap);
If not replydeferred [source] then
sendreply(source); /* source*/

```

reply:

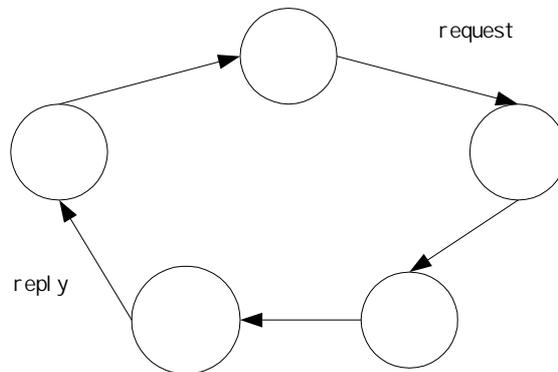
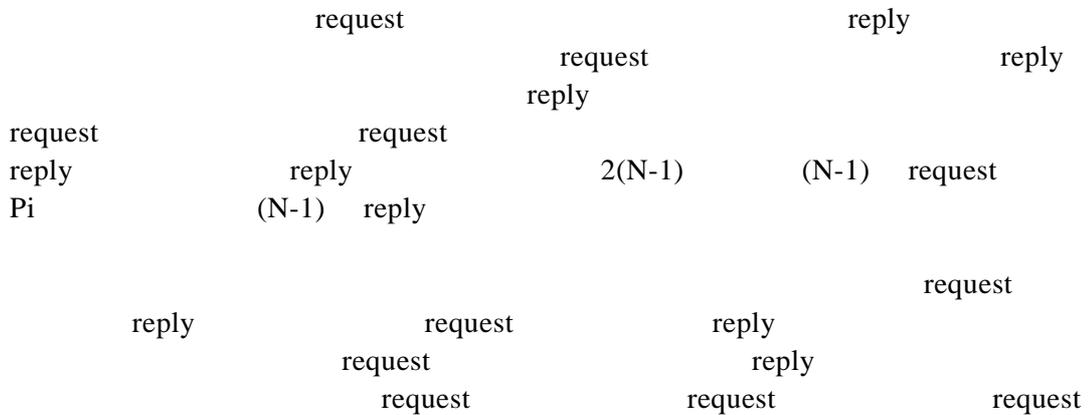
```

replycount:=replycount-1;
end;
procedure Release; /* */
var i:integer;
begin for i:= 1 to n do
if replydeferred[i] then
begin replydeferred[i]:=false;
sendreply(i); /* */
end;
requesting:=false;
end;

```

8-11

Ricart



(3)

Suzuki 1982

Logical Ring

Token

K

K+1

8.3.5

1

A

B

B

C

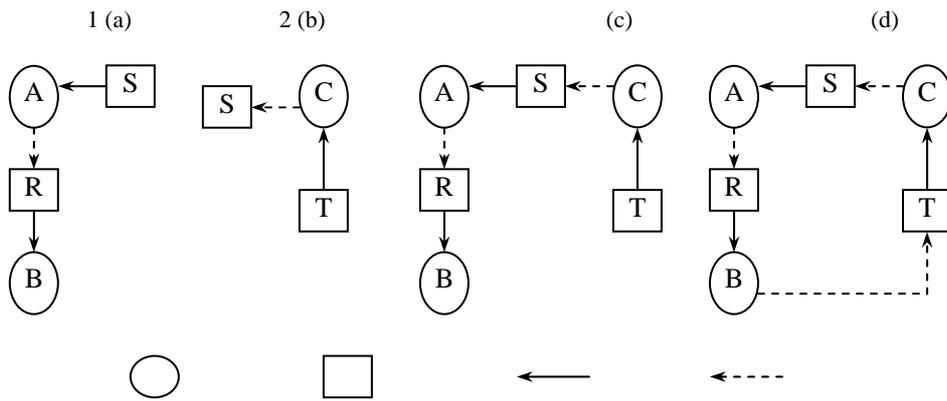
C

A

2

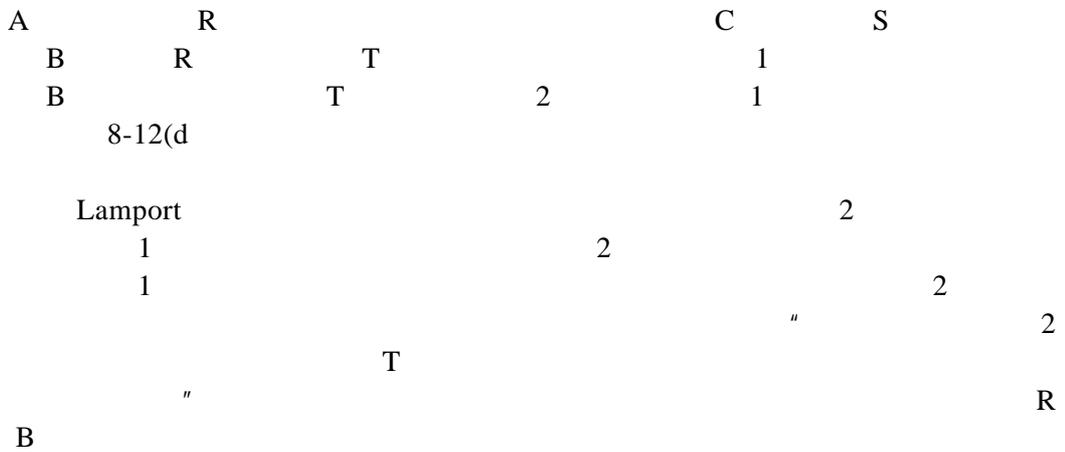
1

A B 1 C 2 R S T
 8-12(a)(b)
 ● A S R R B
 ● B R
 ● C T



8-12

8-12 c



2

- 1
- 2
- 3
- 4

" "

8.3.6

1

-
-

File Service

File server

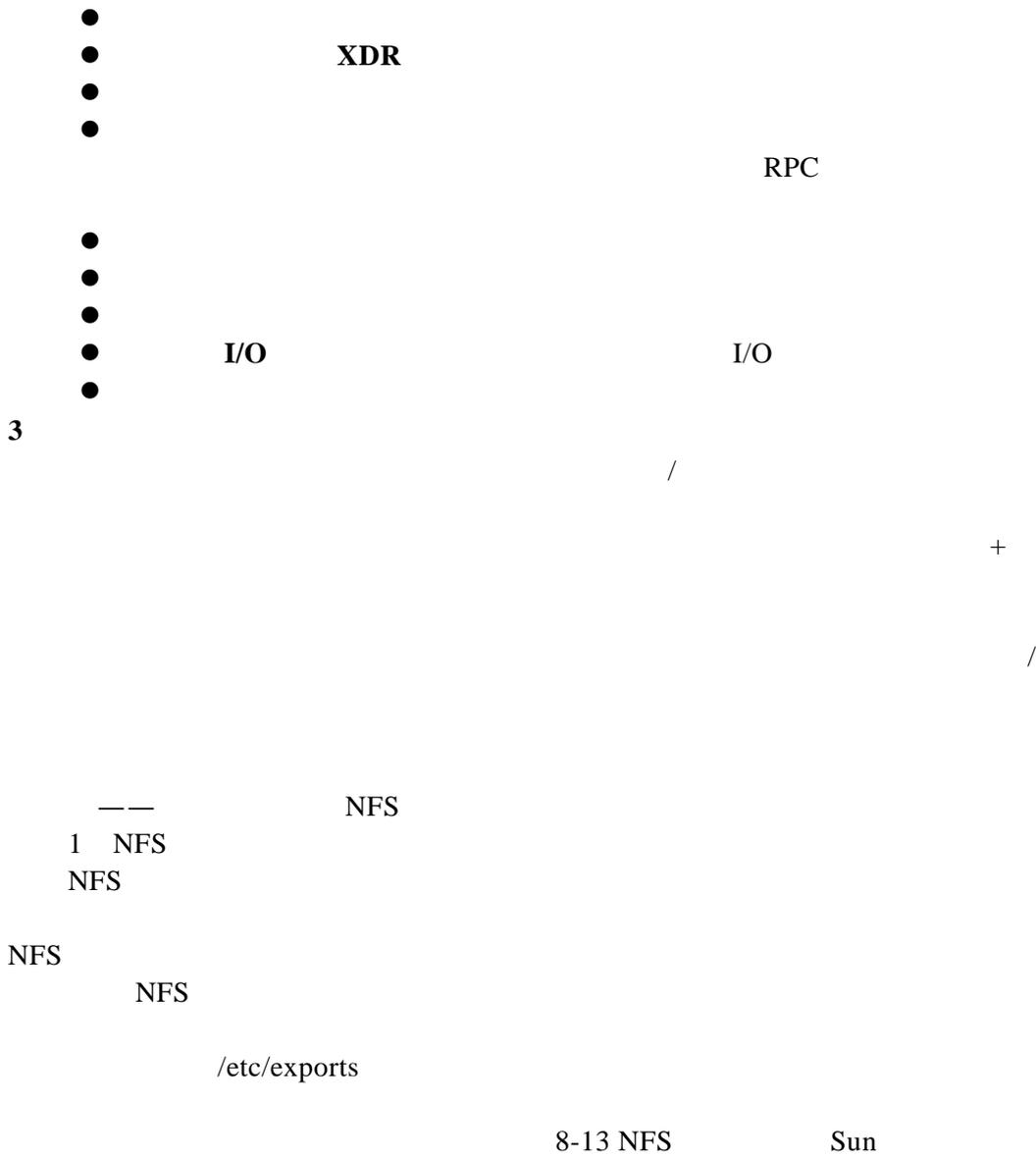
2

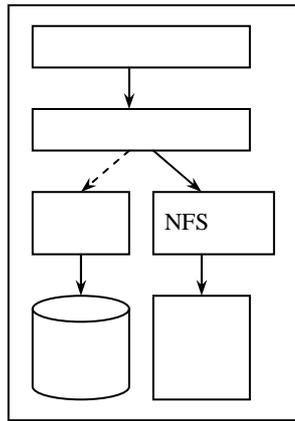
-

Sun

vfs
NFS

vnode





NFS

A

B

NFS

2 NFS
NFS

NFS

/

NFS

i

/etc/rc Shell
Shell

Sun UNIX

/etc/rc
NFS

/etc/rc

NFS

NFS

NFS
open close

UNIX

open close

open
read

lookup

NFS

UNIX

UNIX

NFS

NFS

NFS

UNIX

rwX

NFS

yellow page

NIS

NIS

3 NFS
NFS

SUN NFS
open read close
VFS

VFS
UNIX

i
VFS

UNIX

i

v
i

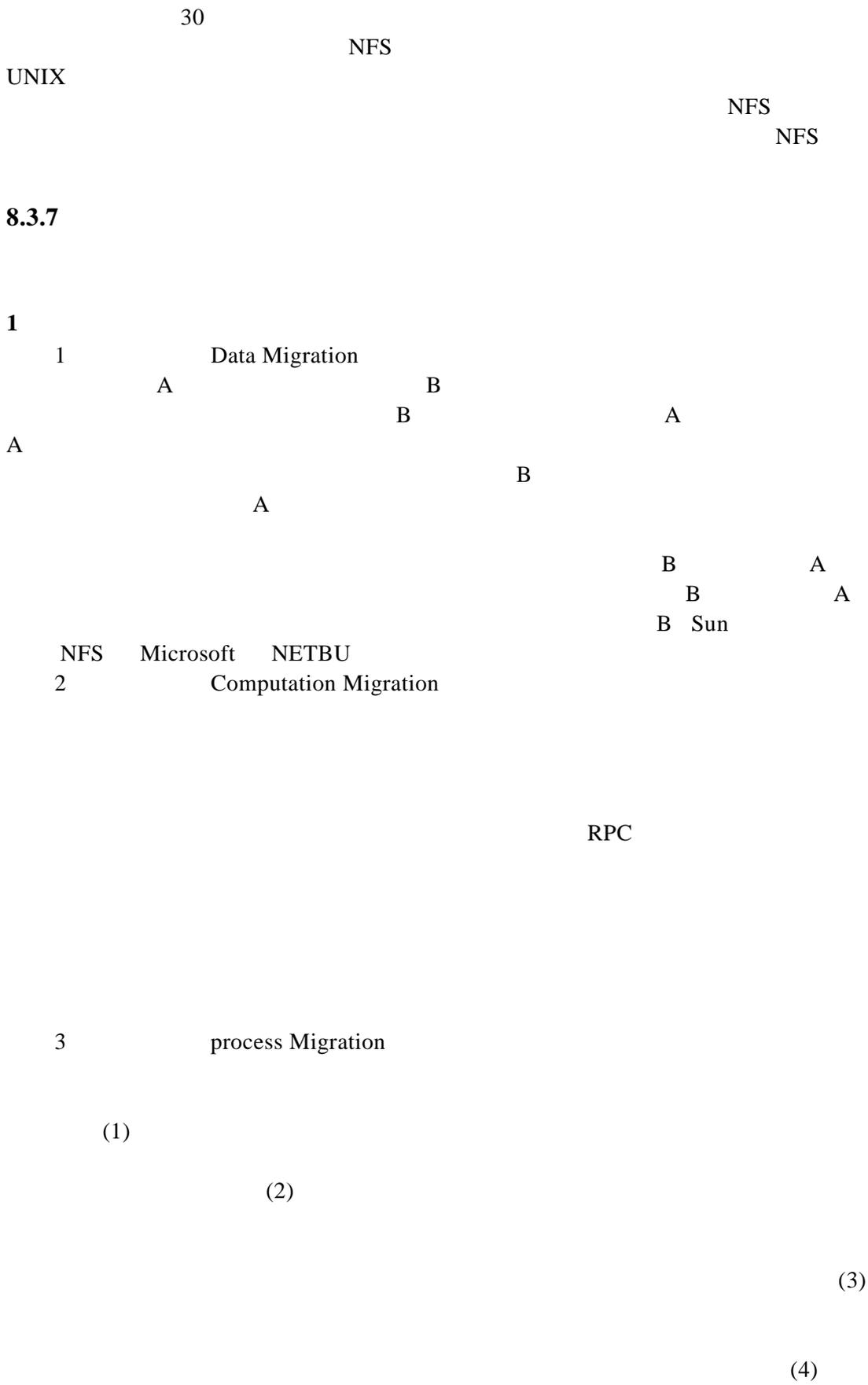
mount open read

v

mount

(1) Mount

			mount	mount	
			v	NFS	
r	i			v	r
v			NFS		r
i	v				
					VFS
(2)	Open				
			v		r
	NFS		NFS		
	r	VFS			v
	r				
Open				VFS	v
(3) /	Read/Write			read	
VFS	v			i	r
					8192
				8K	



(1)

(2)

(3)

IBM

AIX

UNIX

●

●

●

●

" "

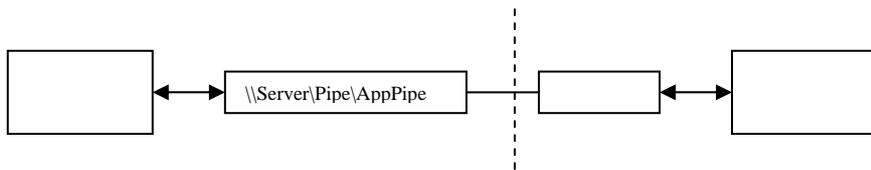
" "

UNC Windows 2000 UNC
 1)

Server \\ServerXPipe\PipeName
 DNS
 mspress microsoft com NetBIOS mspress IP 255.0.0.0
 Pipe "Pipe" PipeName
 \\MyComputer\Pipe\MyServerApp\ConnectionPipe
 CreateNamedPipeWin32
 \\.\Pipe\PipeName "\\." Win32

API

CreateNamedPipe CreateNamedPipe
 ConnectNamedPipe Win32
 CreateNamedPipe
 Win32CreateFile CallNamedPipe
 ConnectNamedPipe
 ConnectNamedPipe
 ReadFile WriteFile Win32
 8-14



8-14

API
 ImpersonateNamedPipeClient
 2)

Win32 API

CreateMailslot

CreateMailslot " \\.\Mailslot\MailslotName"

CreateMailslot

CreateMailslot

CreateMailslot CreateNamedPipe ReadFile

CreateFile

" Server Mailslot MailslotName "

" * Mailslot MailslotName"

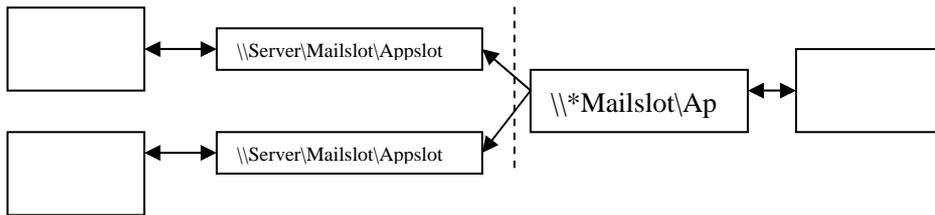
" Domain \MailslotLMailslotName"

WriteFile

425 426

425 426 424 Windows 2000 425

426 8-15



8- 15

3)

Win32

Kernel32.dll Win32 DLL ReadFile WriteFile

Win32 I/O

CreateFile Win32I O

\Winnt\System32\Npfs.sys

\Winnt\System32\Drivers\Msfs.sys

\Device\NamedPipe

\??\Pipe Device Mailslot

\??\Mailslot \\.Pipe\... \.Lmailslot\...

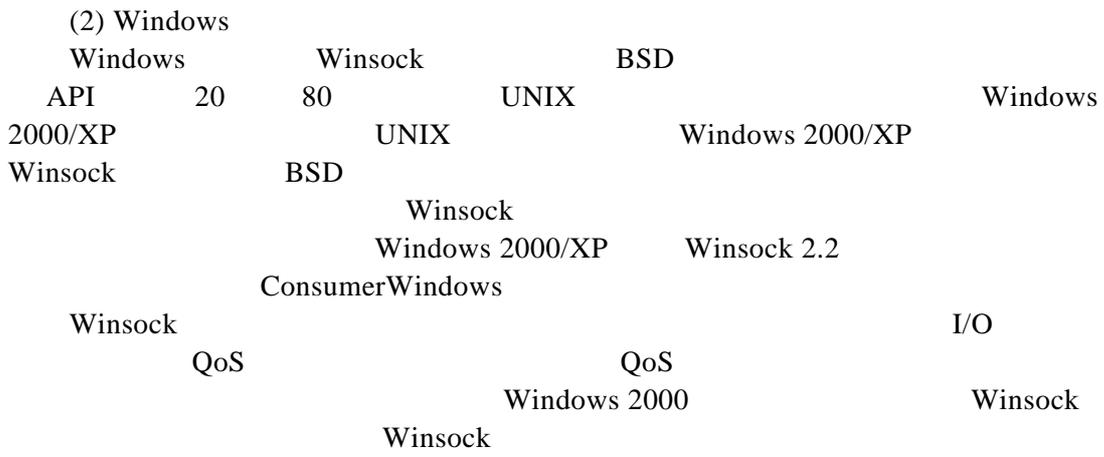
CreateFile \\. \??

CreateNamedPipe CreateMailslot

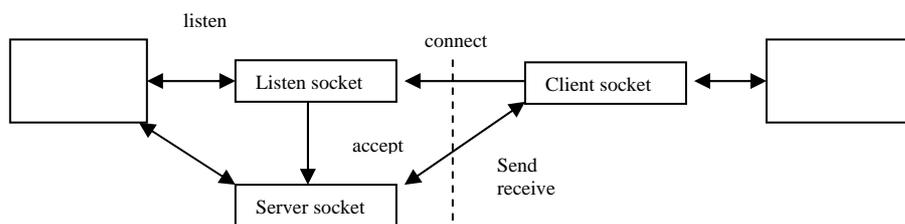
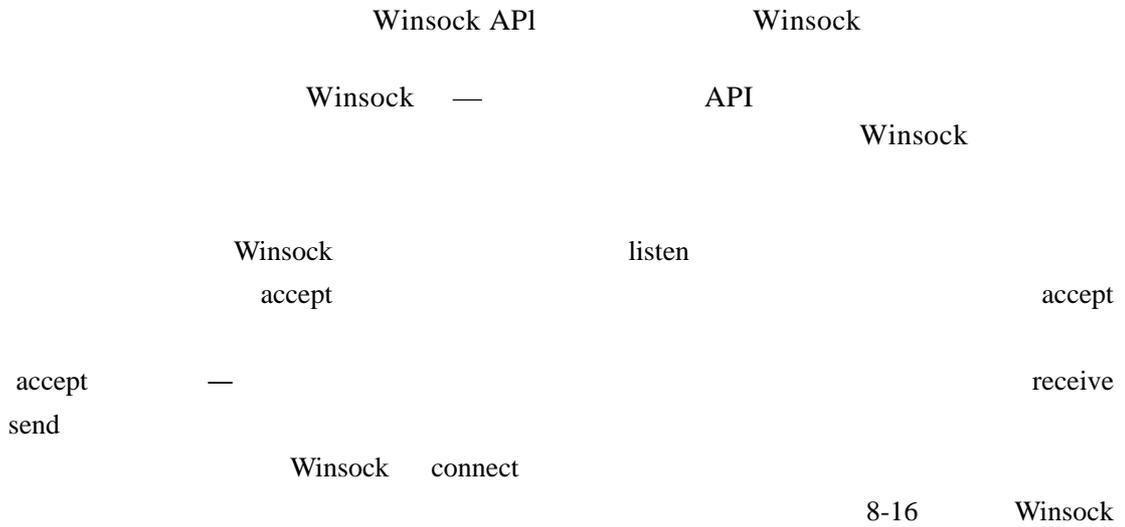
NtCreateNamedPipeFile NtCreateMailslotFile

FSD FSD

- FSD Windows 2000
 - CreateFile FSD
 - Win32 ReadFile WriteFile
 - FSD
 - FSD
- CIFS CIFS IPX TCP/IP NetBEUI FSD



1) Winsock



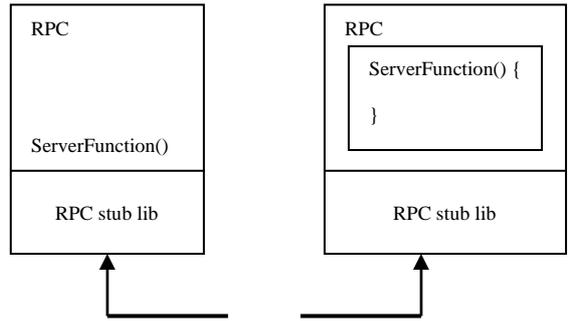
2) Winsock
 Windows Winsock API Windows
 Winsock — Windows Windows I/O
 BSD
 Winsock WSA WSAAccept
 BSD Winsock
 AcceptEx TransmitFile Windows 2000 Web
 AcceptEx accept
 Winsock accept Web
 TransmitFile Windows 2000
 TransmitFile
 Web

3) Winsock
 Windows 2000 Winsock API
 Winsock Winsock Winsock KWinsock
 Winsock SPI
 Winsock connect accept
 Winsock /
 TCP/IP
 Web IP Web 80 80 HTTP
 Web
 Winsock
 gethostbyaddr getservbyname getservbyport

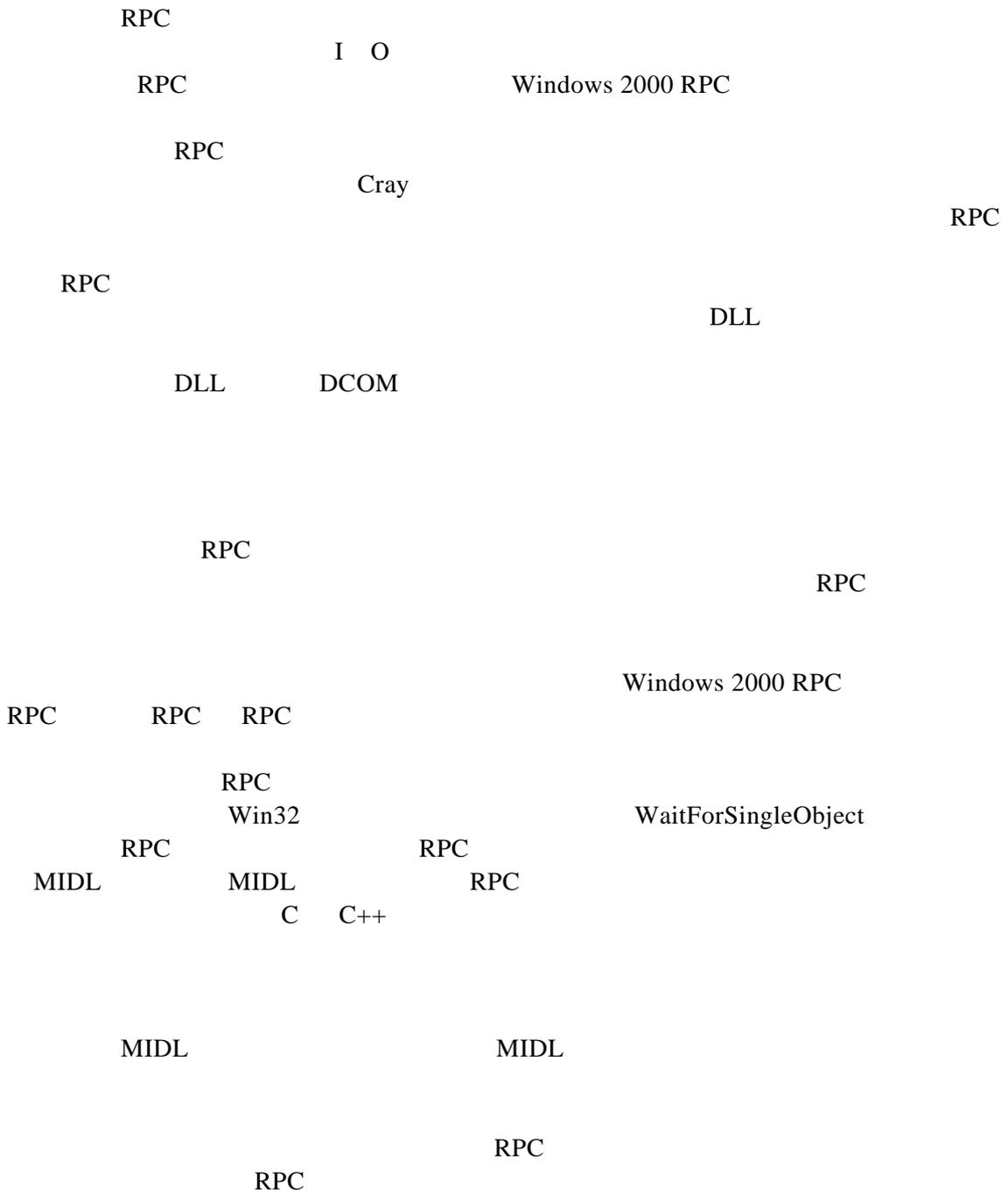
4) Winsock
 Winsock 8-17 API DLL
 (Ws2_32 dll) Ws2_32.dll
 Msafd.dll Winsock
 WinsockHelper
 Wshnetbs.d11 NetBEUI Helper Winsock
 Mswsock.dll TransmitFile AcceptEx WSARcvEx Windows
 2000 TCP IP NetBEUI AppleTalk IPX SPX ATM IrDA
 Help DLL DNS TCP IP IPX SPX

		API	Winsock	Win32 I	O		
Msafd.dll			Afd sys			AFD	TDI
	AFD	TDI IRP				Msafd.d11	AFD
			AFD				

(3)						20	80	
		RPC						
	OSF	RPC		RPC	DCE	OSF	DCE	RPC
unRPC		RPC						
		API						
	1)RPC	RPC						
			I O				Windows 2000	
	È O			O				

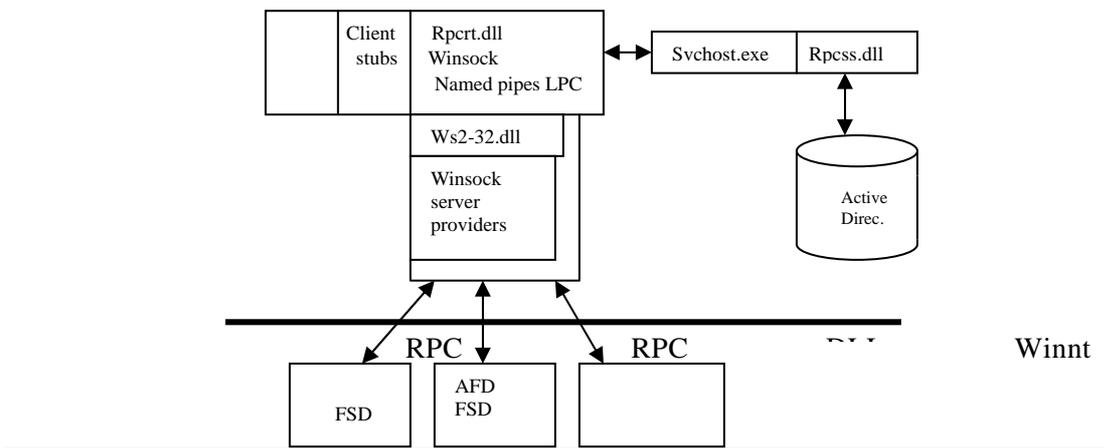


8-18 RPC



RPC RPC RPC
 RPC NetBIOS TCP IP DLL Windows 2000
 DLL Windows 2000 RPC

 RPC NetBIOS WindowsNT4 RPC RPC
 2)RPC Windows 2000 RPC SSP SSP SSP
 RPC RPC SSP RPC Winsock
 SSP Kerberos SSP Kerberos 5
 Windows 2000 SChannel SecureSocketsLayer SSL
 TransportLayerSecurity TLS PCT SSP
 RPC RpcImpersonateClient RpcRevertToSelf
 RpcRevertToSelfEx
 3)RPC 8-19 RPC RPC DLL
 \Winnt\System32\Rpcrt4.dll RPC RPC RPC
 RPC DLL RPC RPC DLL
 RPC DLL RPC RPC DLL
 RPC LPC Winsock API RPC API
 RPC DLL Winsock API RPC API



NetBIOS

NetBIOS

NetBIOS session NetBIOS

NetBIOS NetBIOS MS-DOS NetBIOS MS-DOS MS-DOS NetBIOS

NetBIOS MS-DOS MS-DOS NetBIOS

Windows 2000 Netbios

NetBIOS

3)NetBIOS API

NetBIOS \Winnt\System32\Netapi32.dll Netapi32.dll

NetBIOS NetBIOS emulation driver

Winnt\System32\Drivers\Netbios.sys Win32 DeviceControl

NetBIOS NetBIOS TDI

NetBT TCP/IP NetBIOS NetBIOS

NetBIOS \Winnt\System32\Drivers\Netbt.sys NetBT TCP/IP

NetBEUI TCP IP NetBIOS NetBIOS

NetBEUI NetBIOS NetBEUI NetBIOS

NetBT TCP/IP NetBIOS NwLinkNB

IPX/SPX NetBIOS

3.

Win32 UNC WNet API

CIFS CIFS CIFS

CIFS FSD Workstaion

Novell NetWare Windows 2000

- Multipleproviderrouter MPR DLL
- UNC Win32 WNetAPI MultipleUNCProvider MUP

Win32I O API

Windows 2000 — DNS

1)

Win32 WNet Windows

WNETAPI

provider Windows 2000

WNet

WNet DLL Workstanon

DLL

WNet
 provider interface
 \Winnt\System32\Ntlanman.dllk
 HKLM\SYSTEM\CurrentControlSet\Services\lanmanworkstation\NetworkProvJder
 ProviderPath
 WnetAddConnection API
 HKLM\SYSTEM\CurrentControlSet\Control\Network Provider\Order\ProviderOrder
 MPR

ProviderOrder

WNetAddConnection

2) UNC
 UNC MUP MPR UNC
 I/O
 MPR MUP MUP MPR
 Win32 DLL Kernel32.dlt I/O API DLL MUP
 UNC “\??\UNC” NtCreateFile
 “\??\UNC” \Device\Mup MUP

MUP
IRP

\\WIN2KSERVER\PUBLIC\insidew2k\chapl3.doc
 \\WIN2KSERVER\PUBLIC MUP

MUP

ProviderOrder MUP I/O

3) DNS IP
 DNS IP TCP IP DNS DNS
 DNS ‘ IP ’ IP
 zone DNS
 DNS Windows 2000 Windows 2000
 Windows 2000DNS Win32 Winnt System32 Dns.exe
 Windows 2000 DNS
 Windows 2000 DNS

4.

API API API
API
Windows 2000 DLC NetBEUI TCP IP NWLink
Windows2000
ServicesForMacintosh AppleTalk
● DLC IBM HP
“ ” API DLC
● IBM 1985 NetBEUI NetBEUI LAN
NetBIOS API NetBEUI
WAN NetBIOS NetBEUI
NetBIOS Extended UserInterface NetBIOS
NetBEUI NetBEUI NetBIOS
NBF Windows 2000 NetBEUI Windows
●

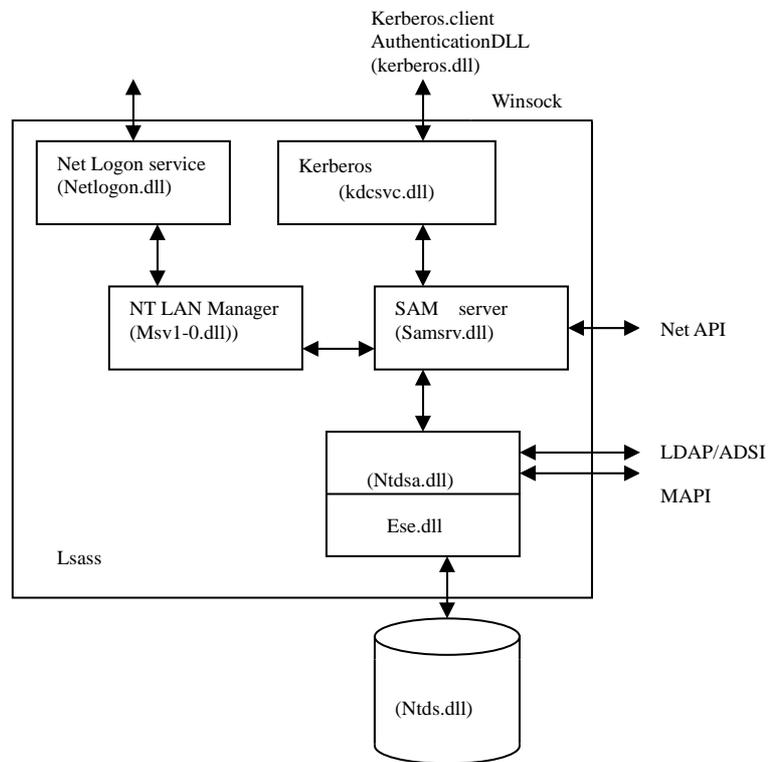
- C++ ADSI API MAPI Windows MicrosoftExchange OutlookAddressBook
- SAM API MSV1_0 \Winnt\System32\Msv1_0.dll NT LAN Kerberos \Winnt\System32\Kdcsc d11
- WindowsNT4 API NetAPI NT4 SAM WinntkNtds Ntds dit

Lsass on-disk Win32 DLL

ESE

MicrosoftExchangeServer5.5

8-20



8-20

3

Windows 2000 AdvancedServer

NDIS

32

IP

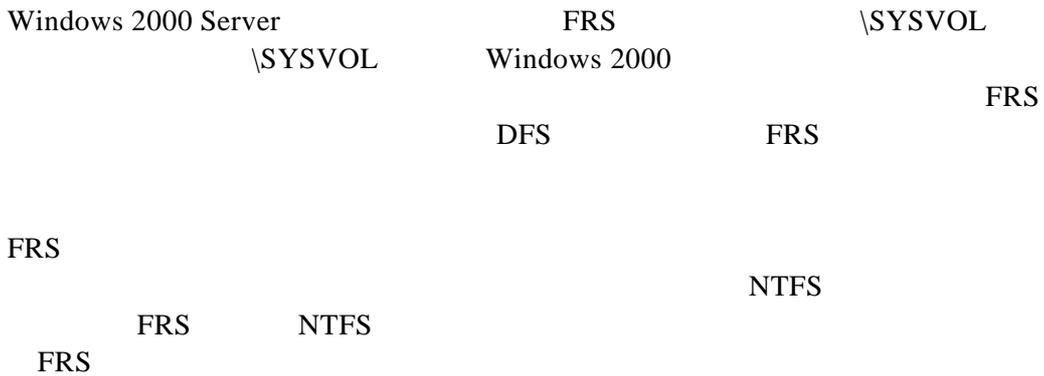
NDIS

TCP IP

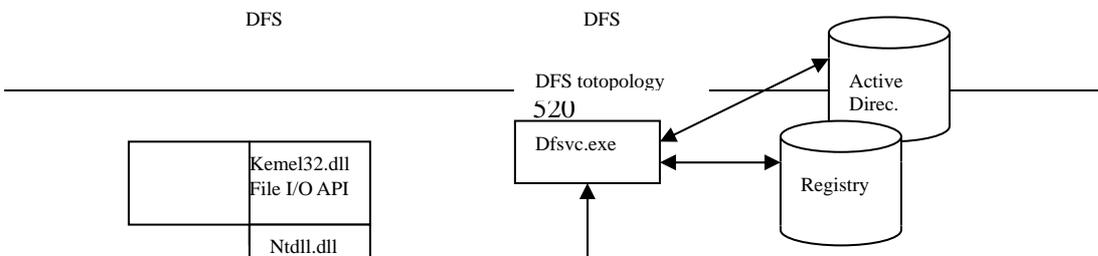
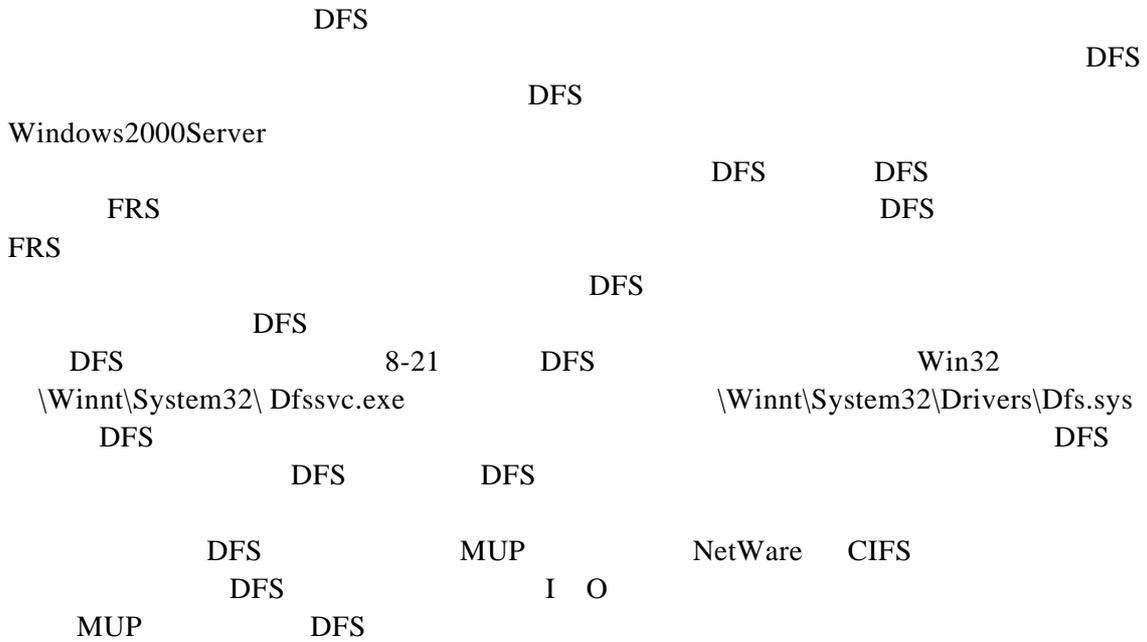
TCP IP

Windows

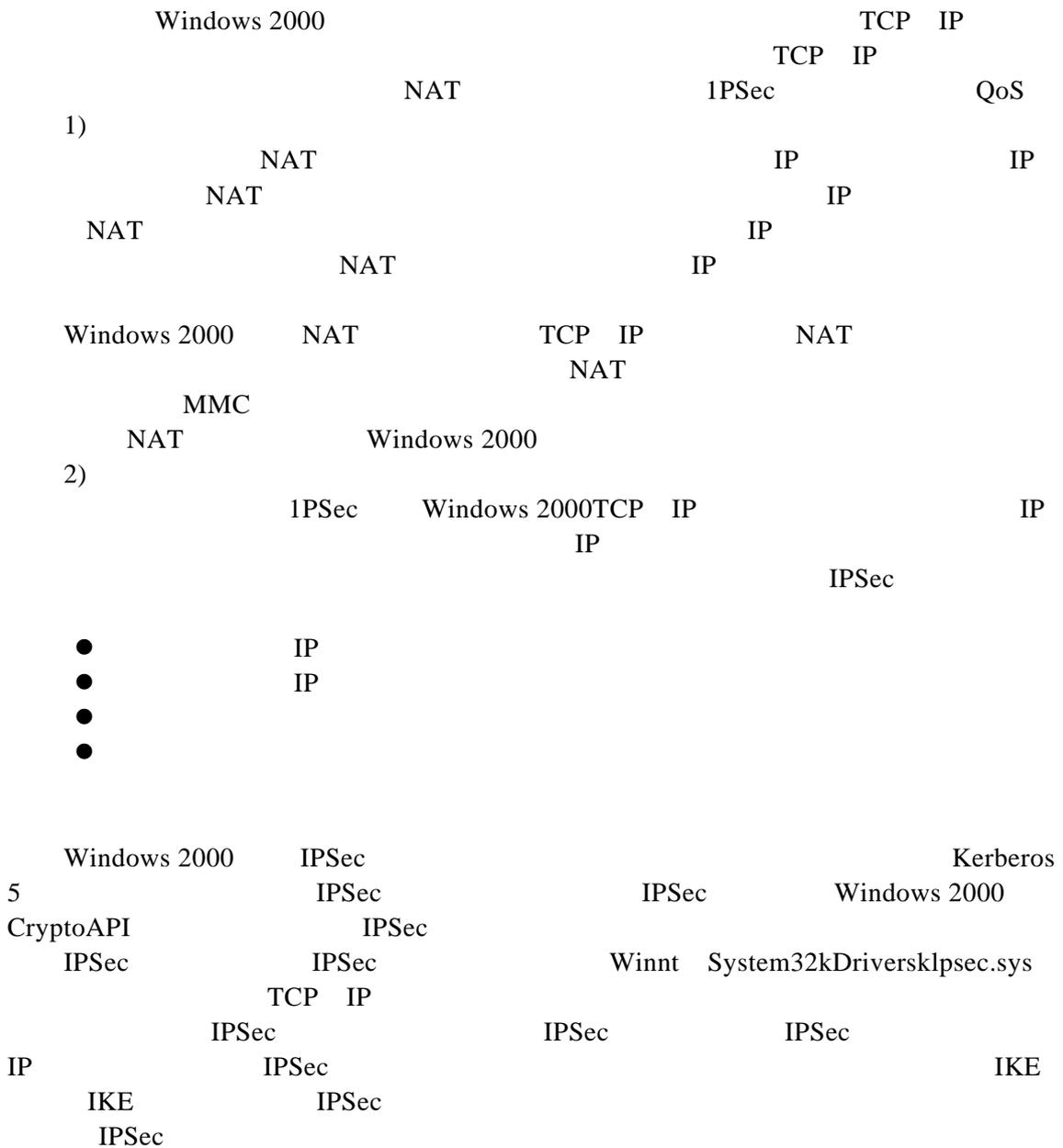
4

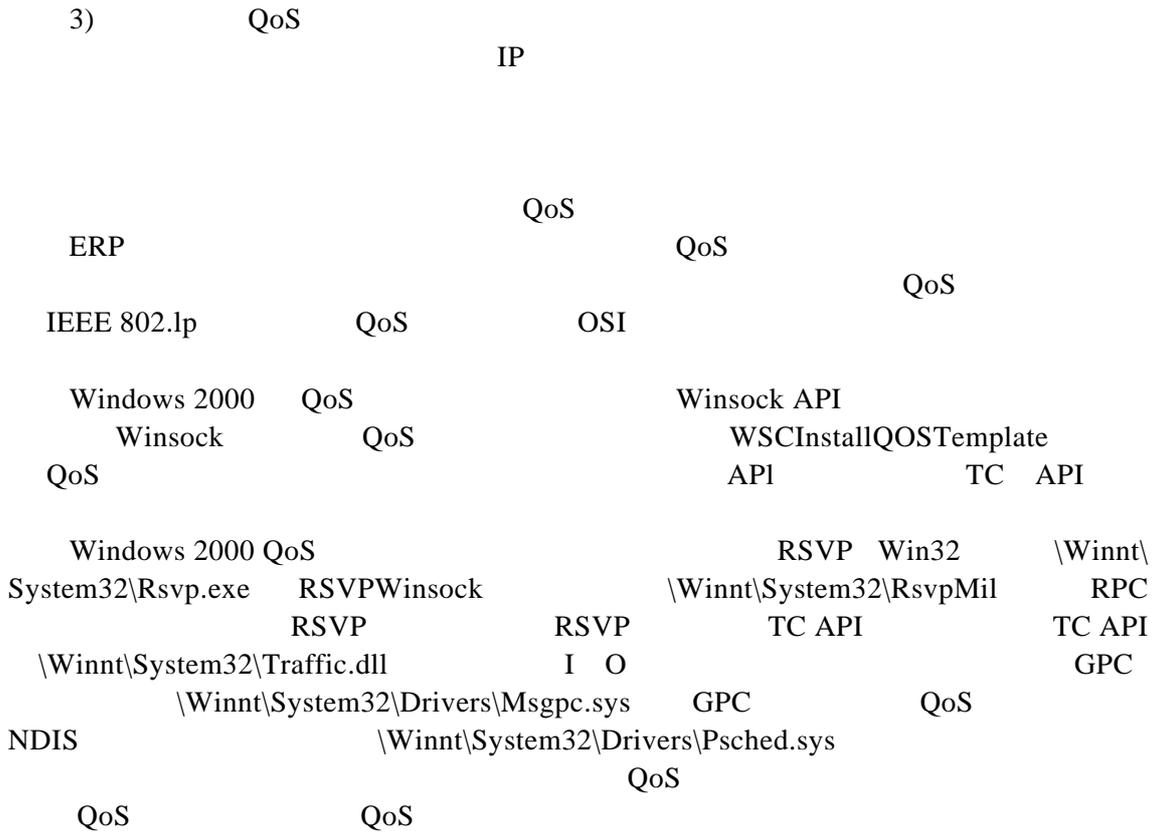


5



6 TCP IP





8.5

OSI/RM
TCP/IP

Windows NT

- 14.
- 15.
- 16. ?
- 17. " ?

4 0 Tc ()<41a913d7115e1d00349817e4/TT2 1 Tf 19 0 TD 7

3 8-8(b)

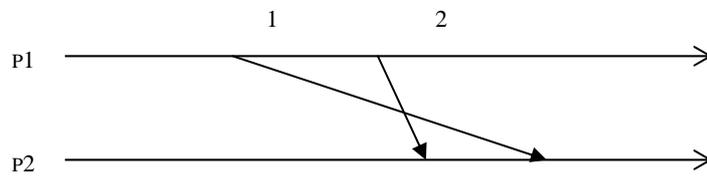
A

A

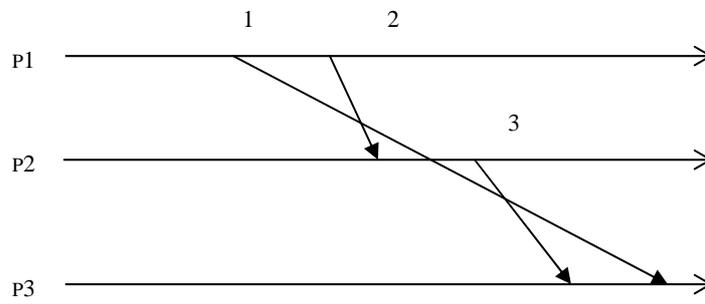
A

4

()



5



1. Operating Systems Internals and Design Principles, William Stalling (Fourth edition) Prentice-Hall International, Inc., 1998
2. Operating System Concepts, Abraham Silberschatz (Bell Labs), Peter Baer Galvin (Cor. Tech. Inc.) (Fifth edition), John Wiley & Sons, Inc., 1997
3. Distributed Systems Principles and Paradigms Andrew S. Tanenbaum, Maarten van Steen, Prentice Hall 2002
4. Applied operating System Concepts, Abraham Silberschatz, Peter Baer Galvin, John Wiley & Sons, Inc., 1998
5. Computer Network (Third Edition) Andrew S. Tanenbaum Prentice-Hall