

НРС



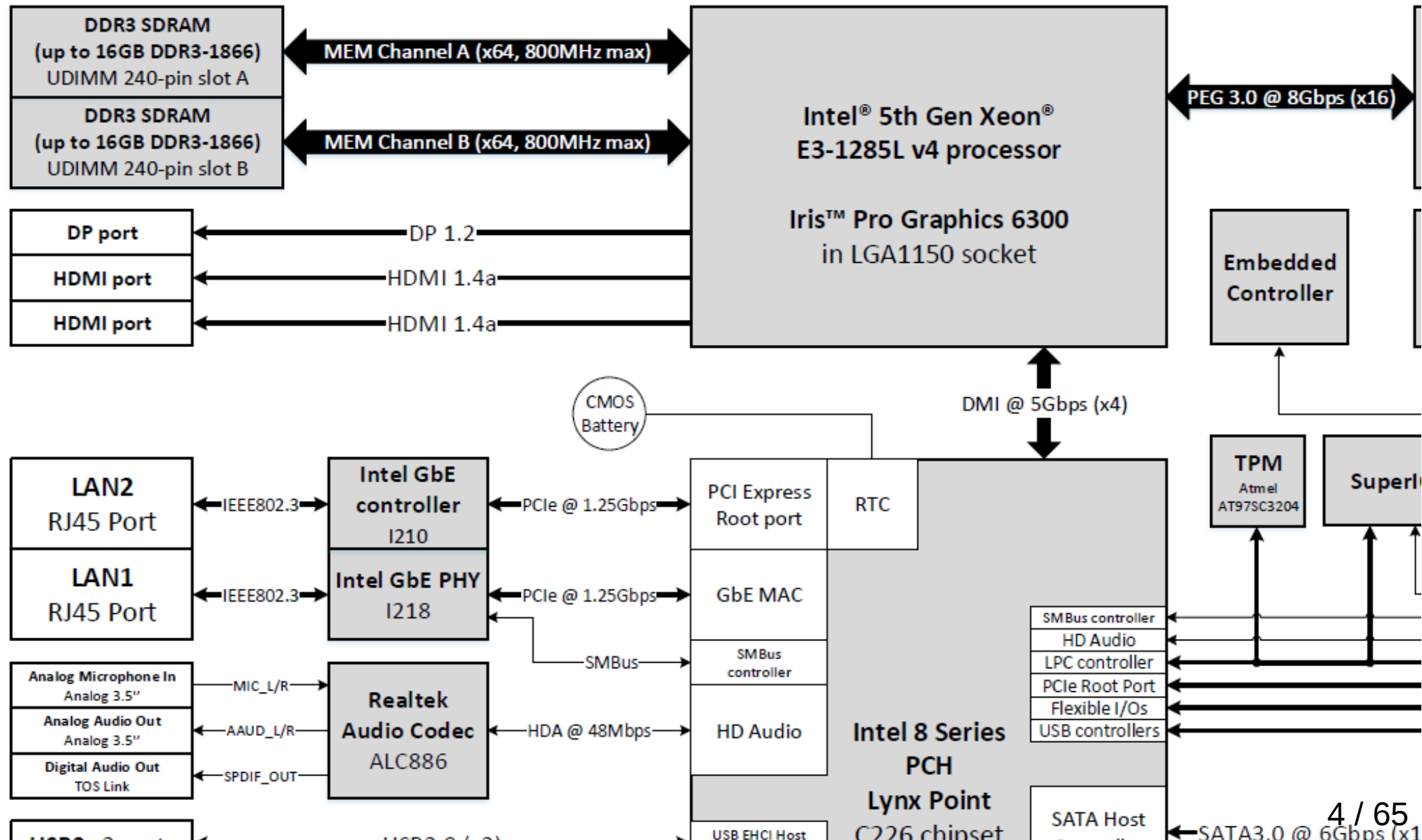
Основные характеристики вычислительного кластера

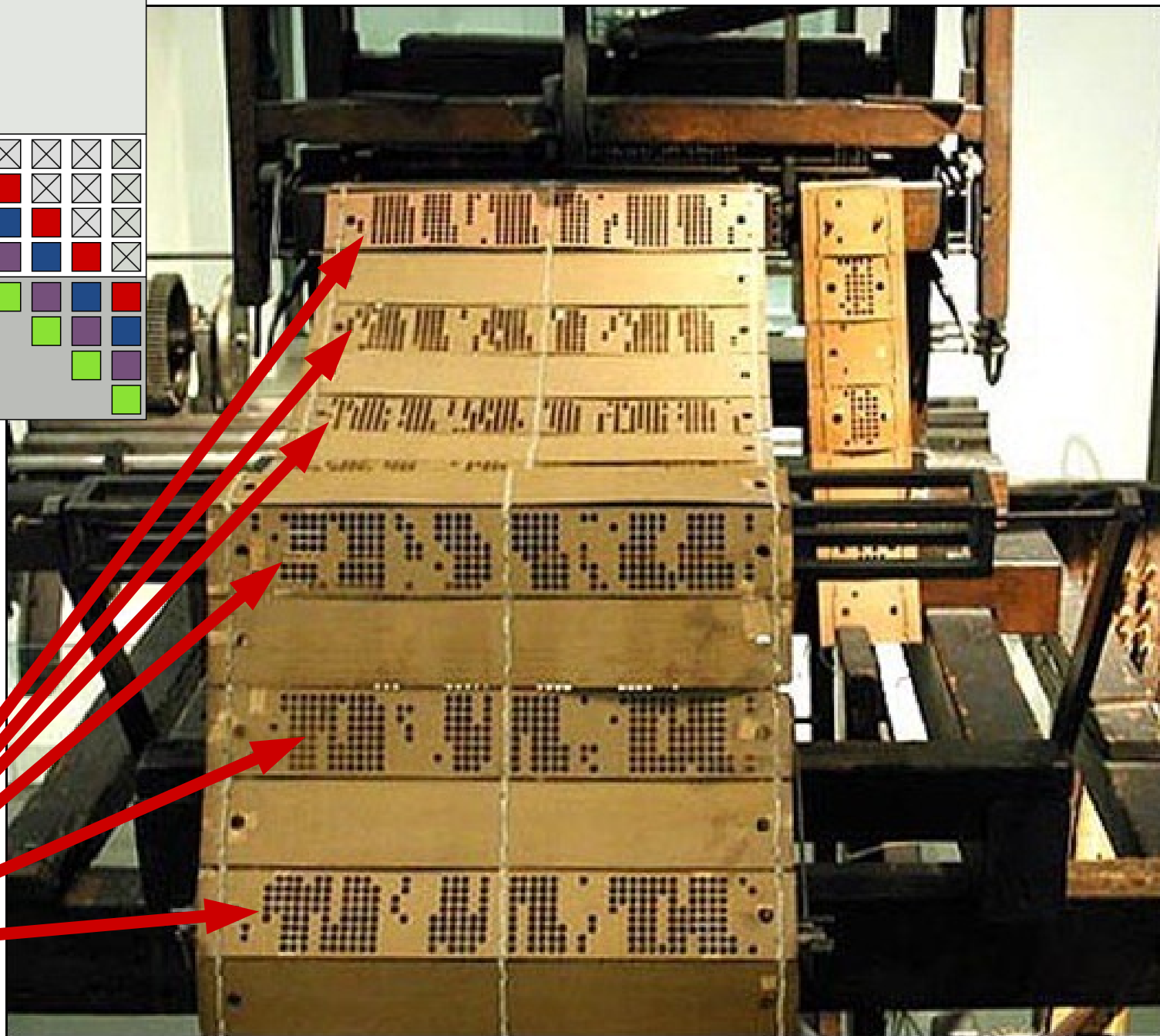
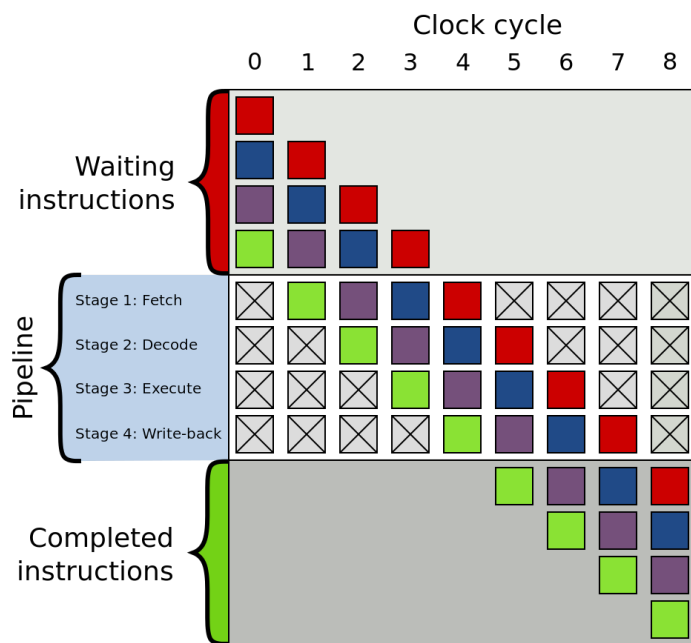
- Пиковая (теоретическая) производительность (Tflops)
- Linpack производительность (Tflops)
- VogoMIPS производительность?
- Пропускная способность Interconnect (Gbps)
- Задержки Interconnect (ns)

Основные характеристики вычислительного кластера

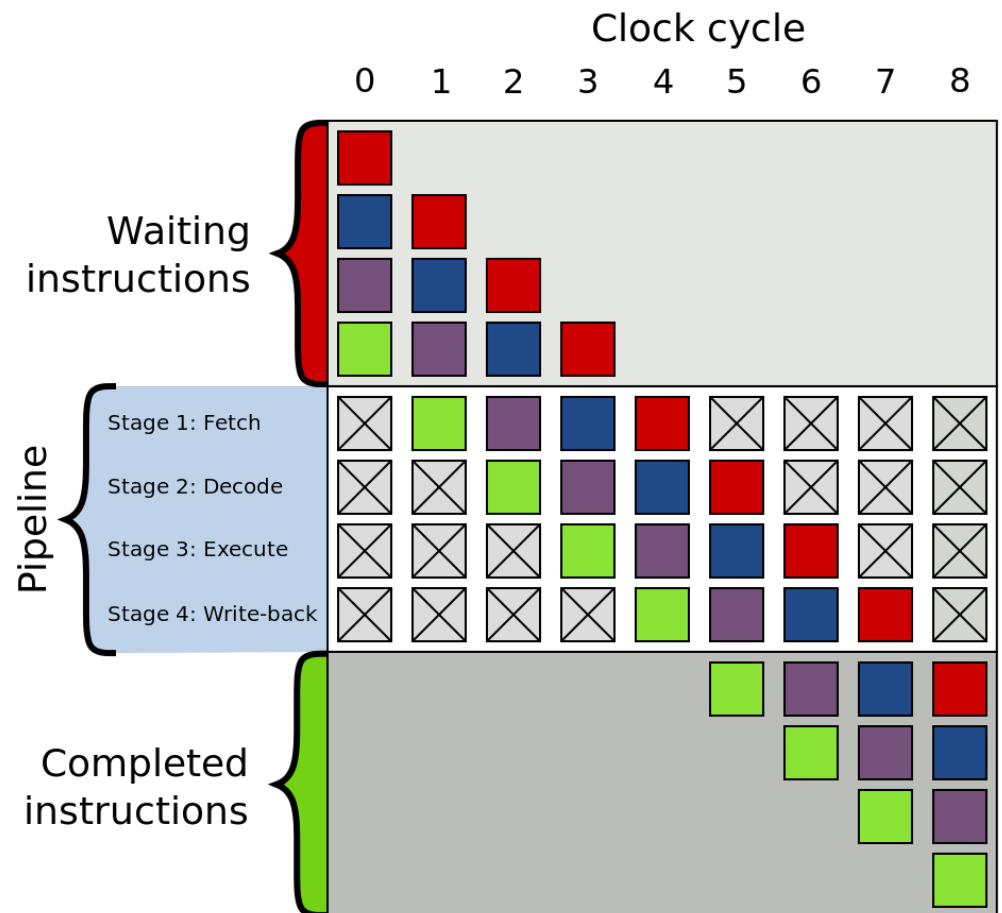
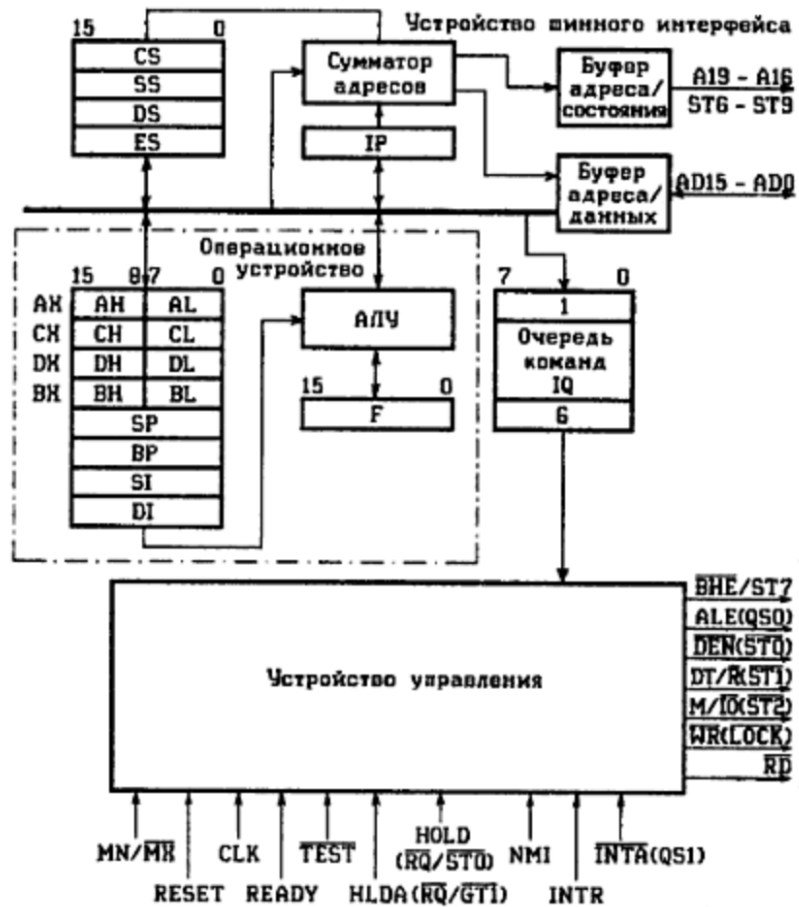
- Пиковая (теоретическая) производительность (Tflops)
- Linpack производительность (Tflops)
- VogoMIPS производительность?
- Пропускная способность Interconnect (Gbps)
- Задержки Interconnect (ns)

Remember?





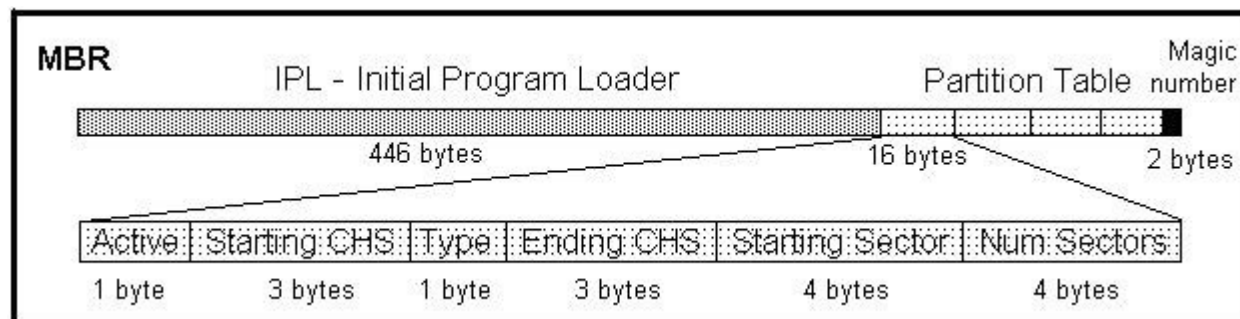
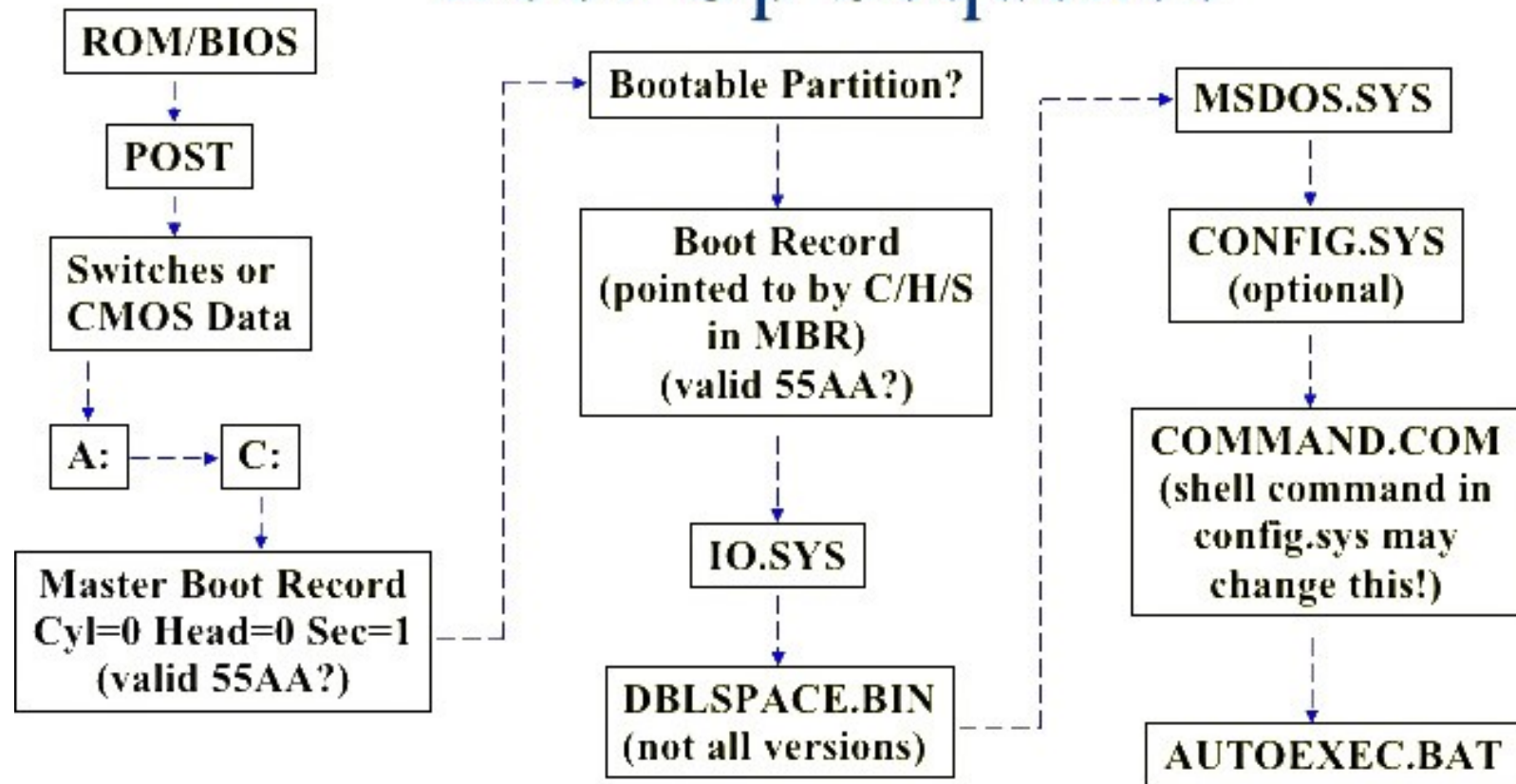
**Waiting
instructions**



We have a hardware, what is next?

We have a hardware, what is next?

Boot-Up Sequence



Disk Operating System

```
Welcome to FreeDOS

CuteMouse v1.9.1 alpha 1 [FreeDOS]
Installed at PS/2 port
C:\>ver

FreeCom version 0.82 pl 3 XMS_Swap [Dec 10 2003 06:49:21]

C:\>dir
Volume in drive C is FREEDOS_C95
Volume Serial Number is 0E4F-19EB
Directory of C:\

FDOS                <DIR>    08-26-04   6:23p
AUTOEXEC BAT        435    08-26-04   6:24p
BOOTSECT BIN        512    08-26-04   6:23p
COMMAND COM        93,963  08-26-04   6:24p
CONFIG SYS          801    08-26-04   6:24p
FDOSBOOT BIN        512    08-26-04   6:24p
KERNEL SYS         45,815  04-17-04   9:19p
        6 file(s)          142,038 bytes
        1 dir(s)    1,064,517,632 bytes free

C:\>_
```

Disk Operating System

```
Welcome to FreeDOS

CuteMouse v1.9.1 alpha 1 [FreeDOS]
Installed at PS/2 port
C:\>ver

FreeCom version 0.82 pl 3 XMS_Swap [Dec 10 2003 06:49:21]

C:\>dir
Volume in drive C is FREEDOS_C95
Volume Serial Number is 0E4F-19EB
Directory of C:\

FDOS                <DIR>    08-26-04   6:23p
AUTOEXEC.BAT        435    08-26-04   6:24p
BOOTSECT.BIN        512    08-26-04   6:23p
COMMAND.COM         93,963  08-26-04   6:24p
CONFIG.SYS          801    08-26-04   6:24p
FDOSBOOT.BIN        512    08-26-04   6:24p
KERNEL.SYS         45,815  04-17-04   9:19p
        6 file(s)        142,038 bytes
        1 dir(s)    1,064,517,632 bytes free

C:\>_
```

- MS DOS memory limit (640KB ought to be enough for anybody)
- No memory protection, first-fit memory allocation scheme
- Single task, but TSR (Terminate and Stay Resident)
- Single user

Disk Operating System

```
Welcome to FreeDOS

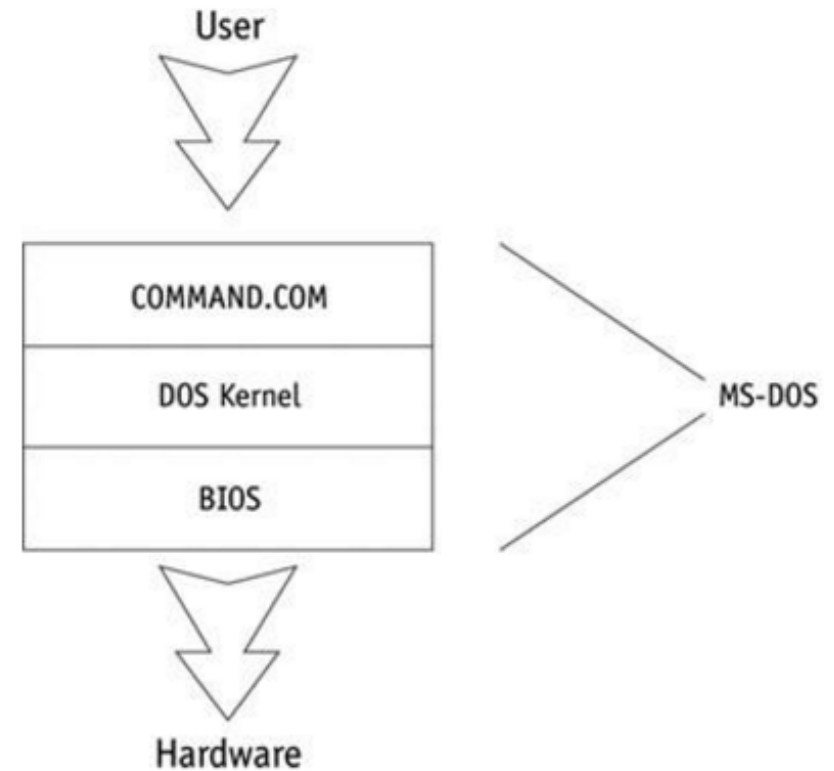
CuteMouse v1.9.1 alpha 1 [FreeDOS]
Installed at PS/2 port
C:\>ver

FreeCom version 0.82 pl 3 XMS_Swap [Dec 10 2003 06:49:21]

C:\>dir
Volume in drive C is FREEDOS_C95
Volume Serial Number is 0E4F-19EB
Directory of C:\

FDOS          <DIR>    08-26-04  6:23p
AUTOEXEC.BAT      435    08-26-04  6:24p
BOOTSECT.BIN      512    08-26-04  6:23p
COMMAND.COM     93,963    08-26-04  6:24p
CONFIG.SYS       801    08-26-04  6:24p
FDOSBOOT.BIN      512    08-26-04  6:24p
KERNEL.SYS     45,815    04-17-04  9:19p
        6 file(s)      142,038 bytes
        1 dir(s)    1,064,517,632 bytes free

C:\>_
```



BIOS

- **BIOS (Basic Input/Output System)**
 - Direct interface with I/O devices
 - Contains device drivers
 - Controls data flow to and from each device (except disk drives)
 - Receives I/O operation status information
 - Passes to processor
 - Handles small differences among I/O units
 - No need to write device driver for manufacturer printer

DOS kernel

- **DOS kernel**
 - Contains routines to interface with disk drives
 - Read into memory
 - Initialization time from MSDOS.SYS file
 - Resides in boot disk
 - Microsoft proprietary program
 - Accessed by application programs
 - Provides hardware-independent services
 - System functions
 - Memory management, file and record management
 - Provides transparency
 - Compensates for manufacturer variations
 - Manages file storage and retrieval
 - Dynamically allocates and deallocates secondary storage as needed

Command processor

- **Command processor (shell)**
 - Sends prompts to user
 - Accepts typed commands
 - Executes commands
 - From **system prompt**
 - Issues appropriate responses
 - Resides in COMMAND.COM file
 - Stored in two different main memory sections
 - Appears on public directory
 - Weakness: not interpretive

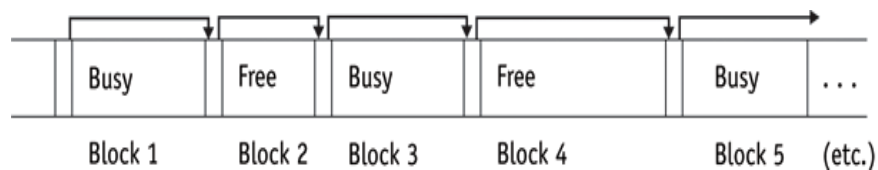
- Main memory structure

- **ROM**

- Very small in size
- Contains program
- Contains section of BIOS with startup process (bootstrapping)
- Initializes computer
- Retrieves resident code and loads into RAM

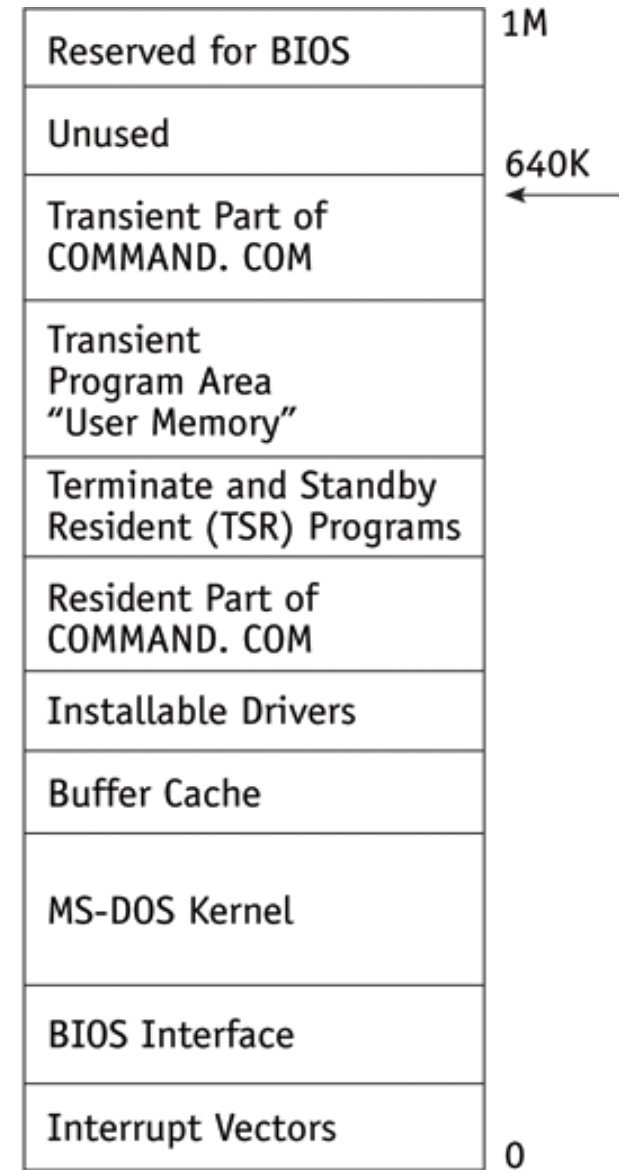
- **RAM**

- Part of main memory
- Where programs are loaded and executed



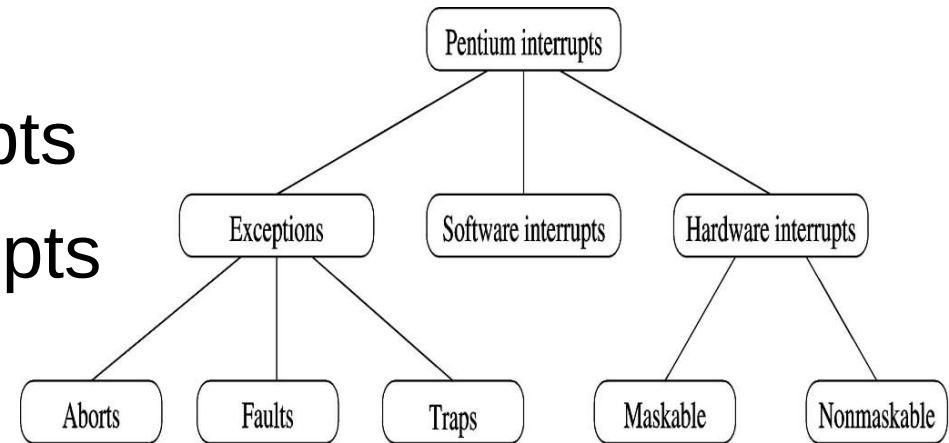
(figure 14.4)

The linked list of memory blocks.

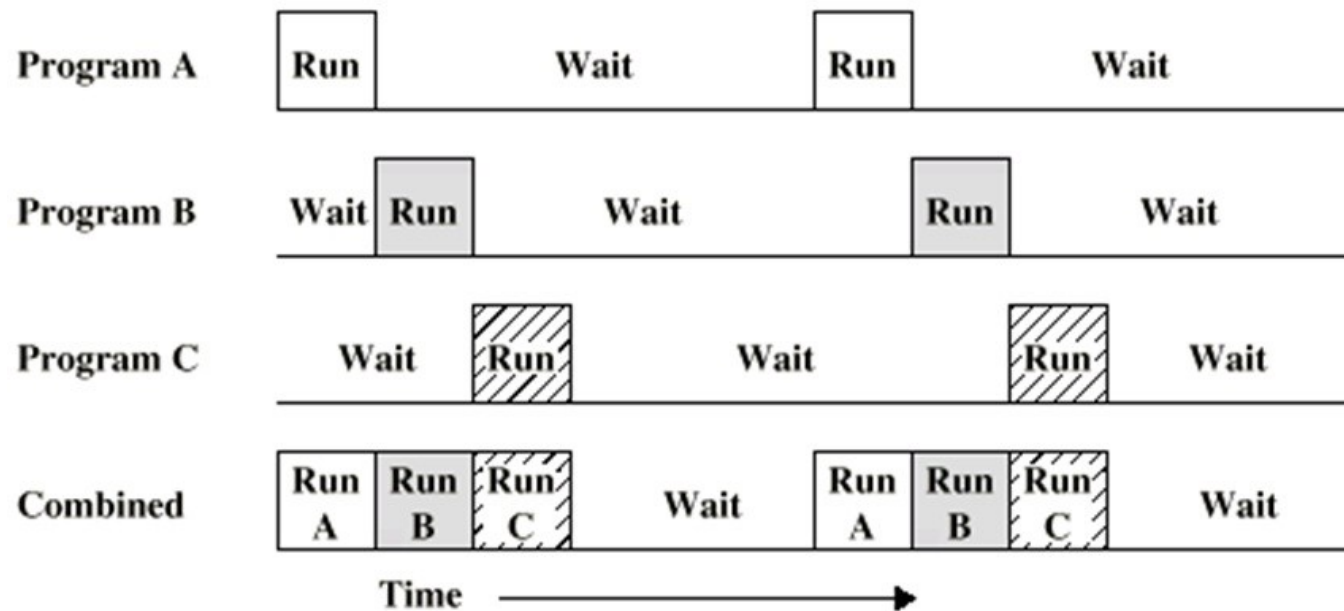


Single user, single task

- Interrupts:
 - Internal hardware interrupts
 - External hardware interrupts
 - Software interrupts



Multitasking:

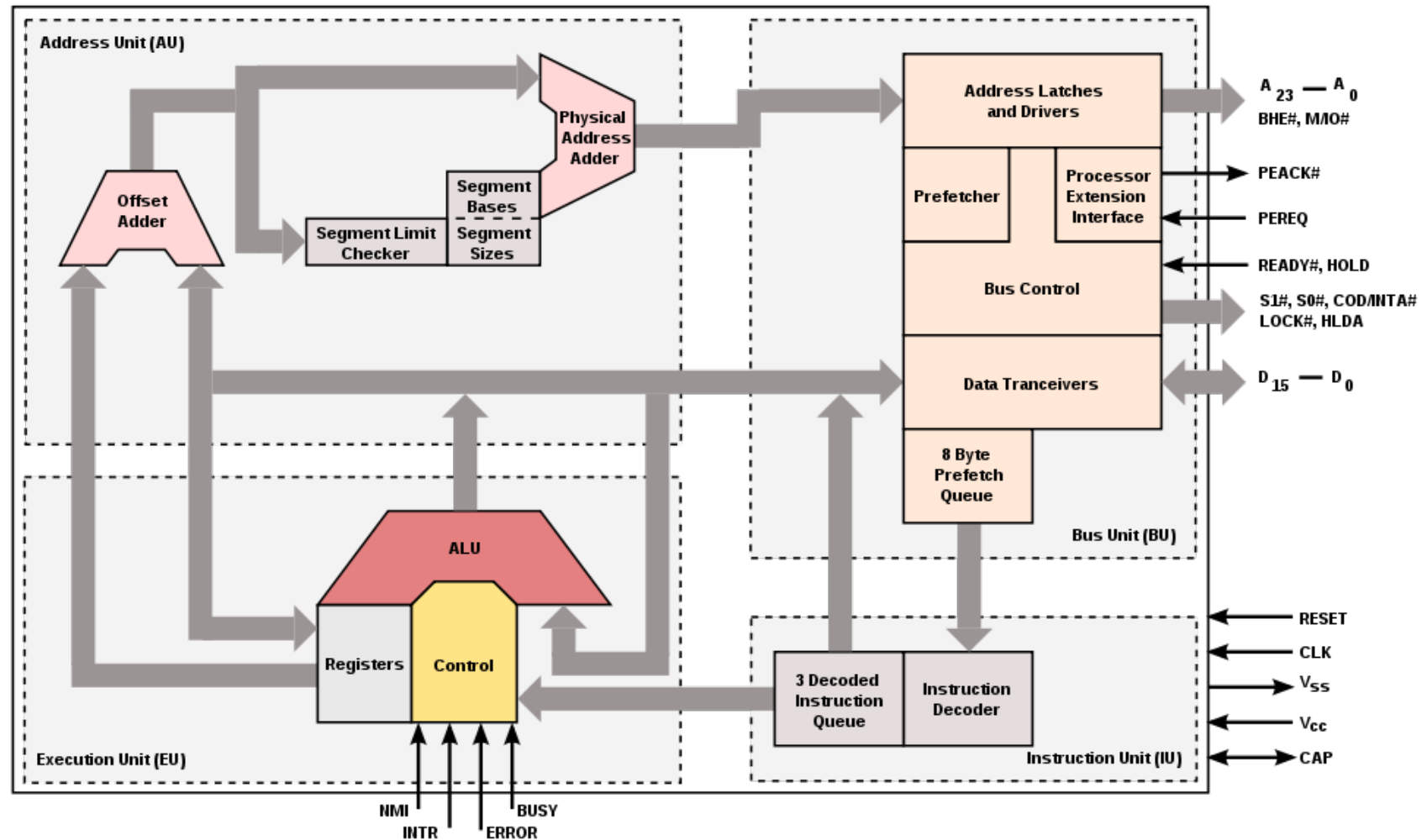


- Single user, single task
- Unsafe, unsecure
- A lot of code duplication
- etc...

Required to add hardware features!

- 80286: protected mode, virtual memory, 20MHz

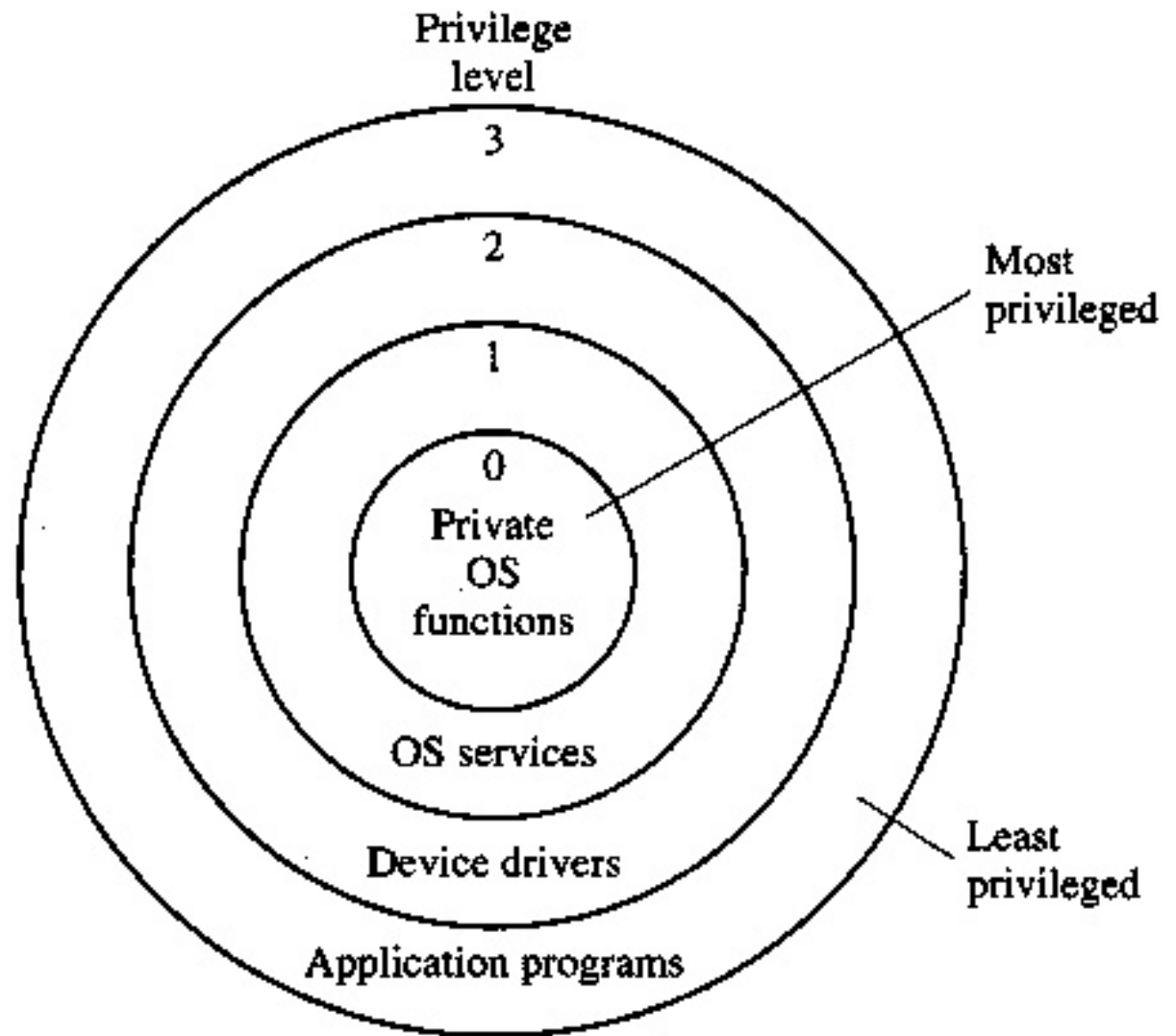
Intel 80286 architecture



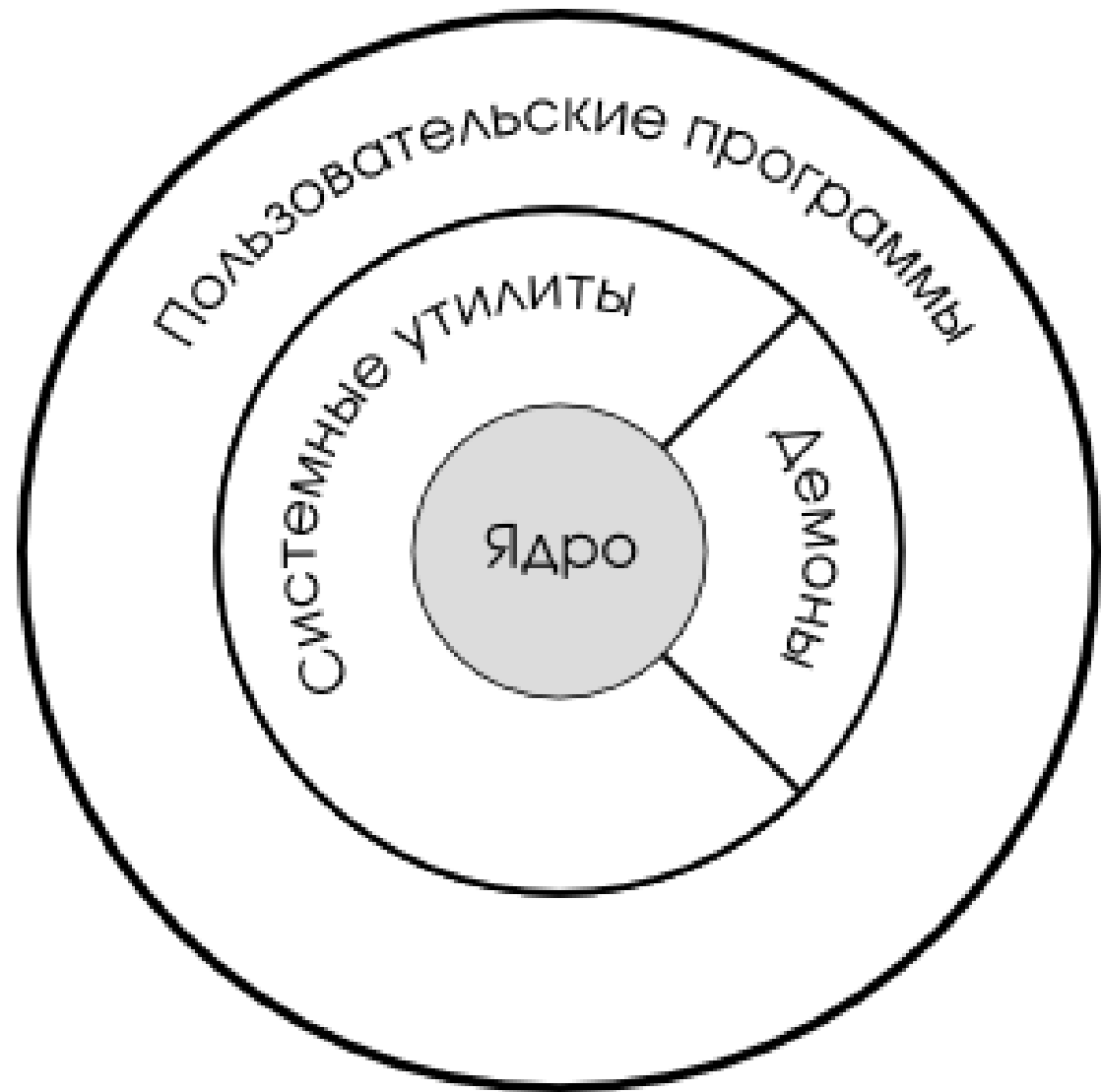
- protected mode, virtual memory

- protected mode, virtual memory

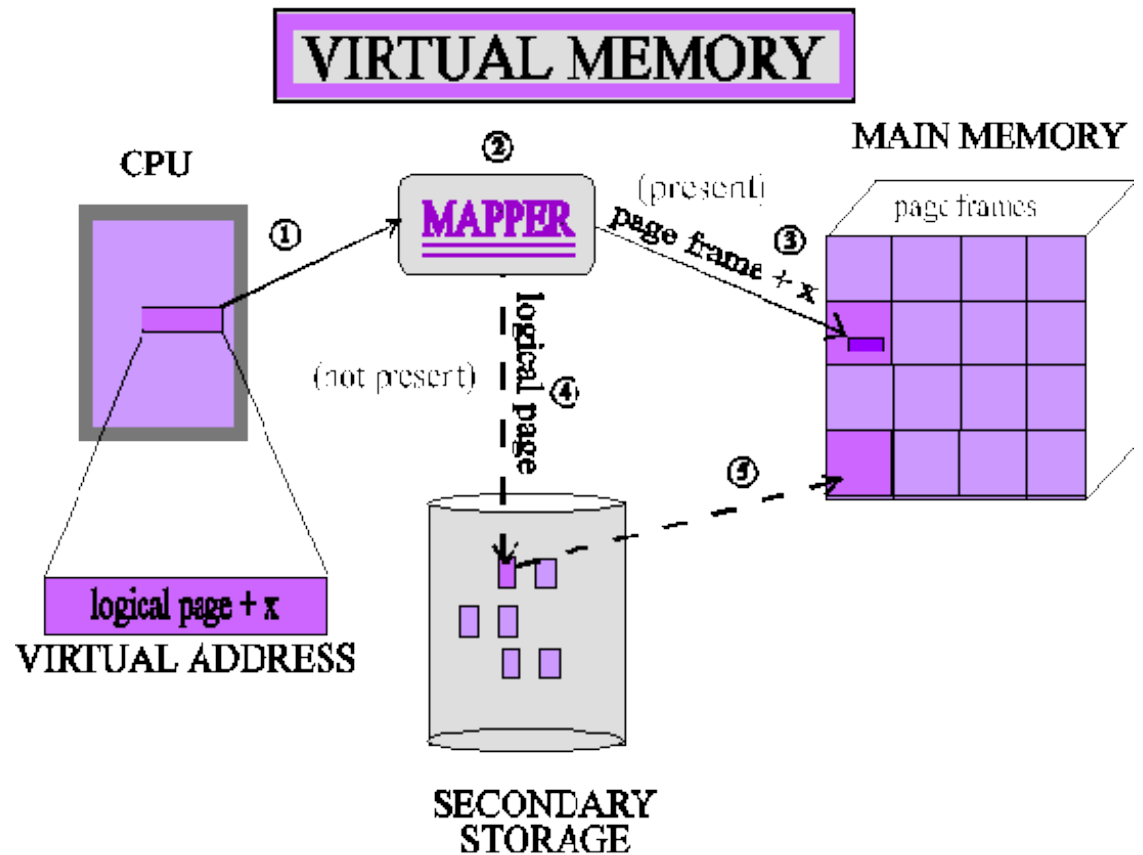
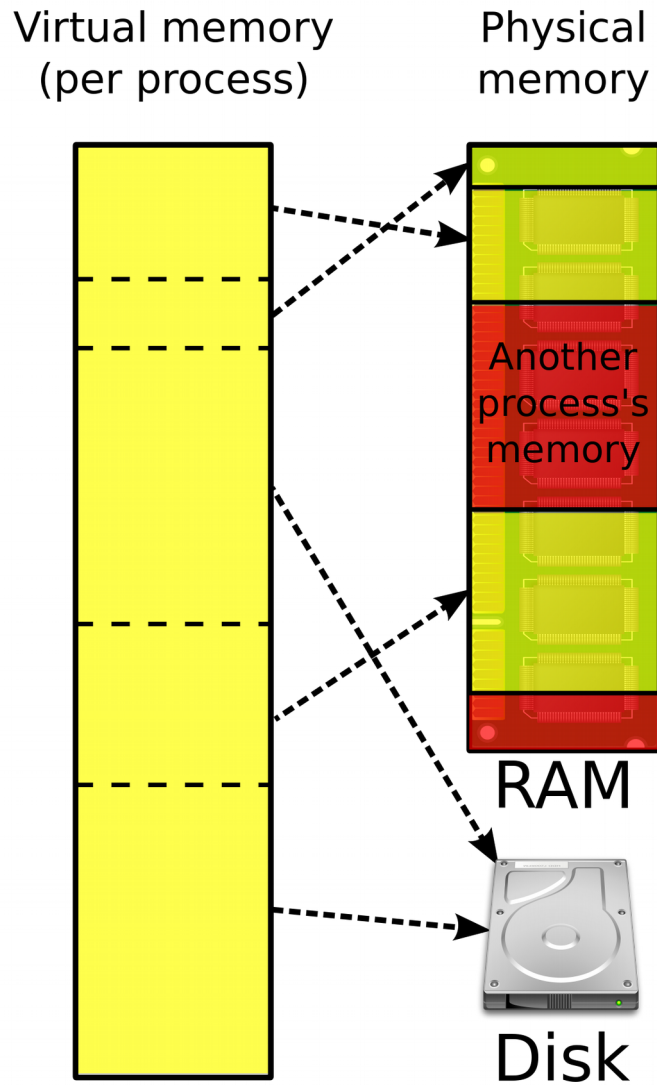
Hardware
level:

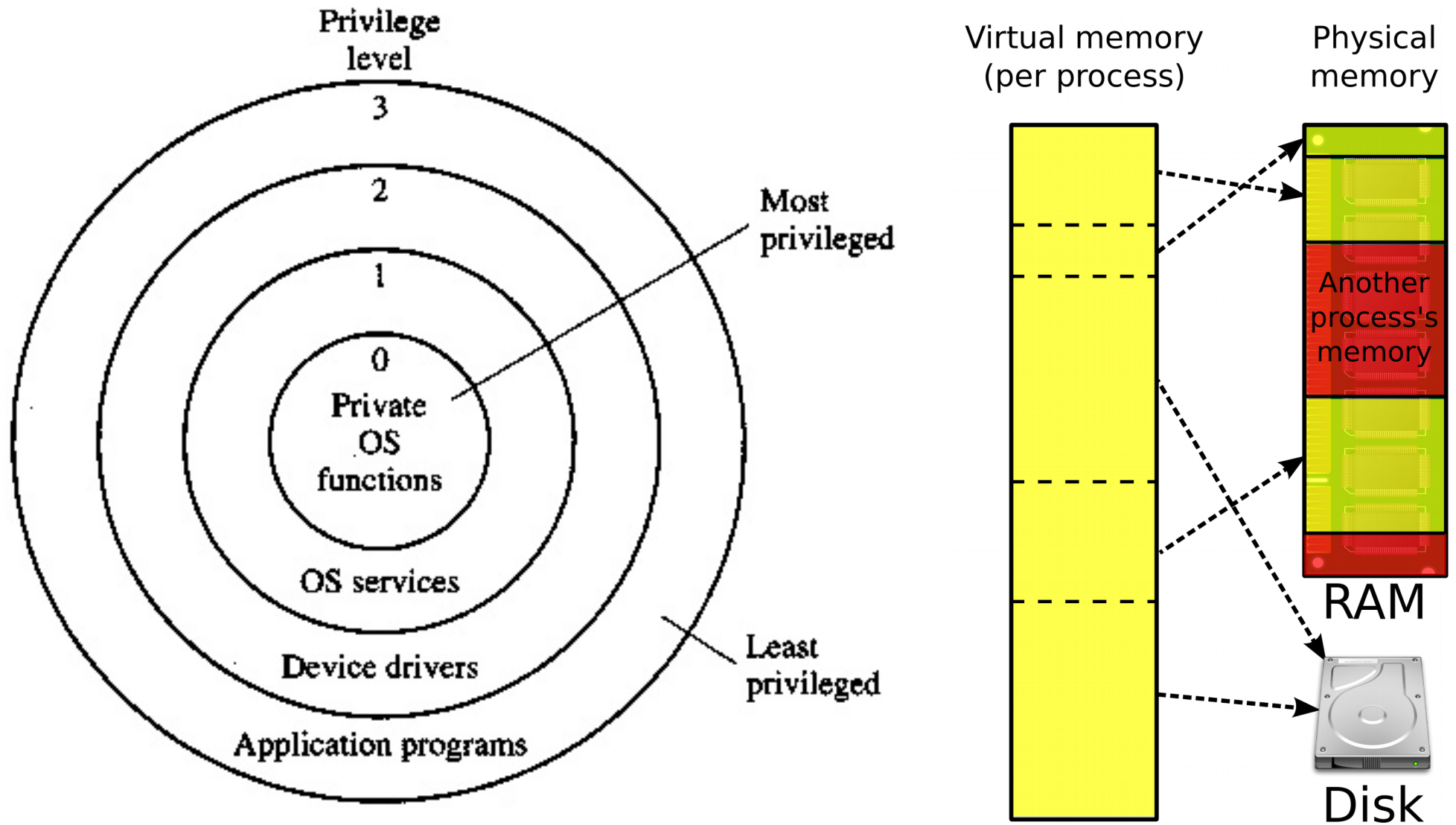


Software
level:



- protected mode, virtual memory

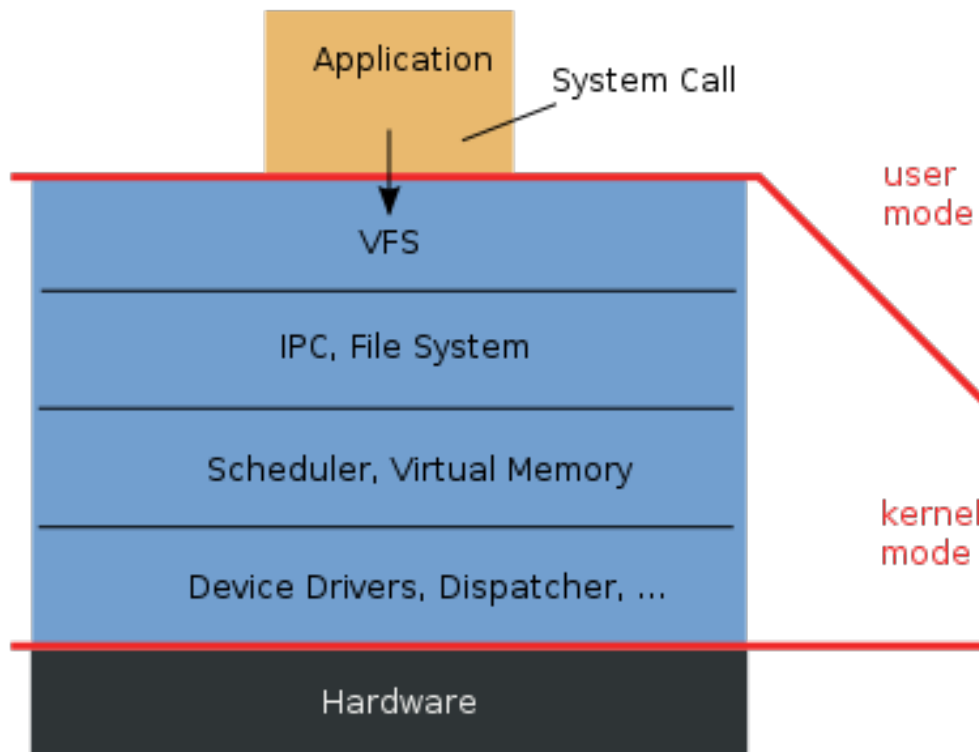




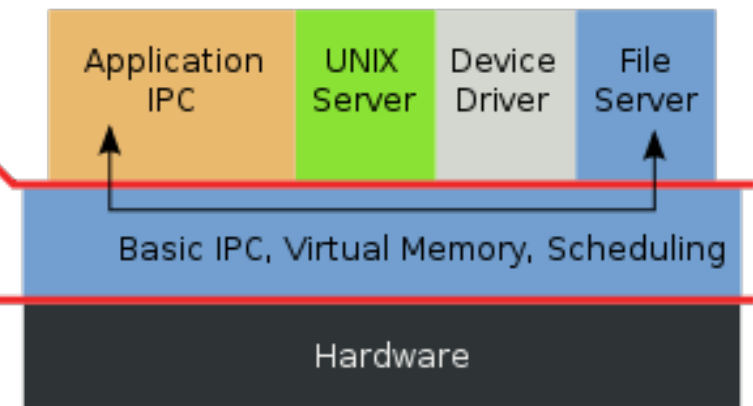
Again: we have a hardware,
what is next?

Monolithic vs Microkernel

Monolithic Kernel
based Operating System

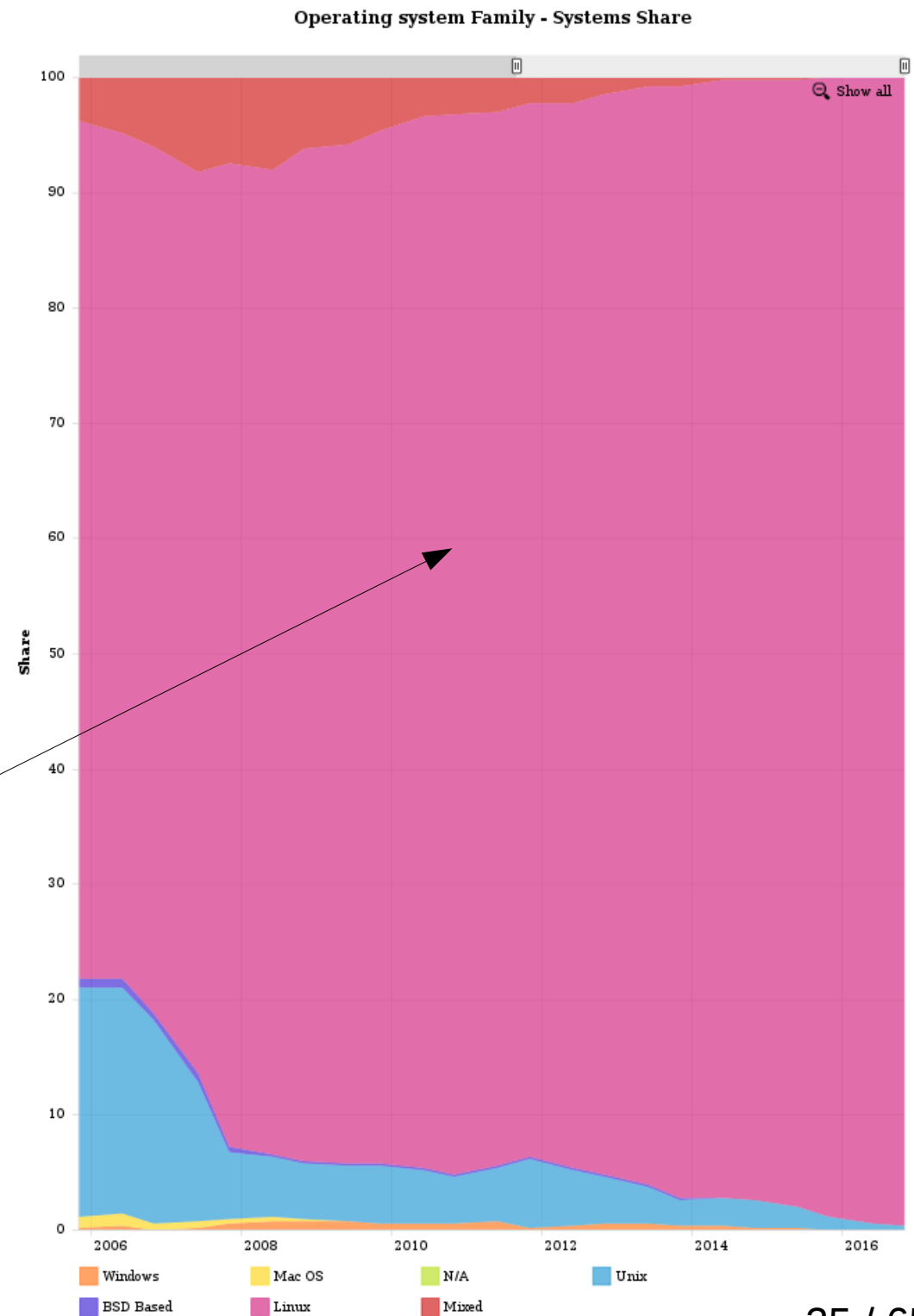


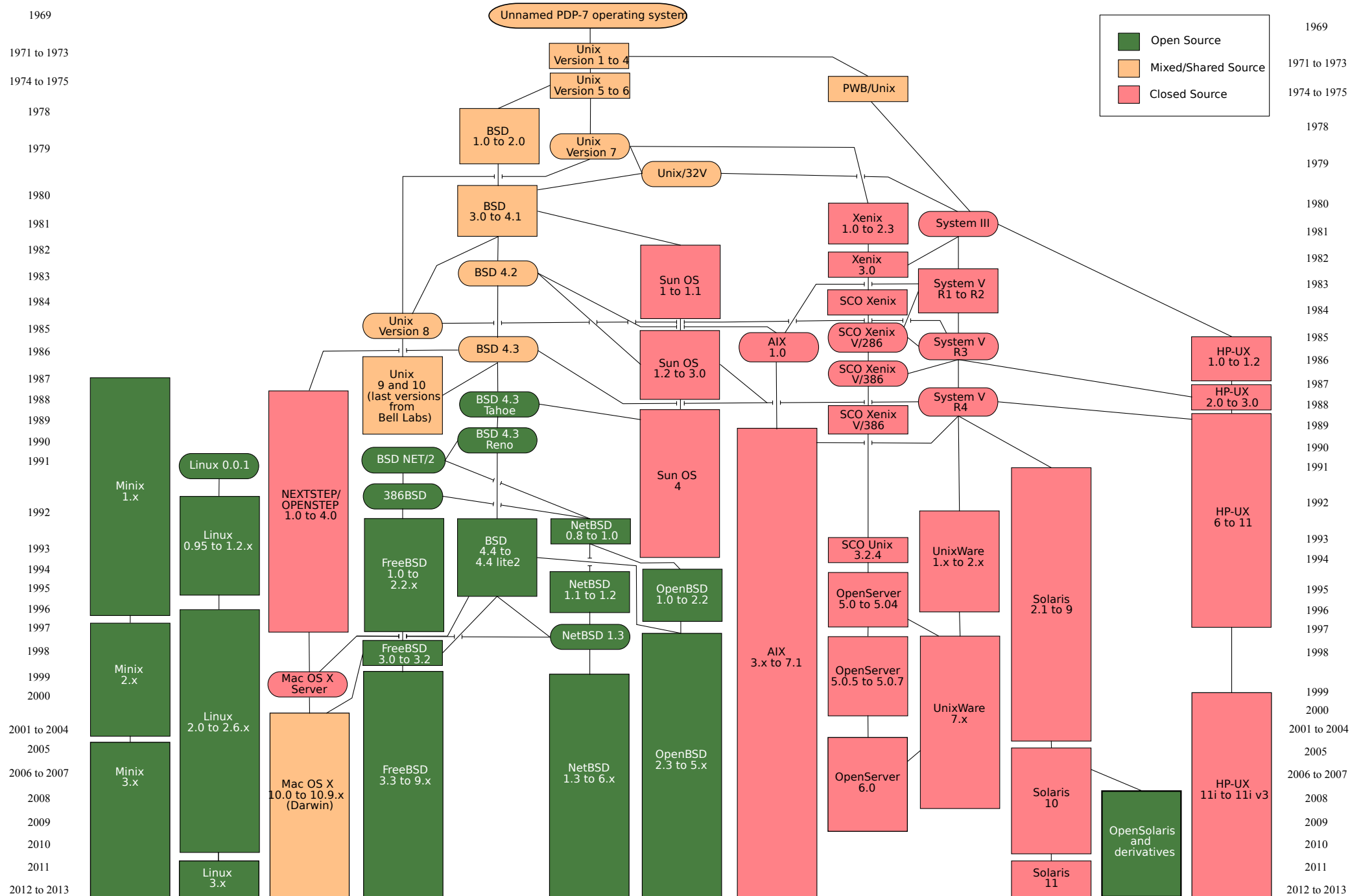
Microkernel
based Operating System



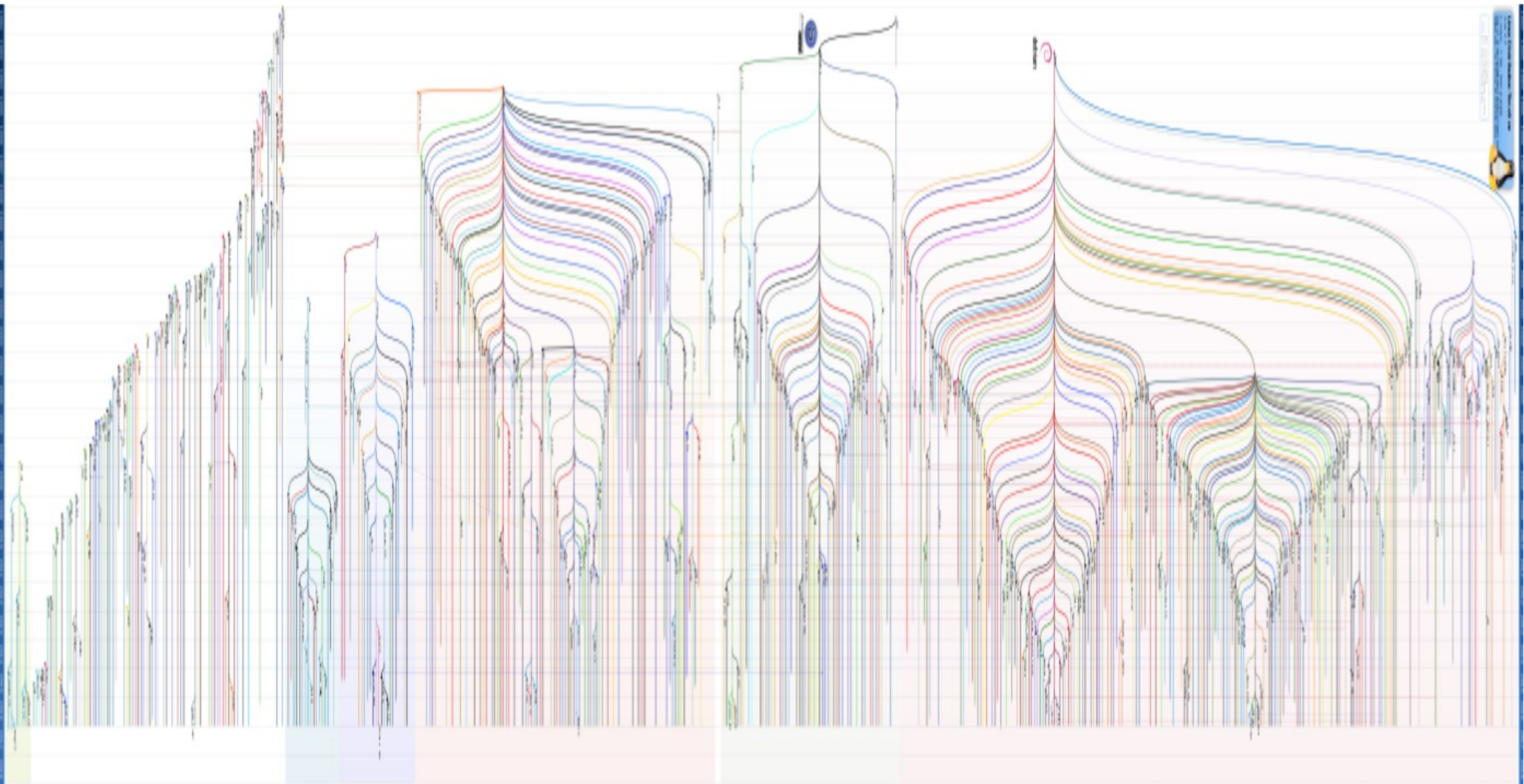
Tenenbaum-Torvalds debate

Linux (in HPC market share)





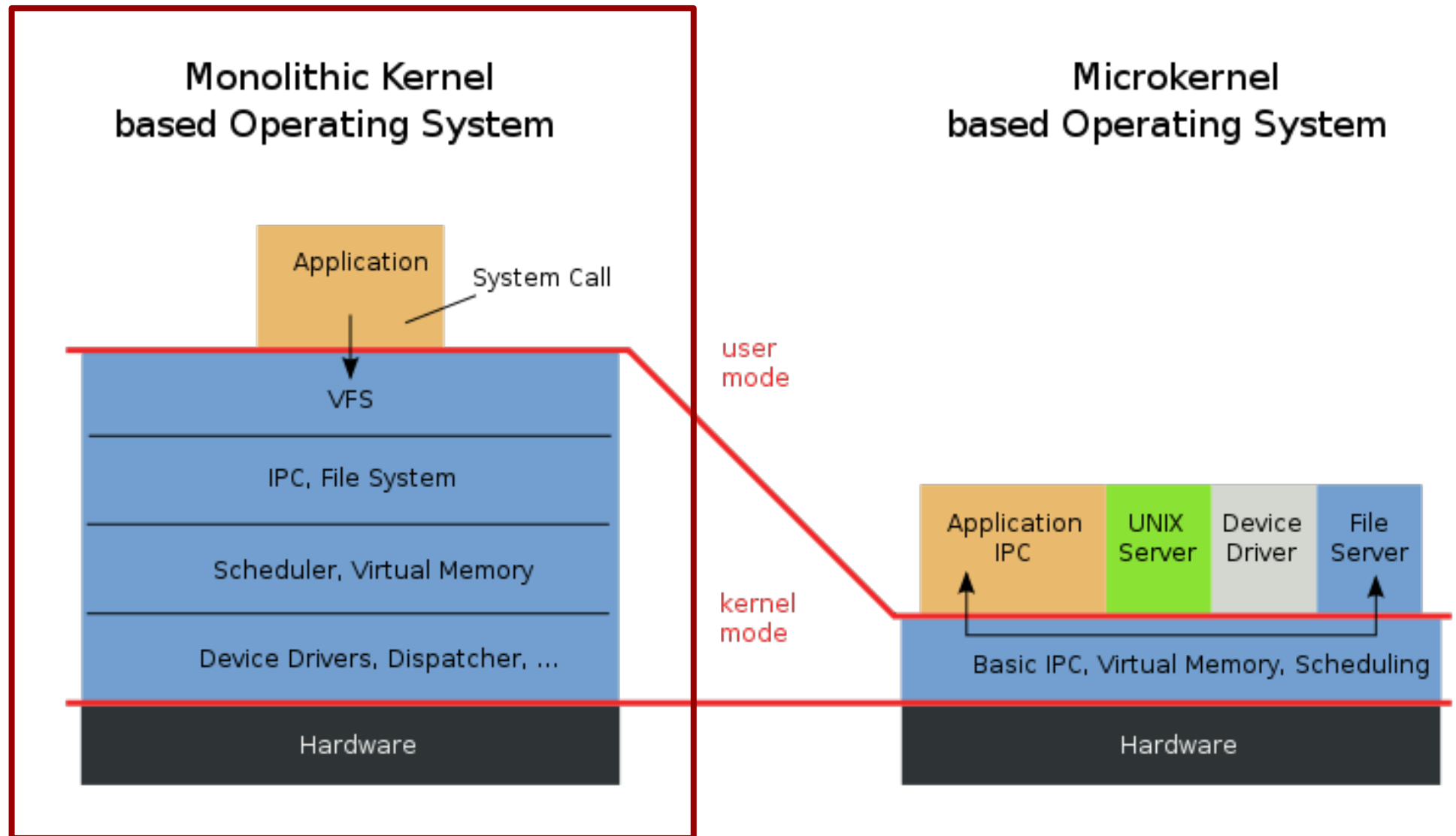
Linux distributions

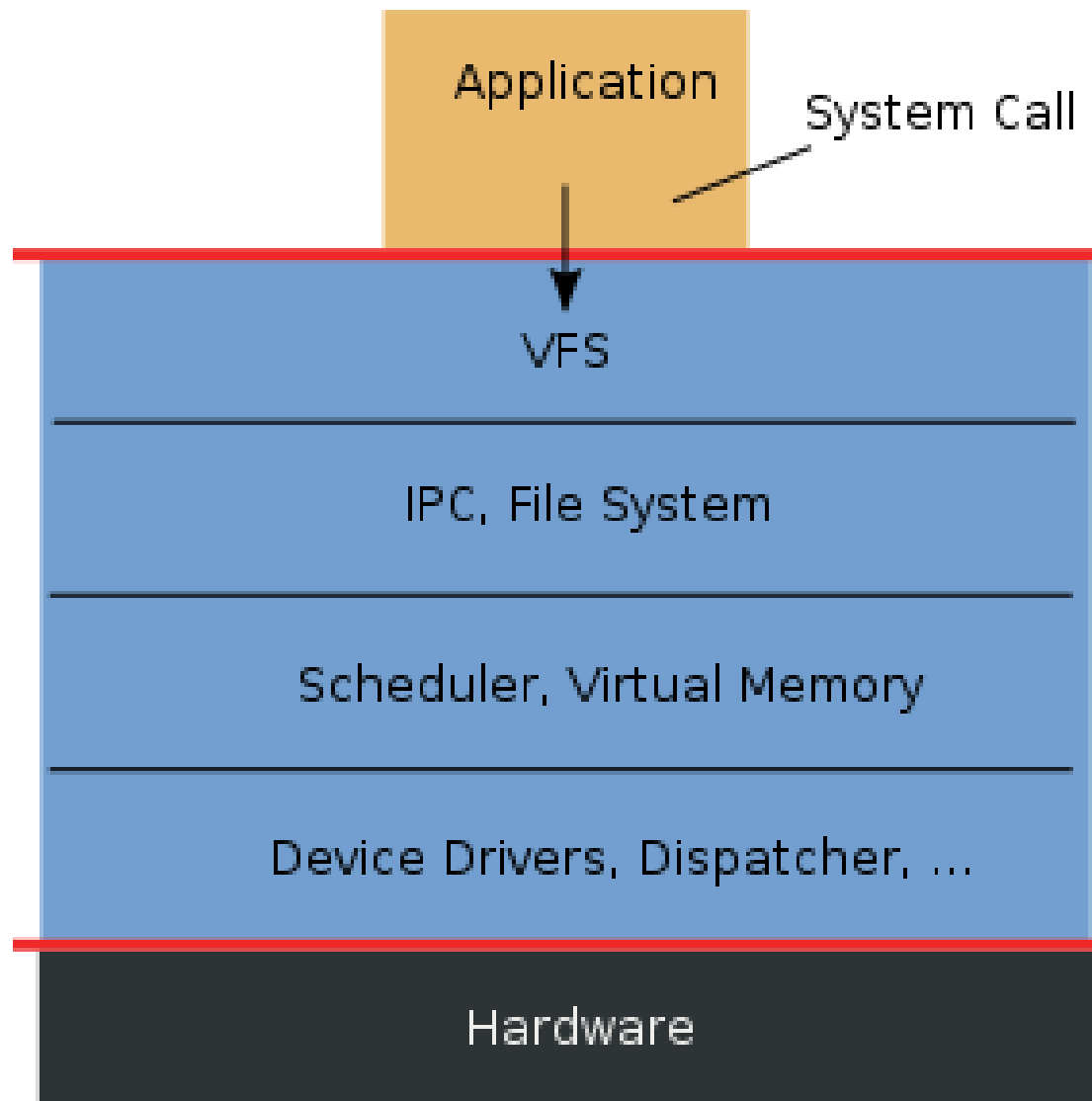


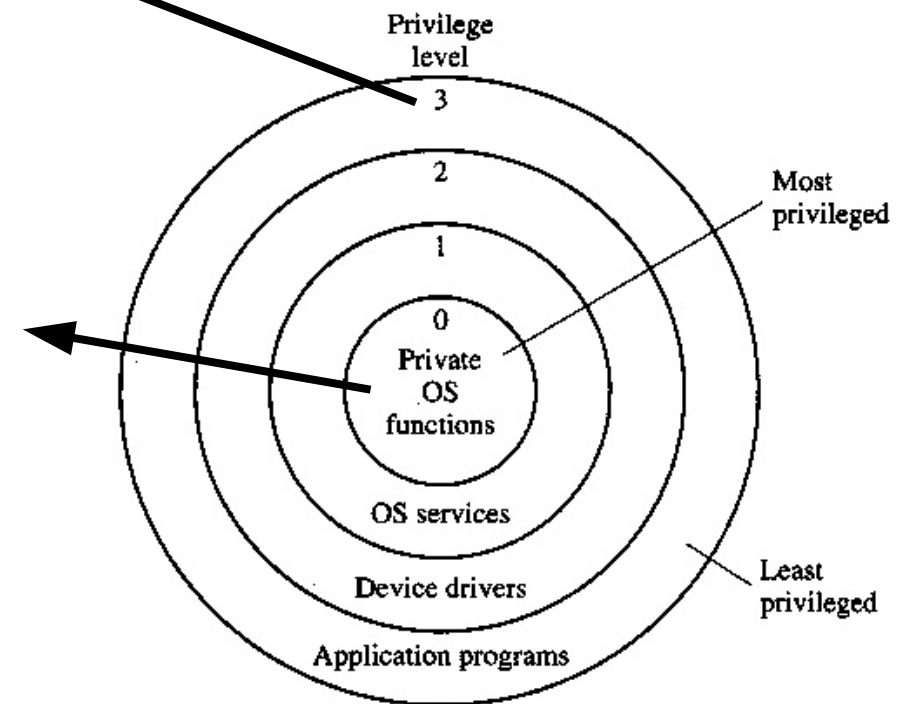
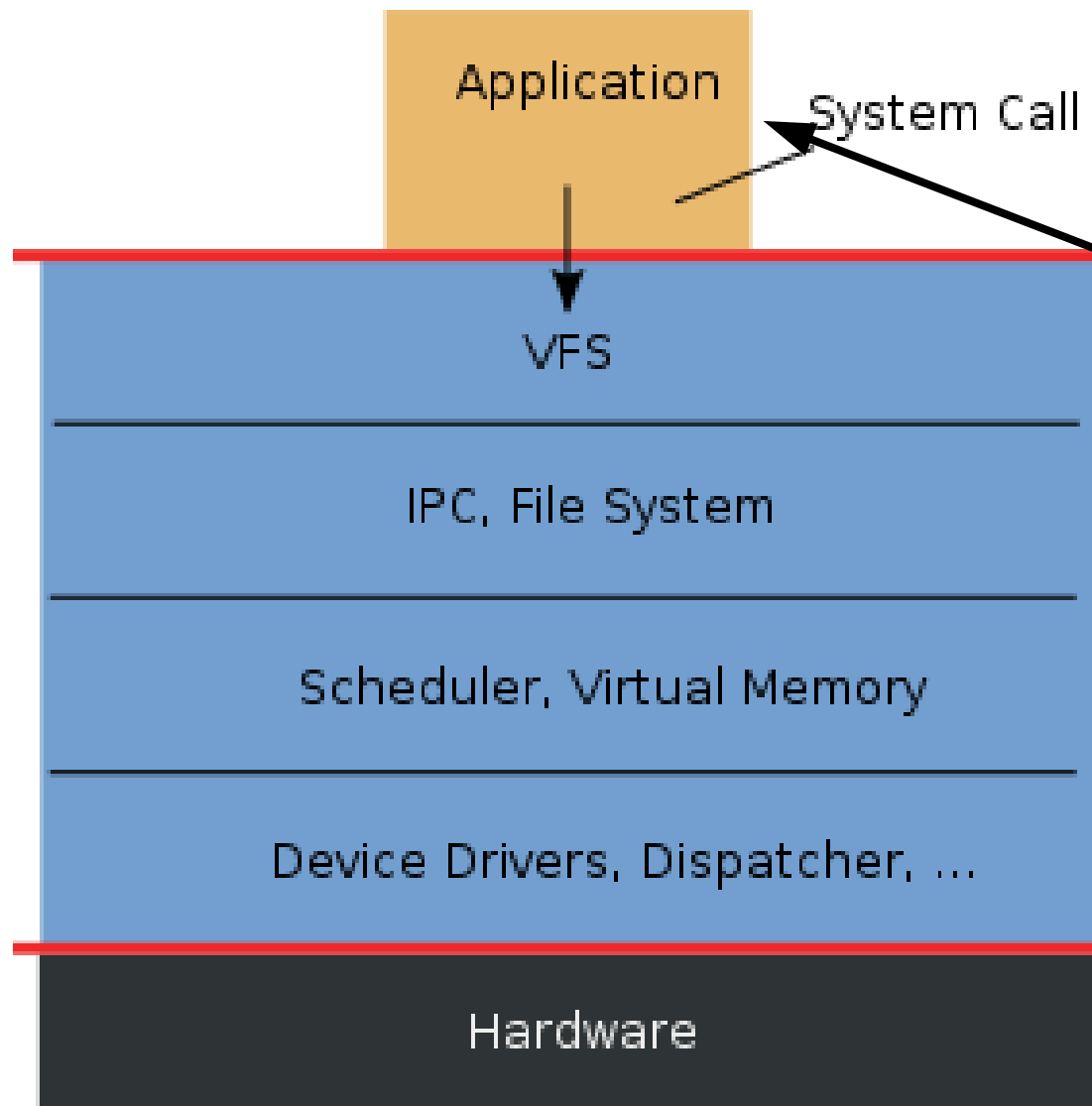
Must read

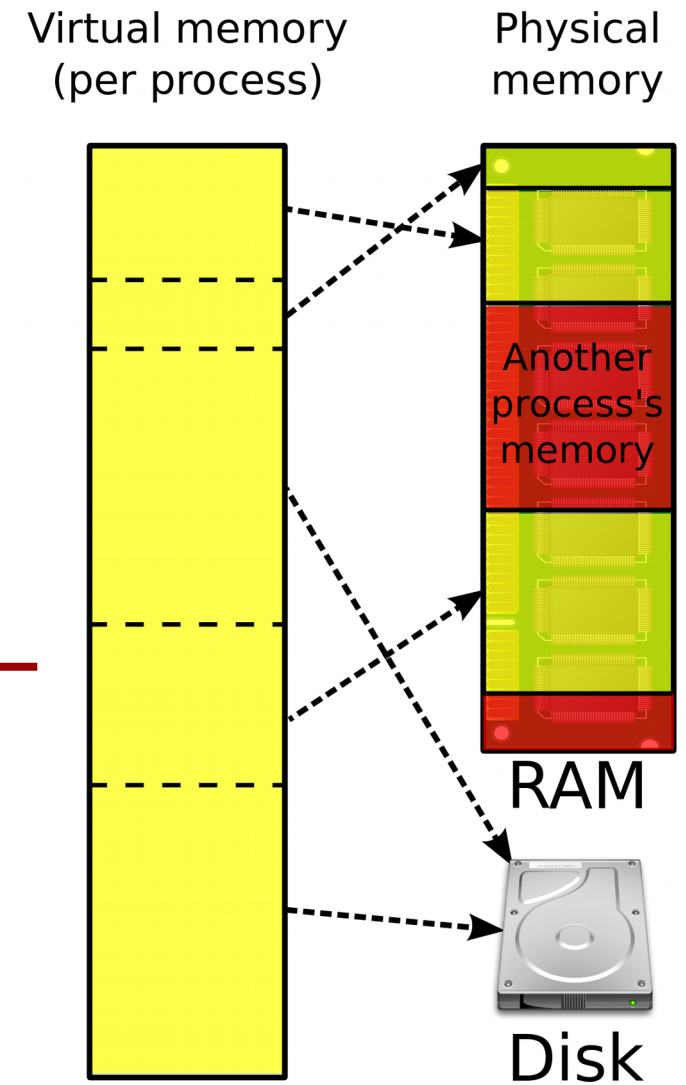
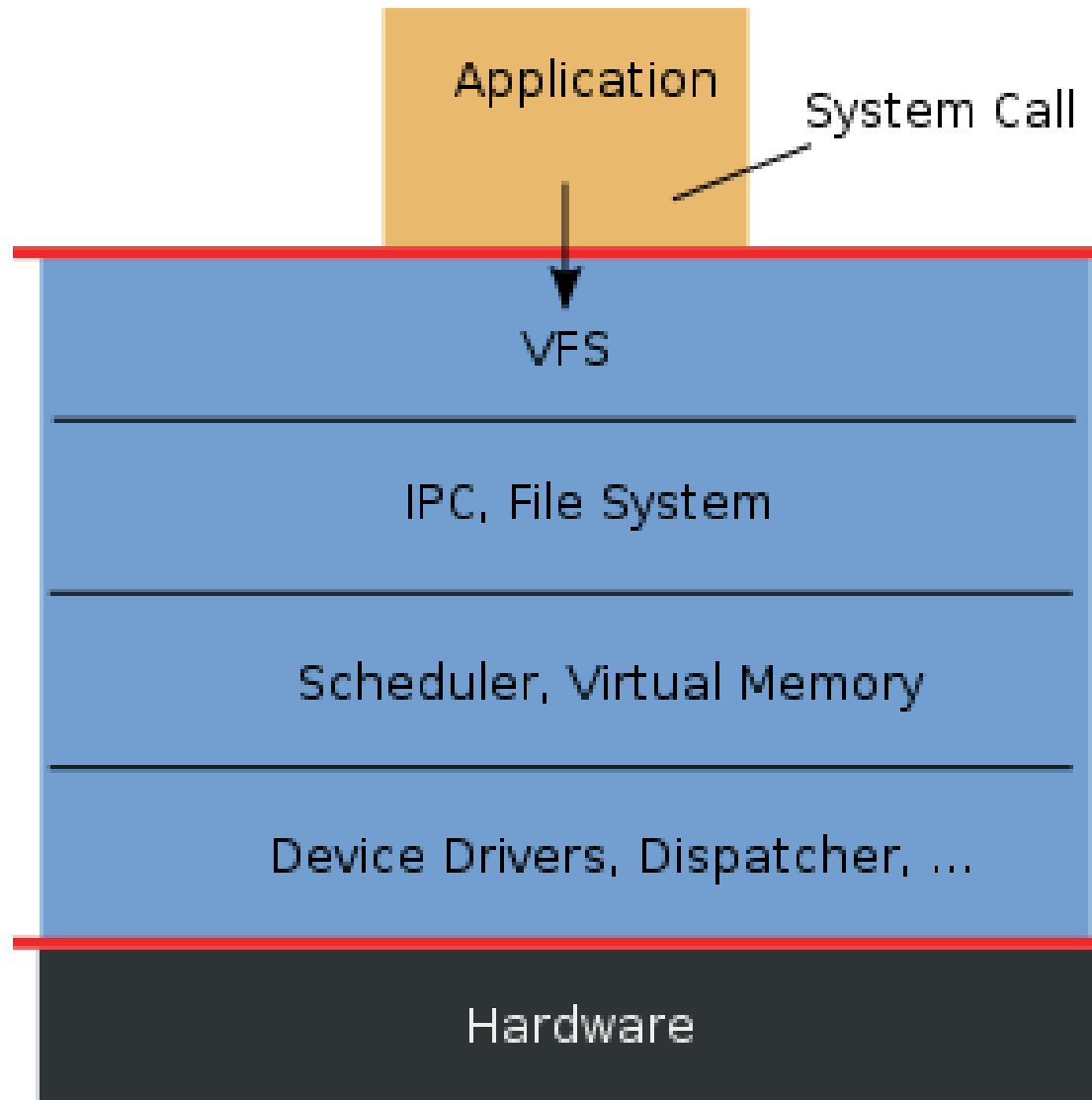
- Tenenbaum-Torvalds debate
- GNU C Library debates about Versioned Interfaces

Monolithic vs Microkernel

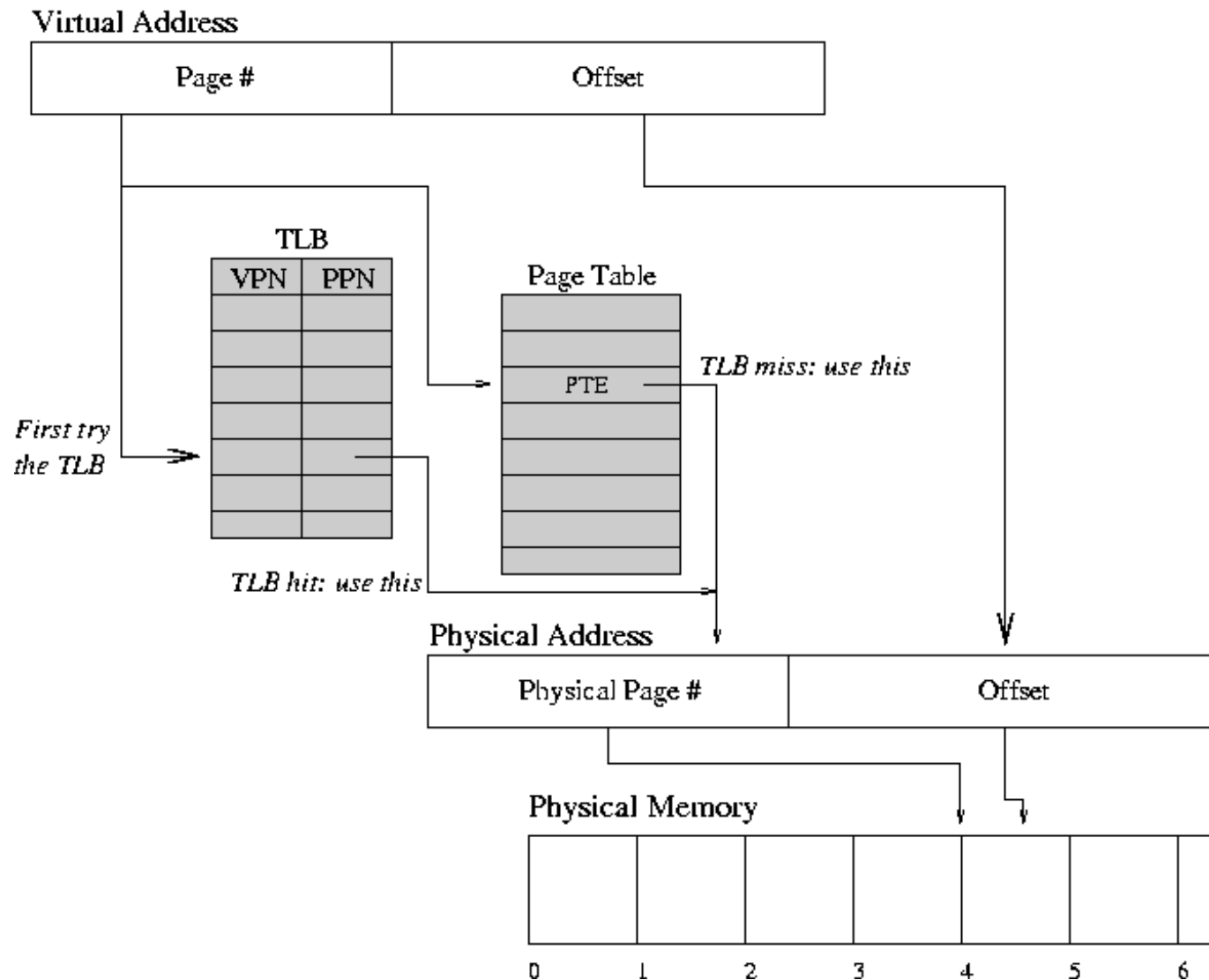




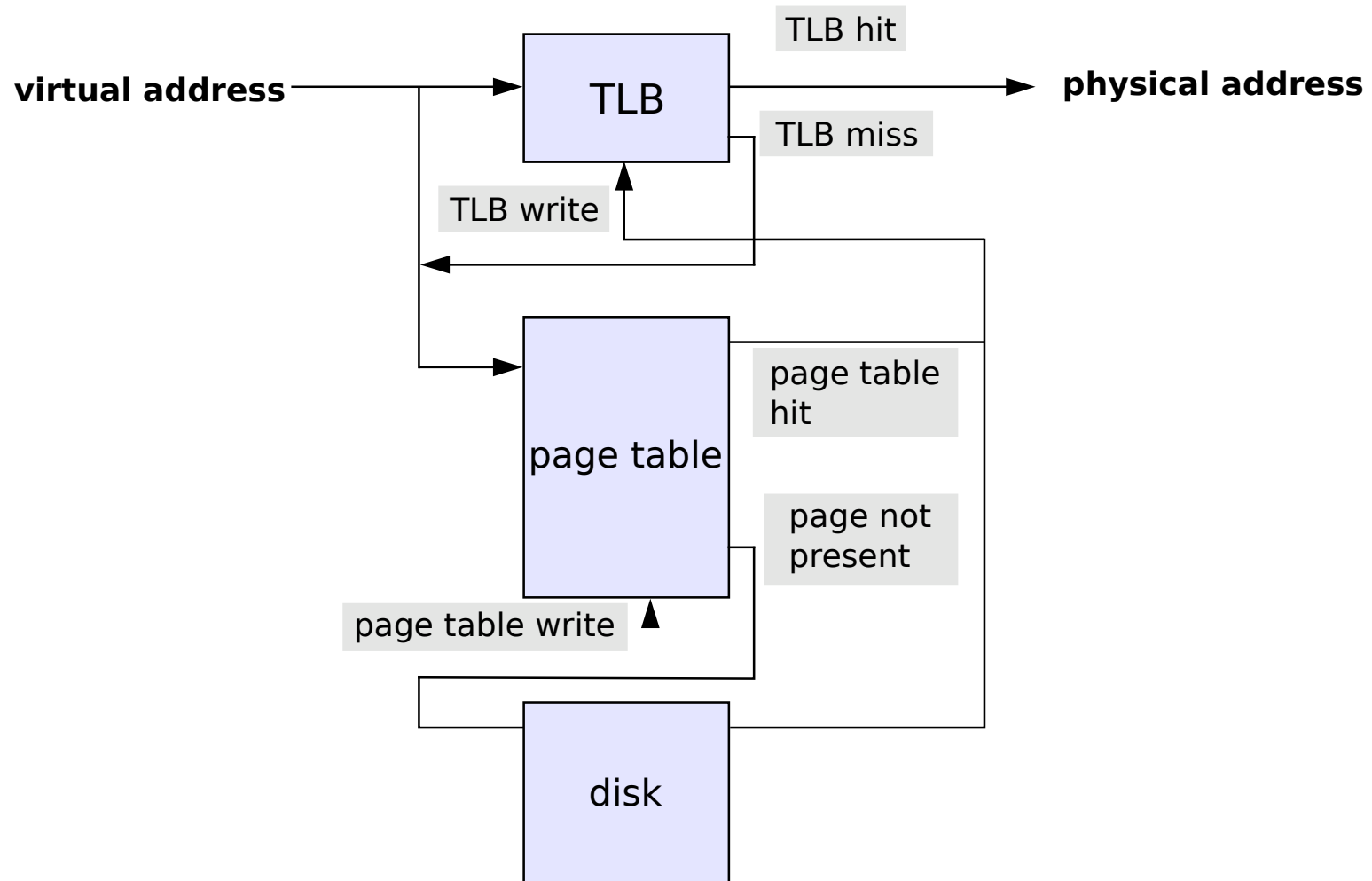




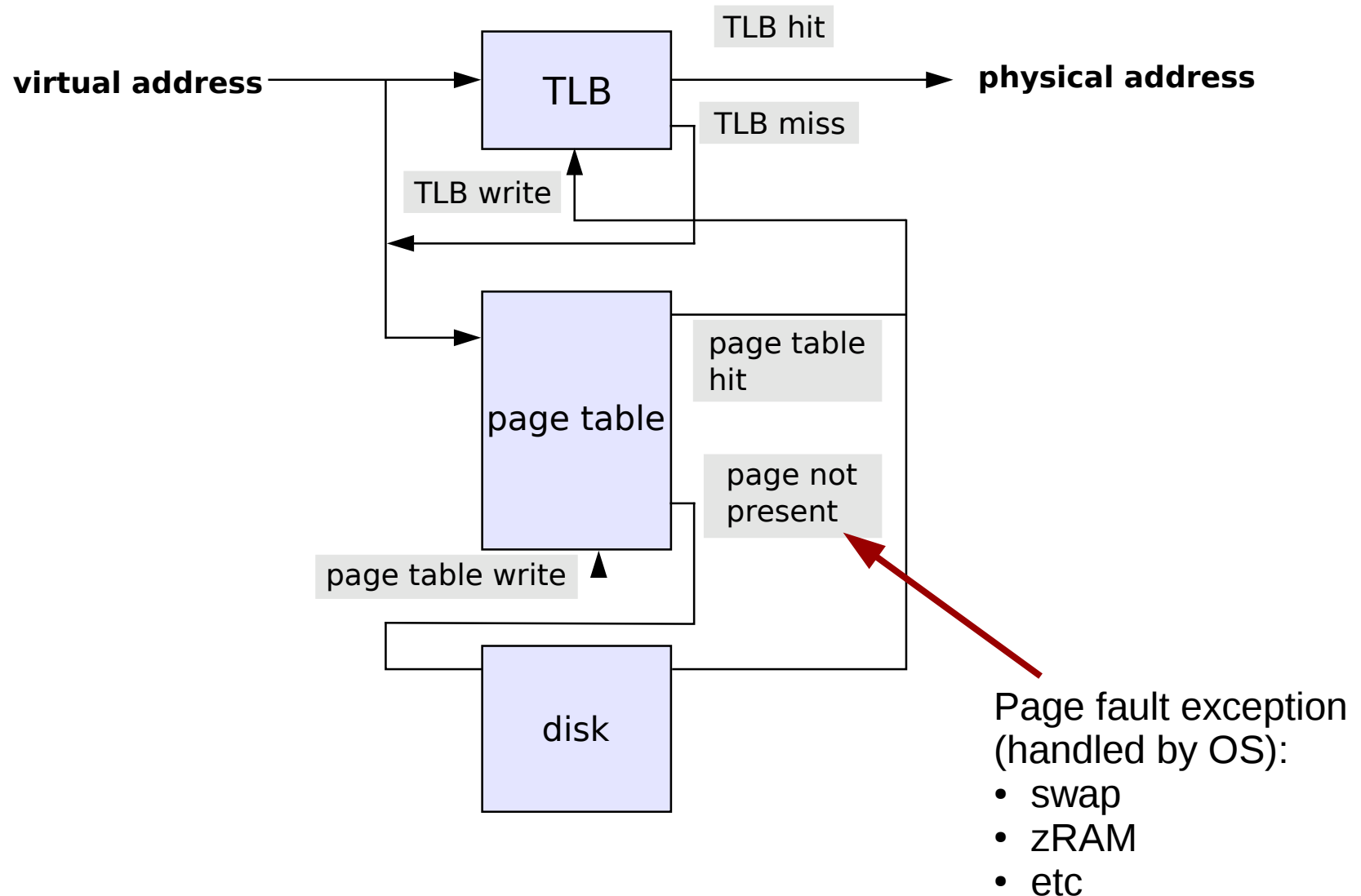
Virtual memory



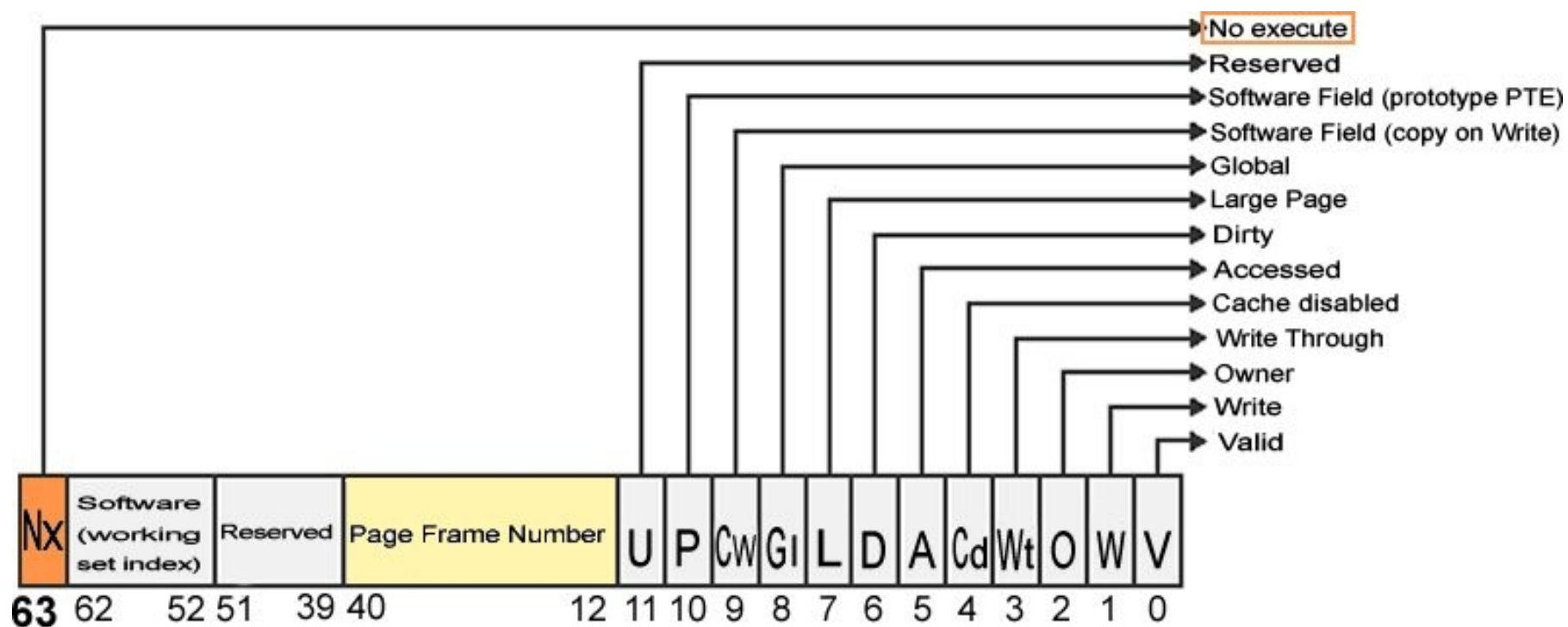
Virtual memory

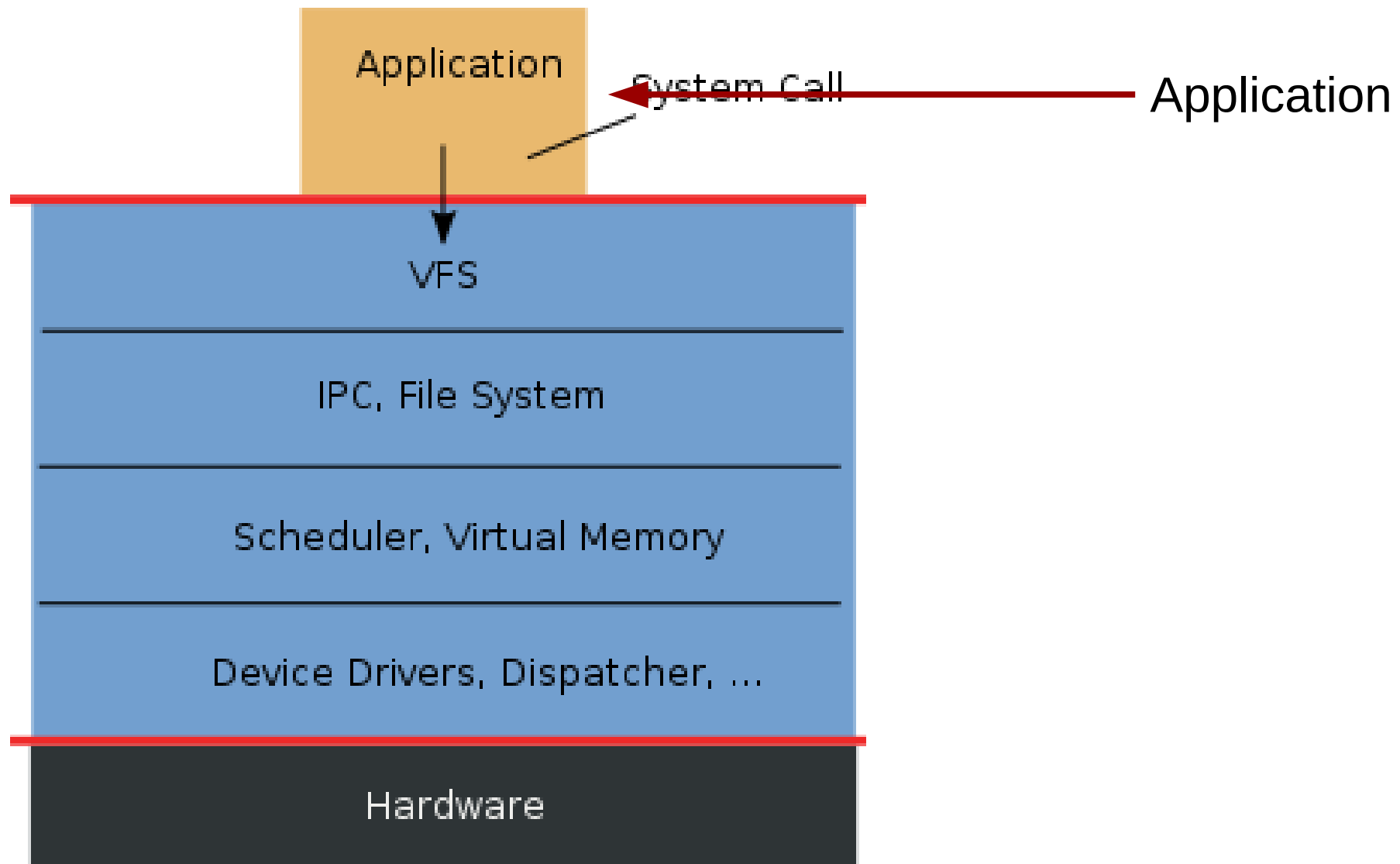


Virtual memory

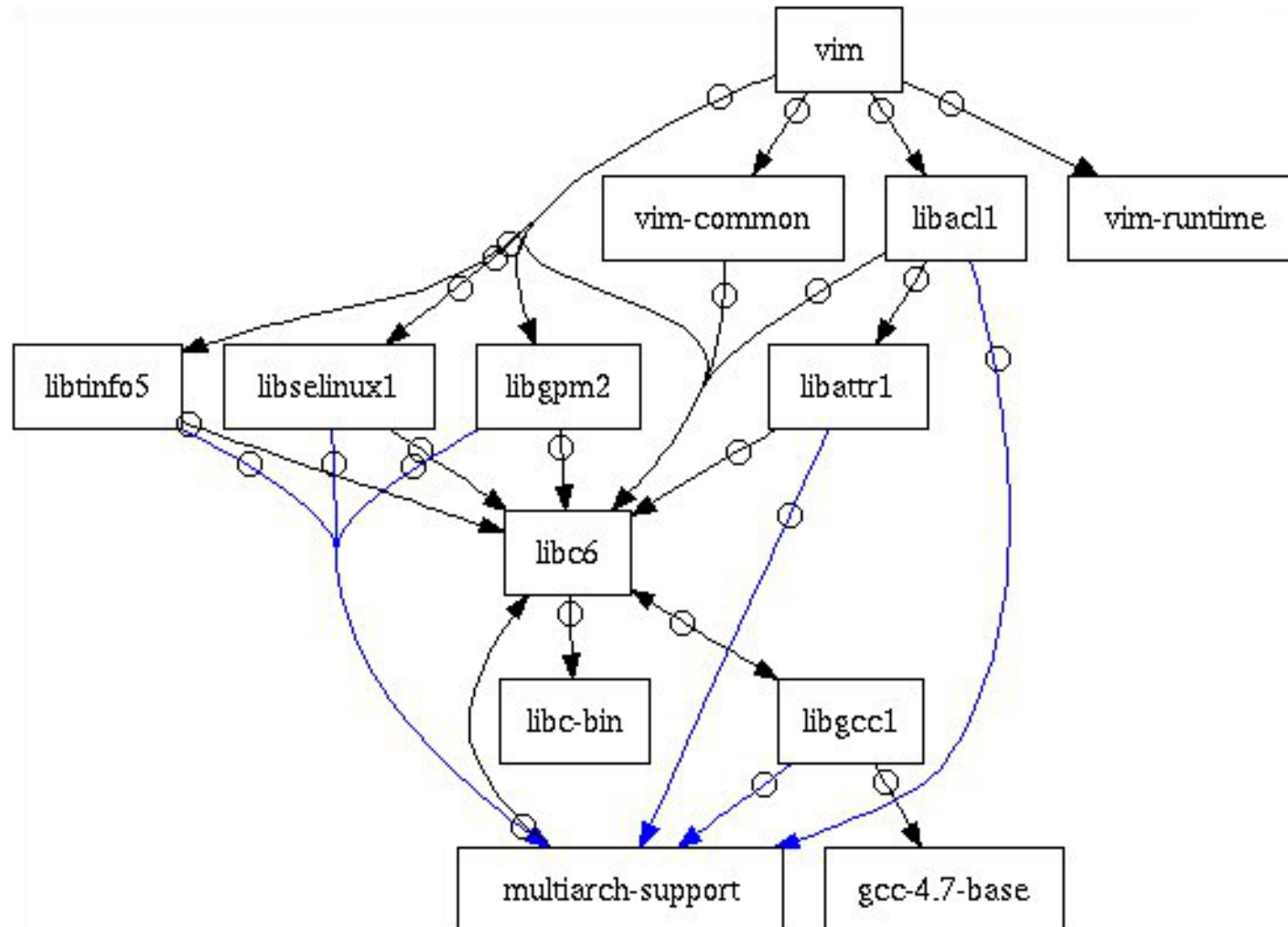


Virtual memory



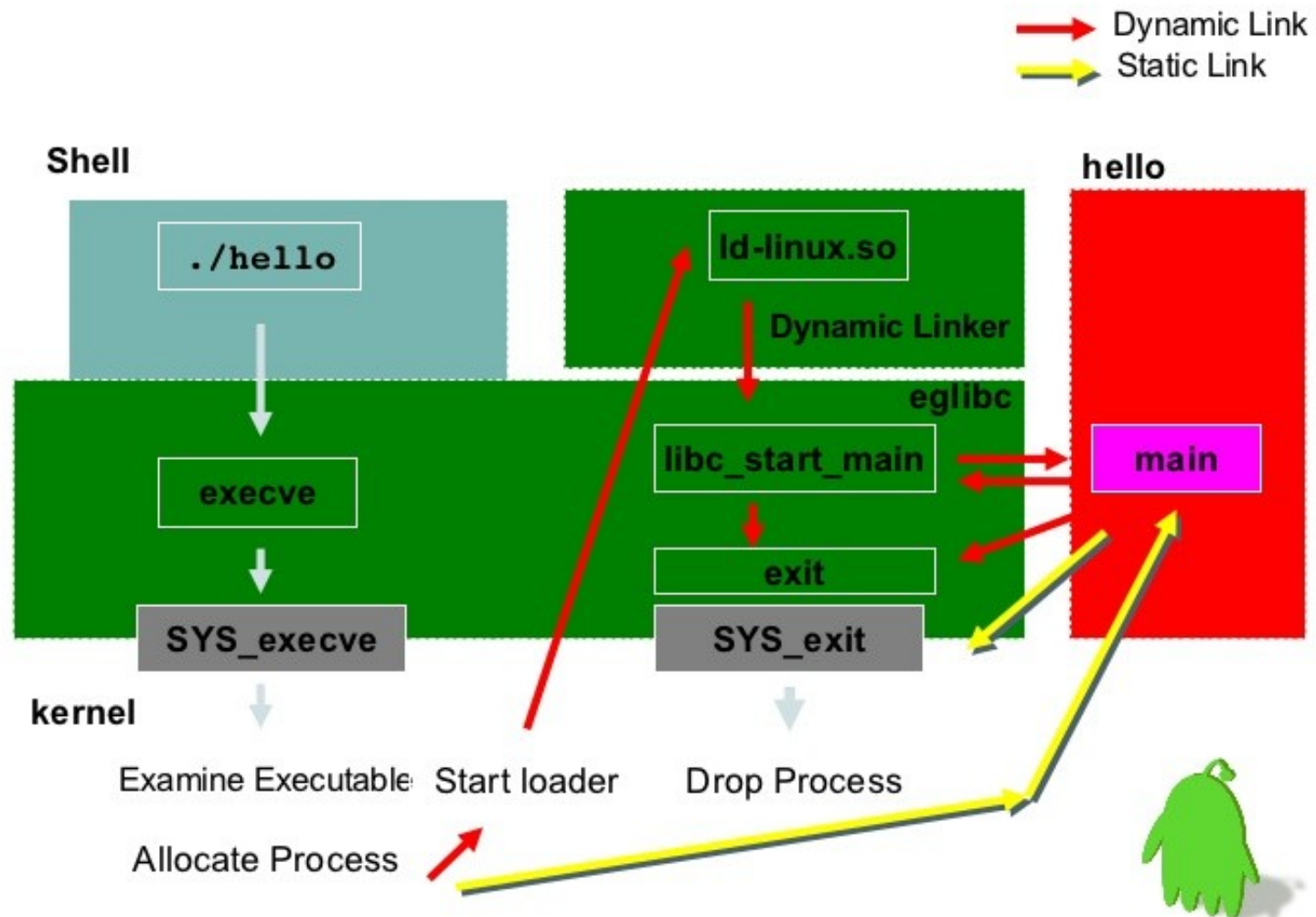


Application

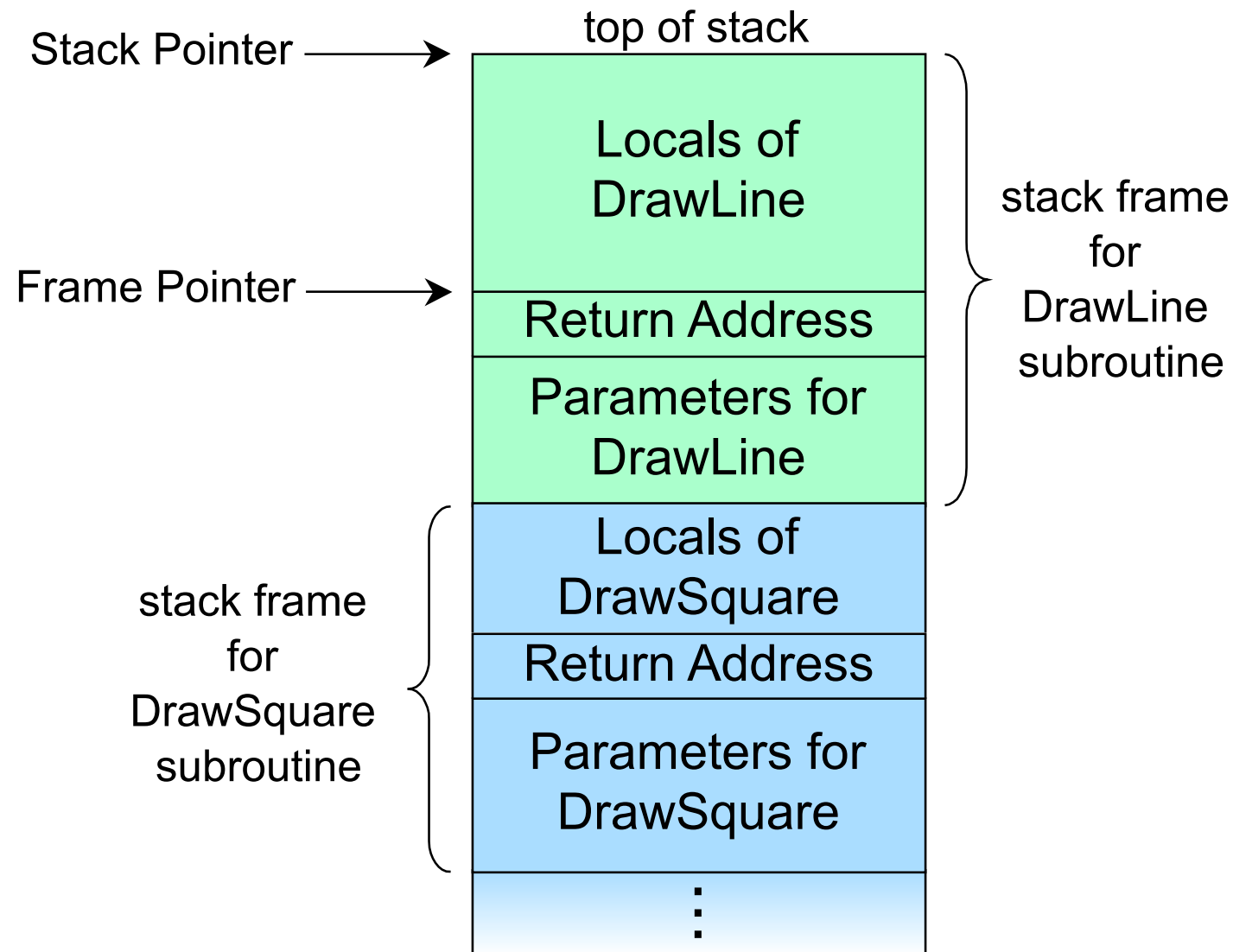


Application

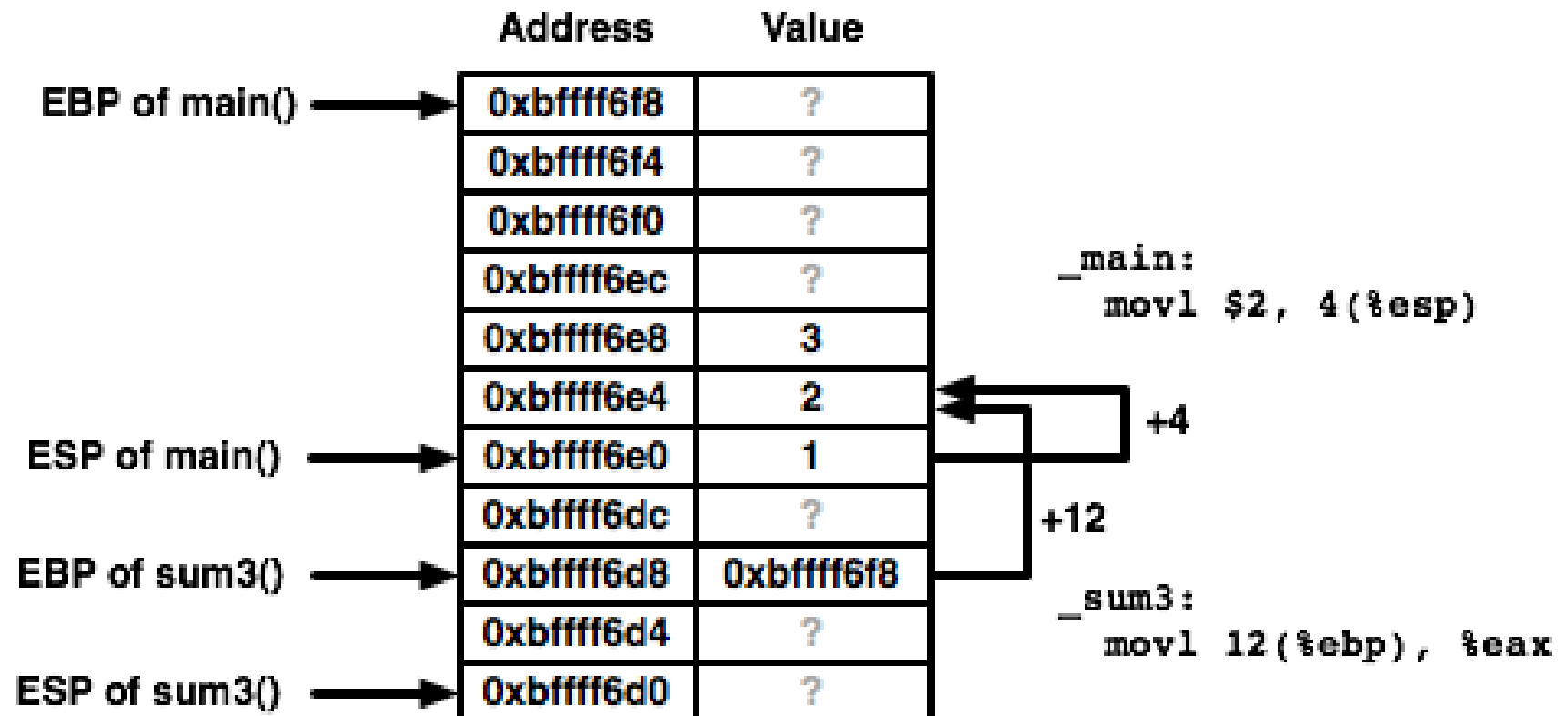
Execution flow of Hello World! (GNU/Linux)

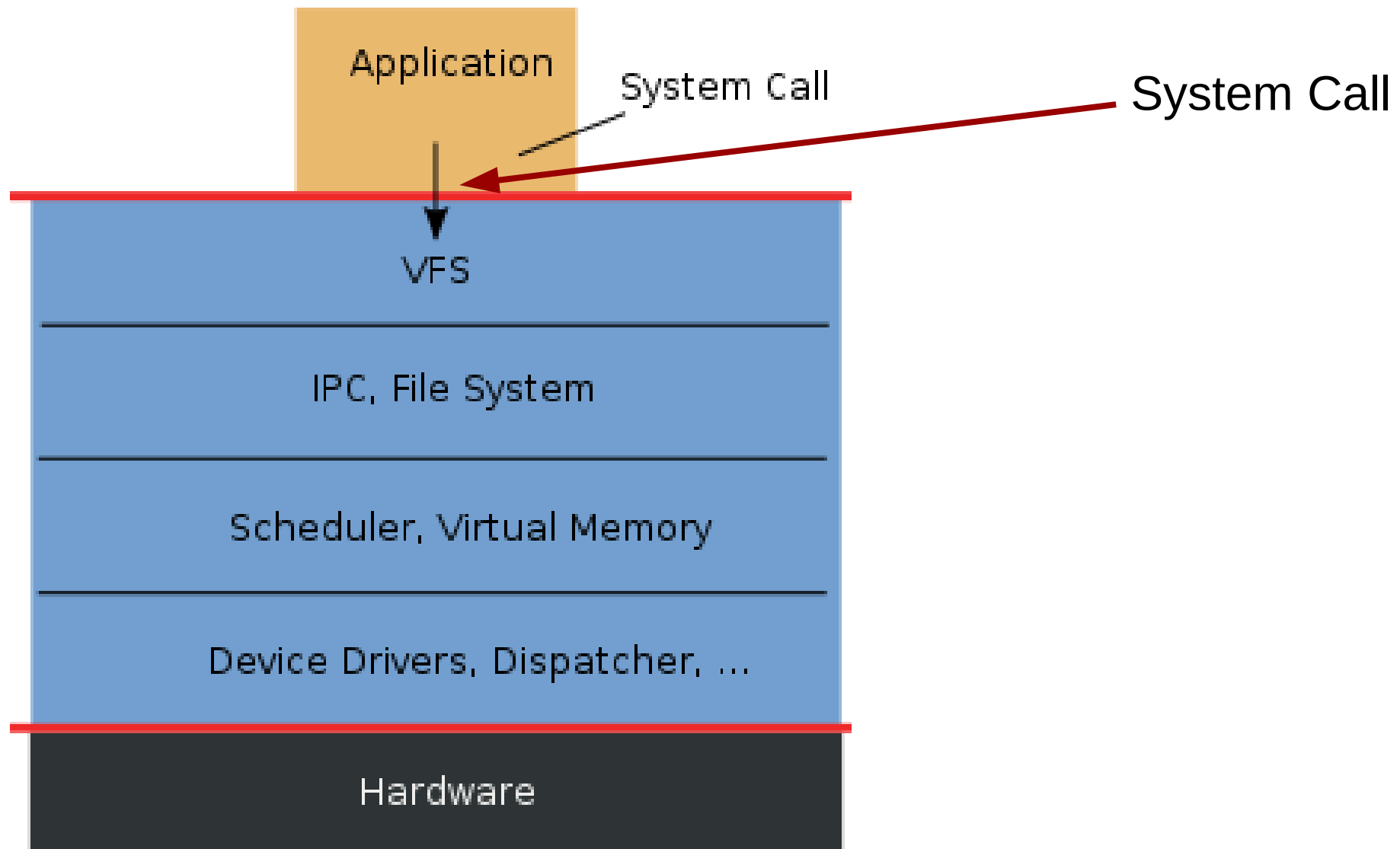


Application

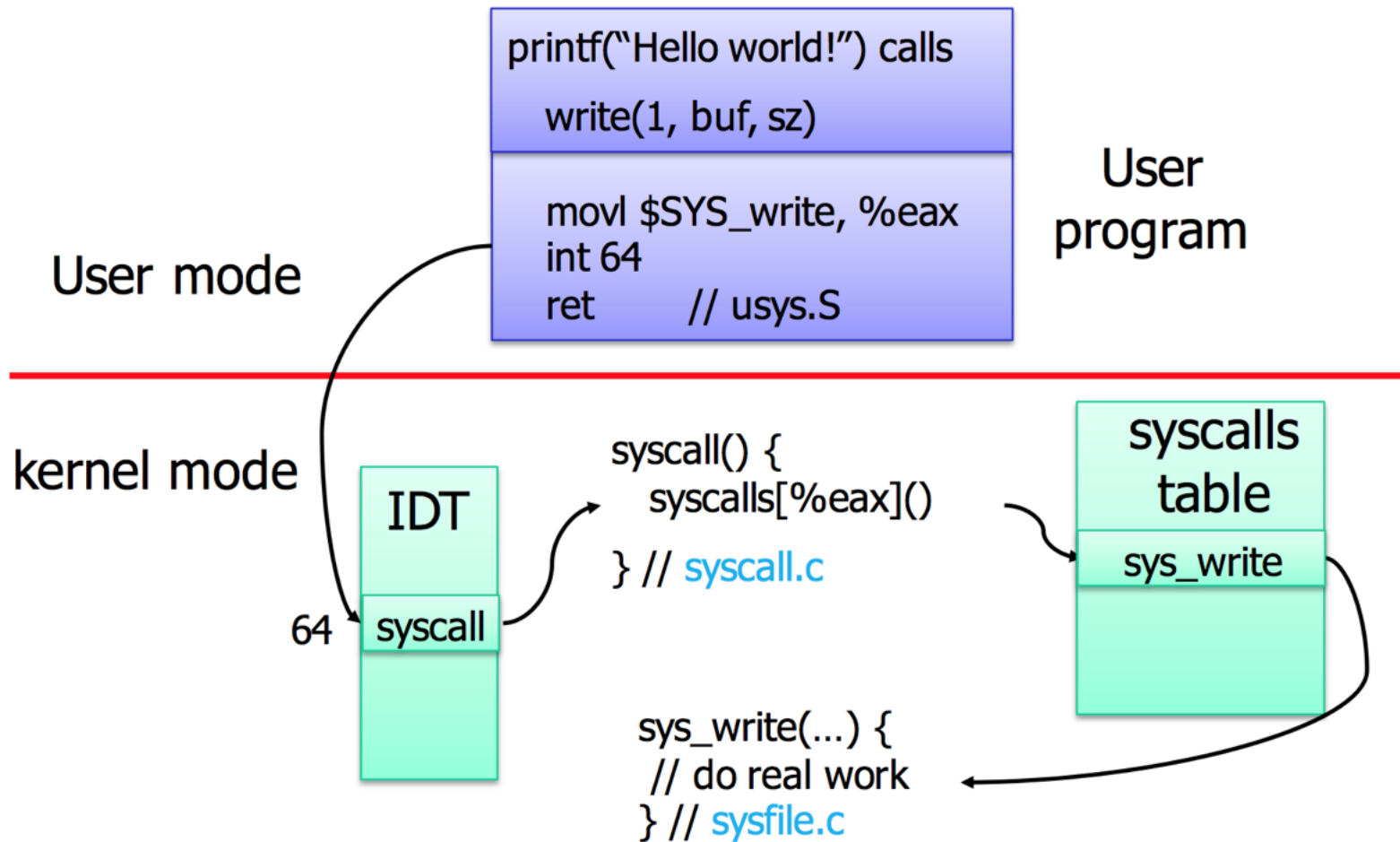


Application

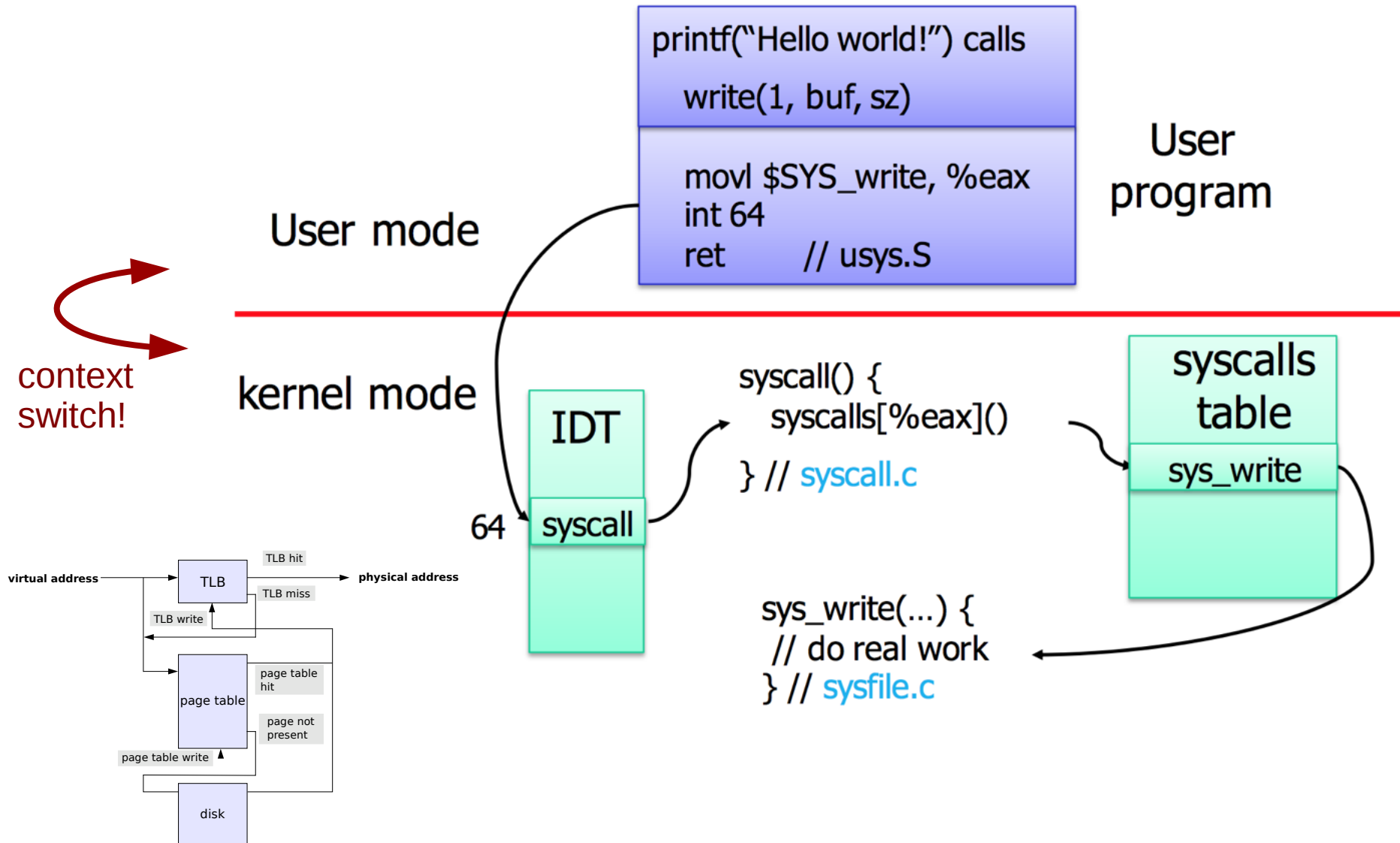


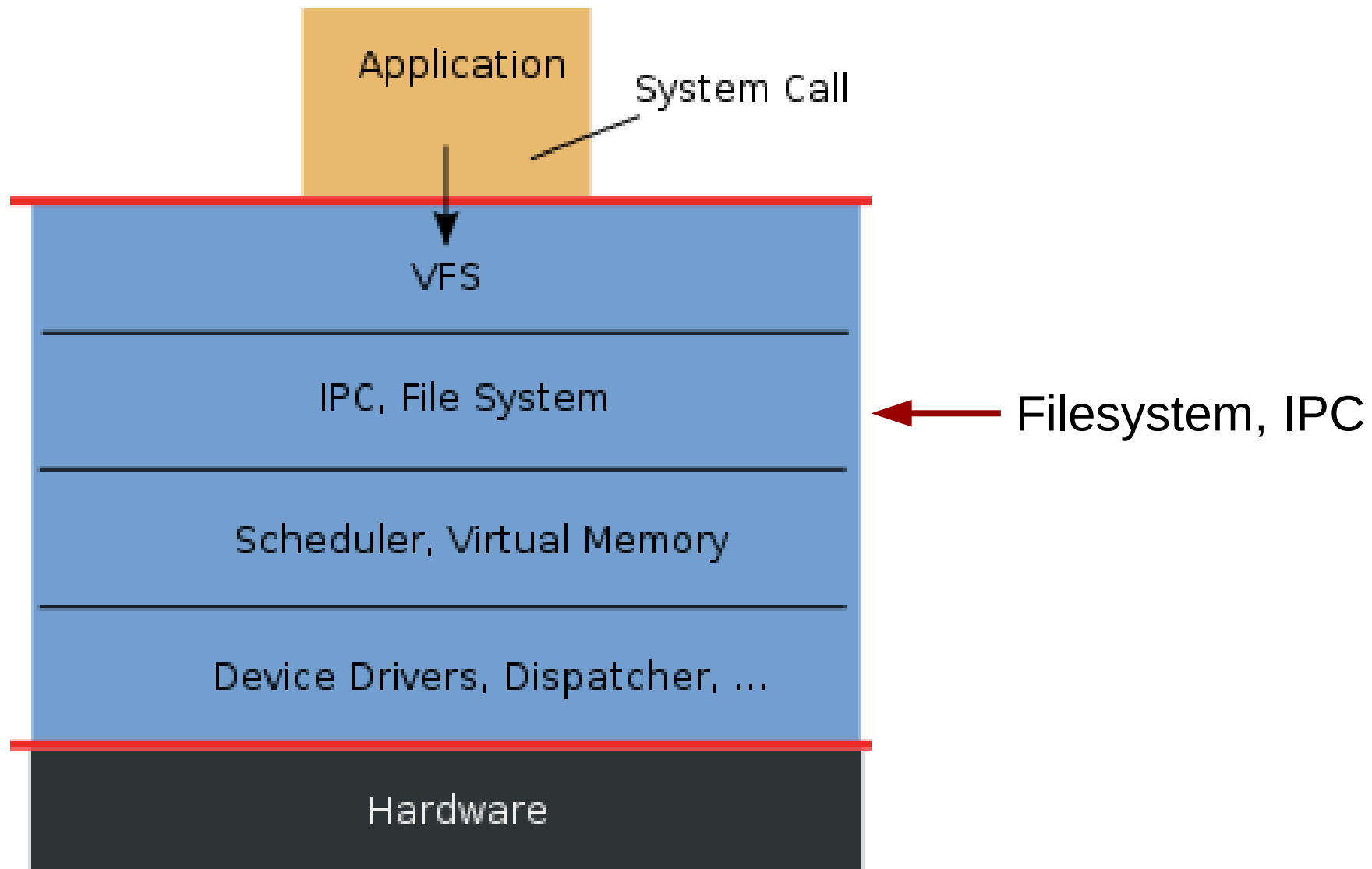


System call

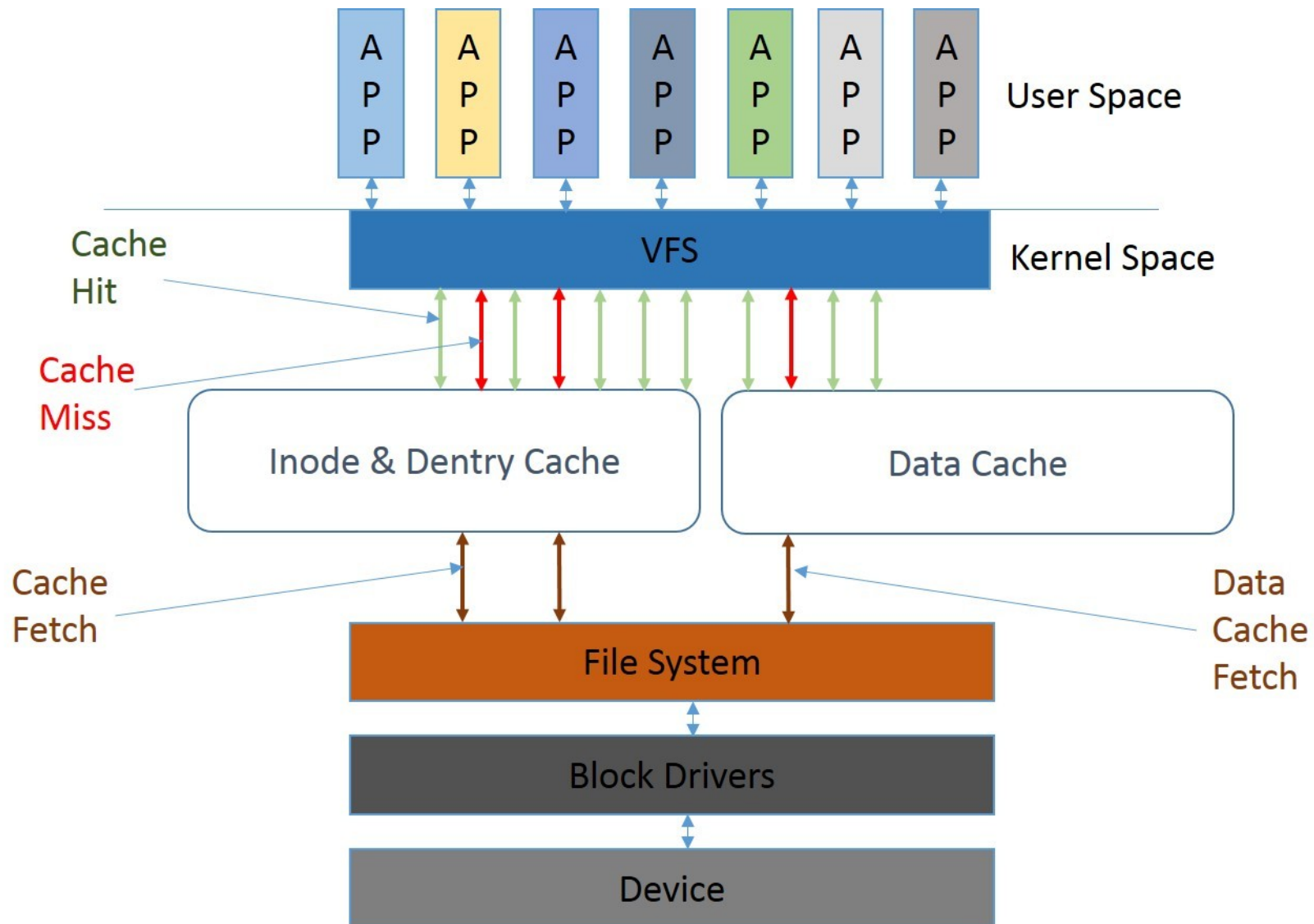


System call

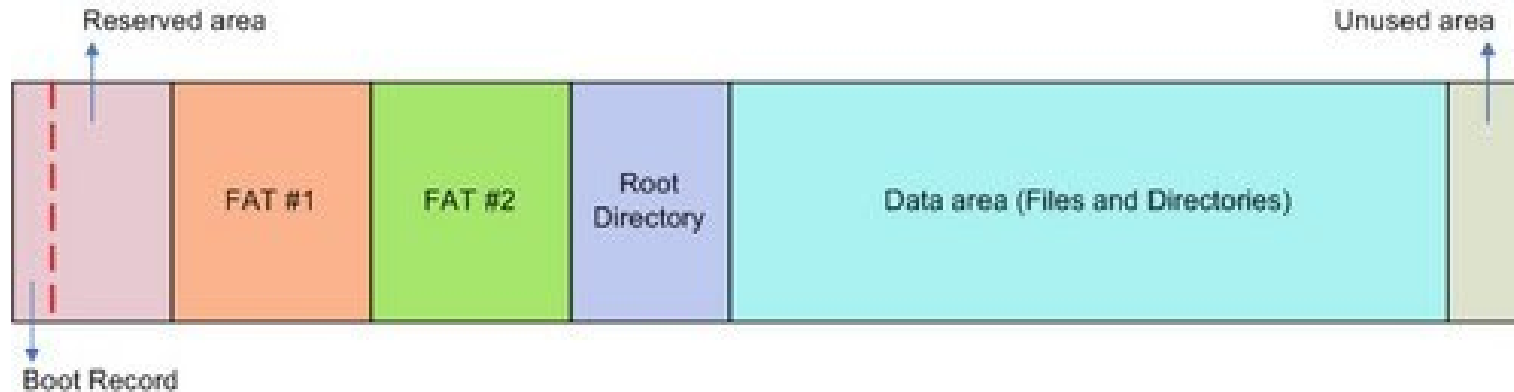




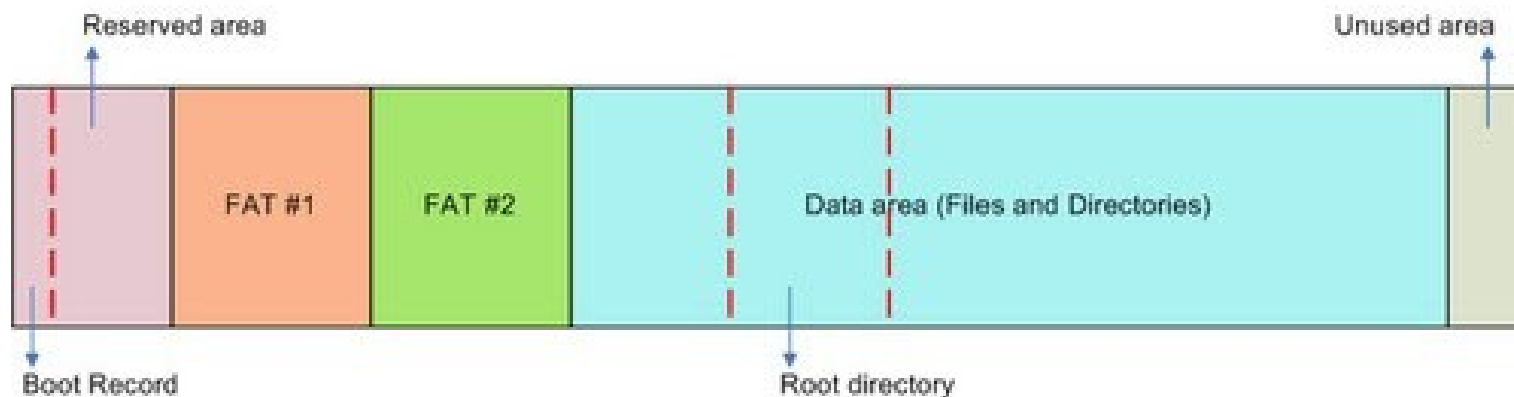
Filesystem: VFS



Filesystem (FAT32)



The structure of FAT16 file system



The structure of FAT32 file system

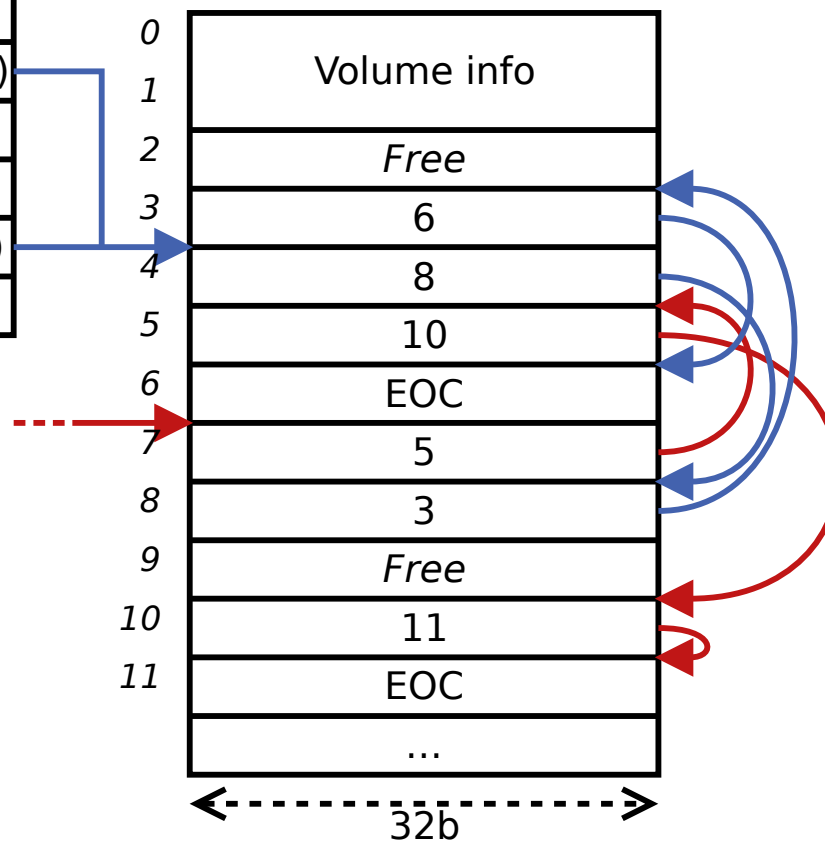
by *lprinceps*, <lprinceps@gmail.com>

Filesystem (FAT32)

Directory table entry (32B)

Filename (8B)
Extension (3B)
Attributes (1B)
Reserved (1B)
Create time (3B)
Create date (2B)
Last access date (2B)
First cluster # (MSB, 2B)
Last mod. time (2B)
Last mod. date (2B)
First cluster # (LSB, 2B)
File size (4B)

File allocation table



Filesystem (FAT32)

XXXXXXXX	XXXXXXXX	00000009	00000004
00000005	00000007	00000000	00000008
FFFFFFFF	0000000A	0000000B	00000011
0000000D	0000000E	FFFFFFFF	00000010
00000012	FFFFFFFF	00000013	00000014
00000015	00000016	FFFFFFFF	00000000
00000000	00000000	00000000	00000000
00000000	00000000	00000000	00000000
00000000	00000000	00000000	00000000
00000000	00000000	00000000	00000000
00000000	00000000	00000000	00000000
00000000	00000000	00000000	00000000
00000000	00000000	00000000	00000000
00000000	00000000	00000000	00000000
00000000	00000000	00000000	00000000
00000000	00000000	00000000	00000000

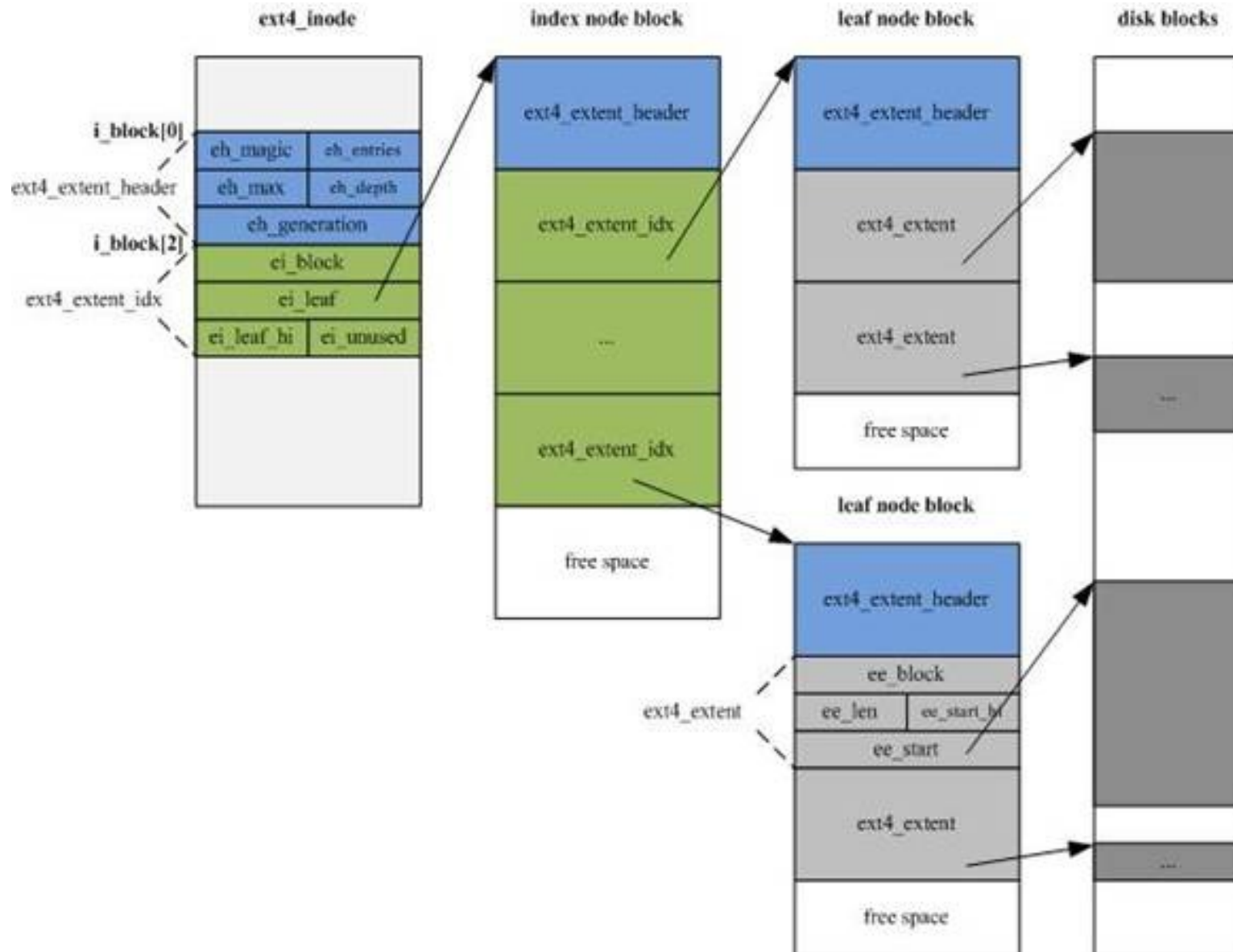
Root Directory:
2, 9, A, B, 11

File #1:
3, 4, 5, 7, 8

File #2:
C, D, E

File #3:
F, 10, 12, 13, 14, 15, 16

Filesystem (ext4)



Filesystem (ext4)

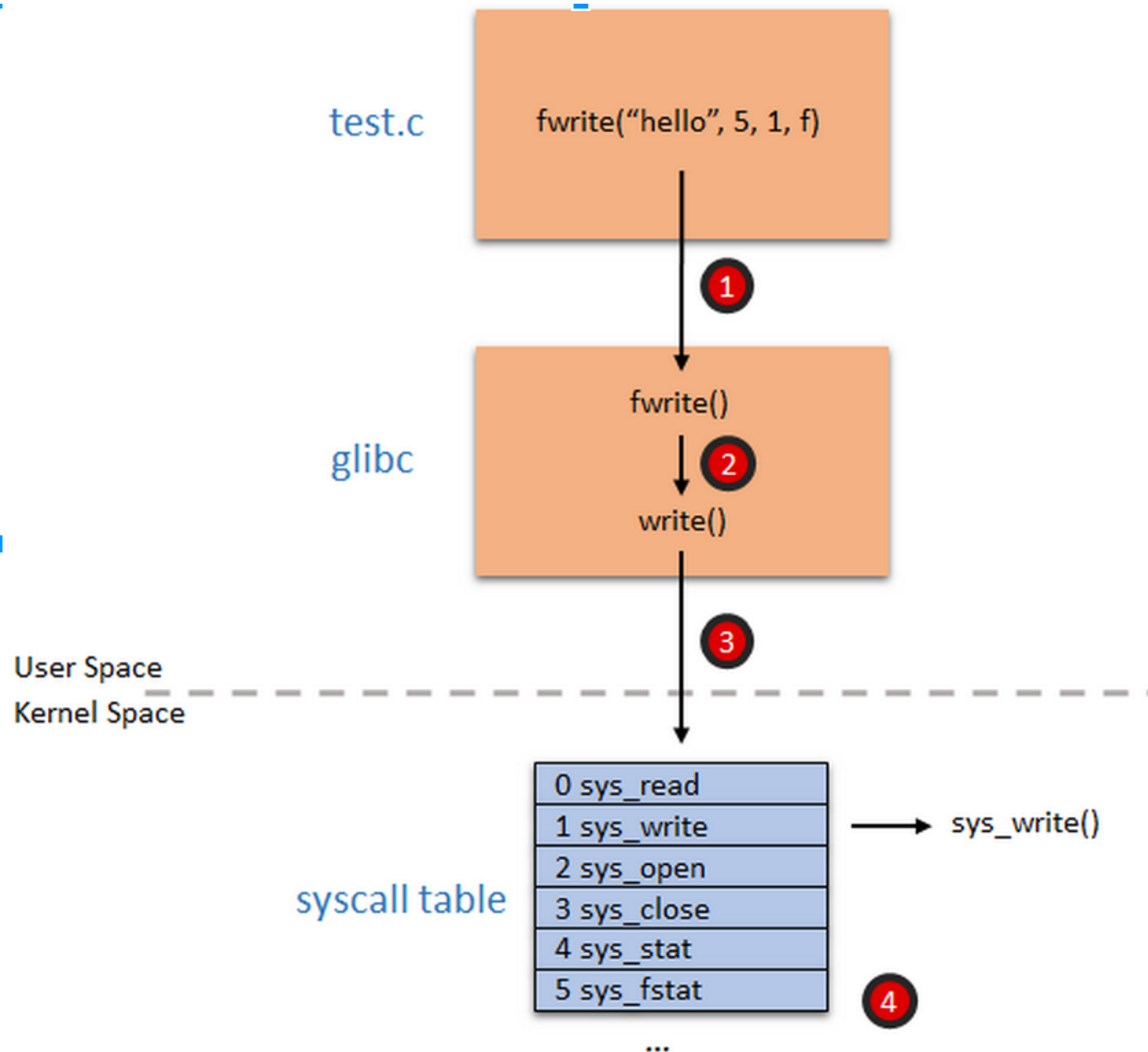
TxB	Journal		TxE	File System	
	Contents (metadata)	(data)		Metadata	Data
issue complete	issue complete	issue complete			
-----			issue complete		
-----				issue complete	issue complete

Table 42.1: Data Journaling Timeline

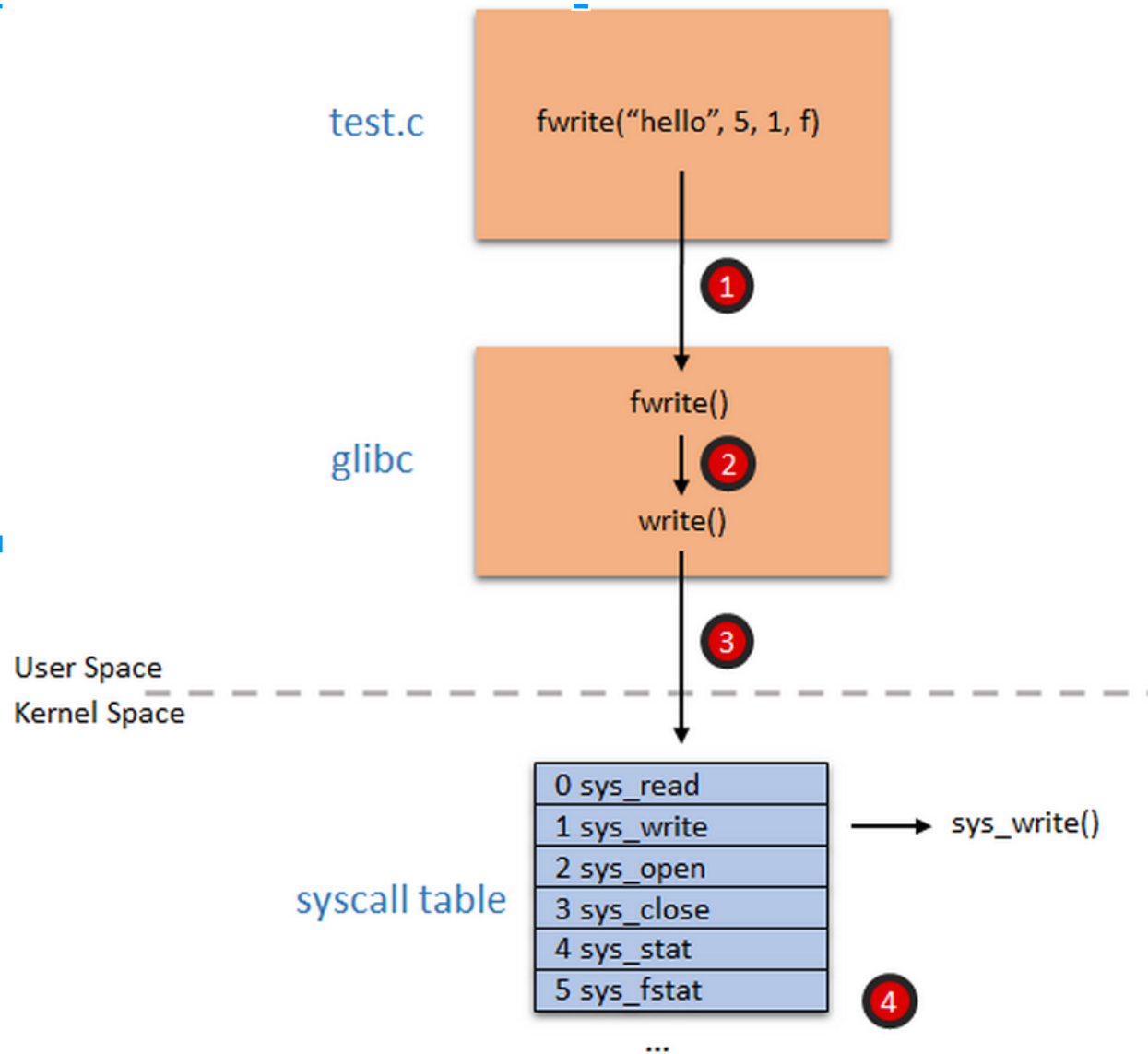
Filesystem

- Copy on Write
- Transparent compression
- SSD-cache / TRIM
- Network filesystems
- Distributed filesystems
- Pseudo file systems
(wikipediaFS, PiFS, ...)
- ...

Filesystem

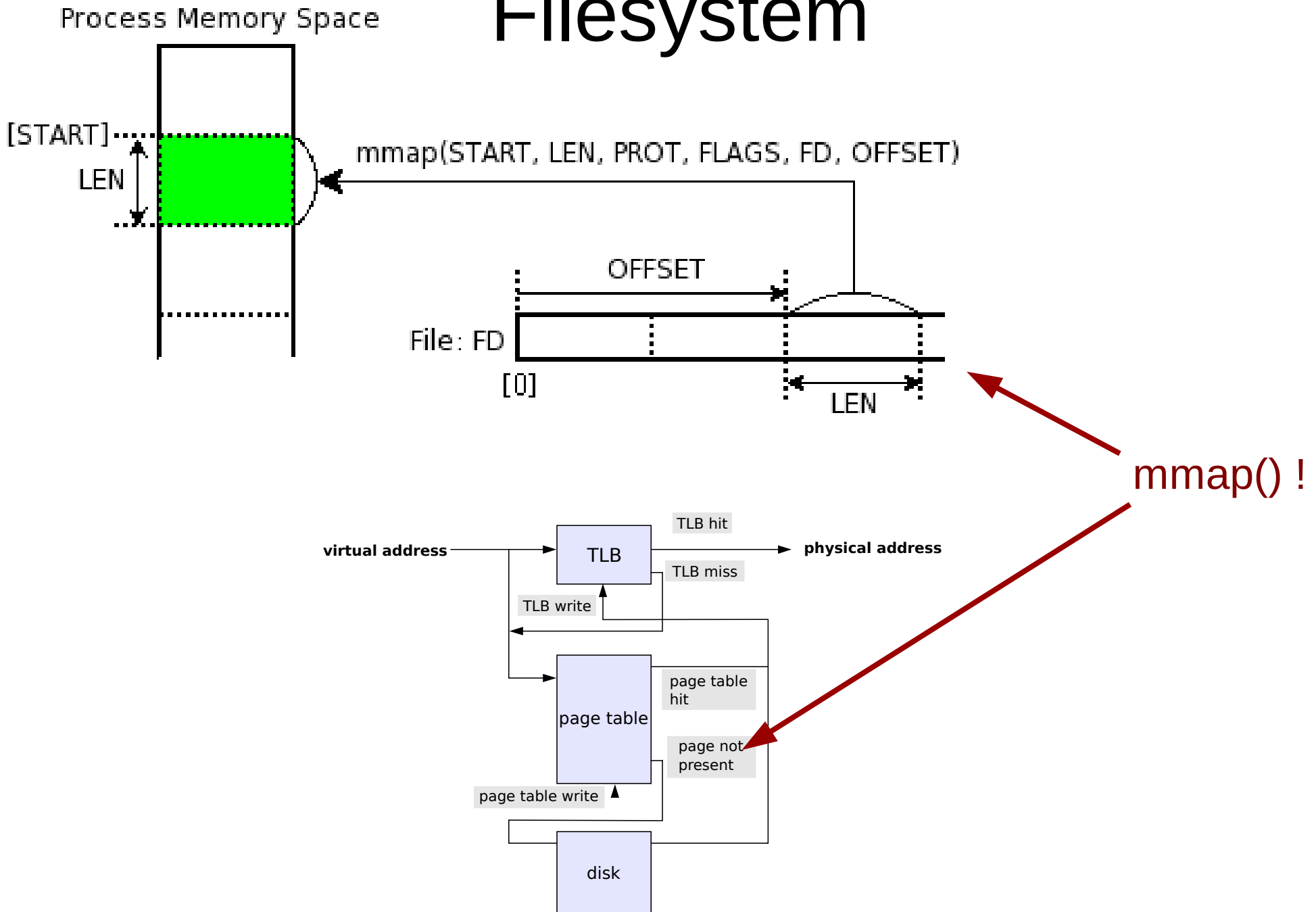


Filesystem



`mmap()` !

Filesystem



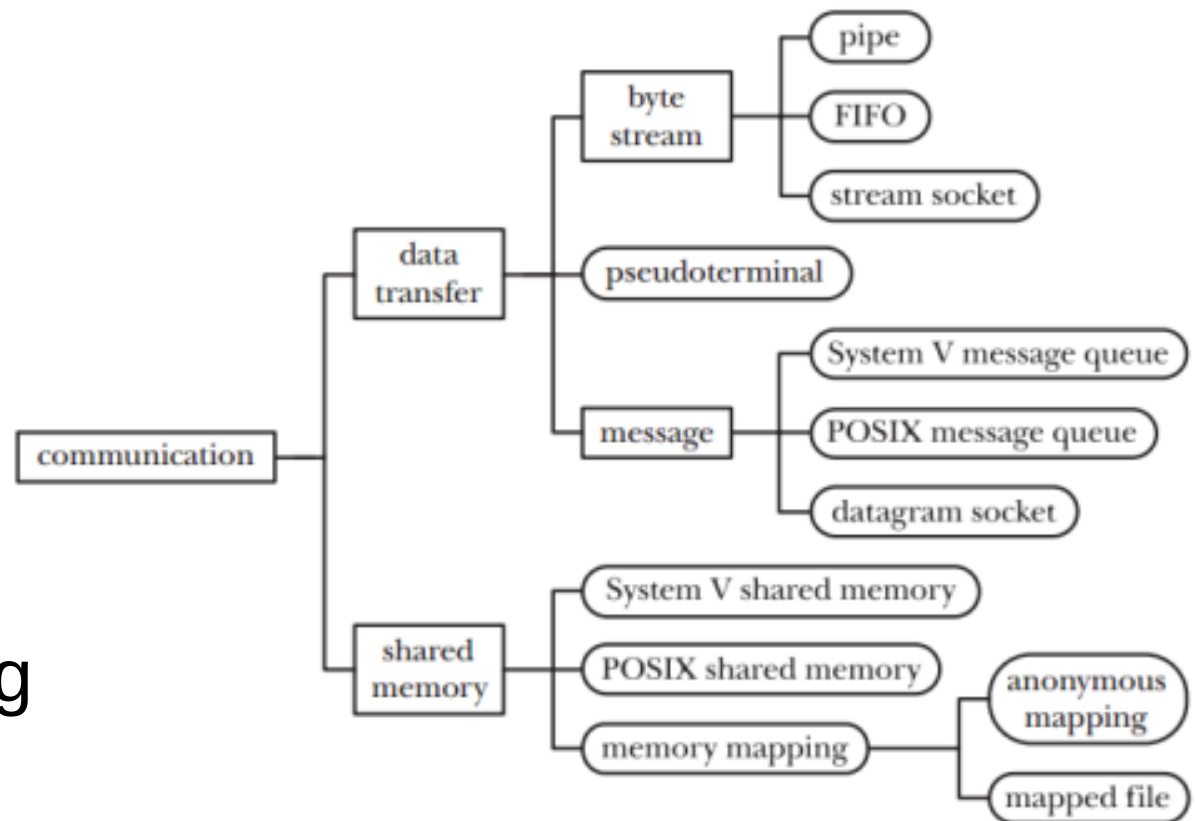
IPC

- IPC == Inter-Process Communication:
 - Files
 - Signals, Traps
 - Sockets
 - Message queue
 - Pipe
 - Shared memory
 - Message passing
 - ...

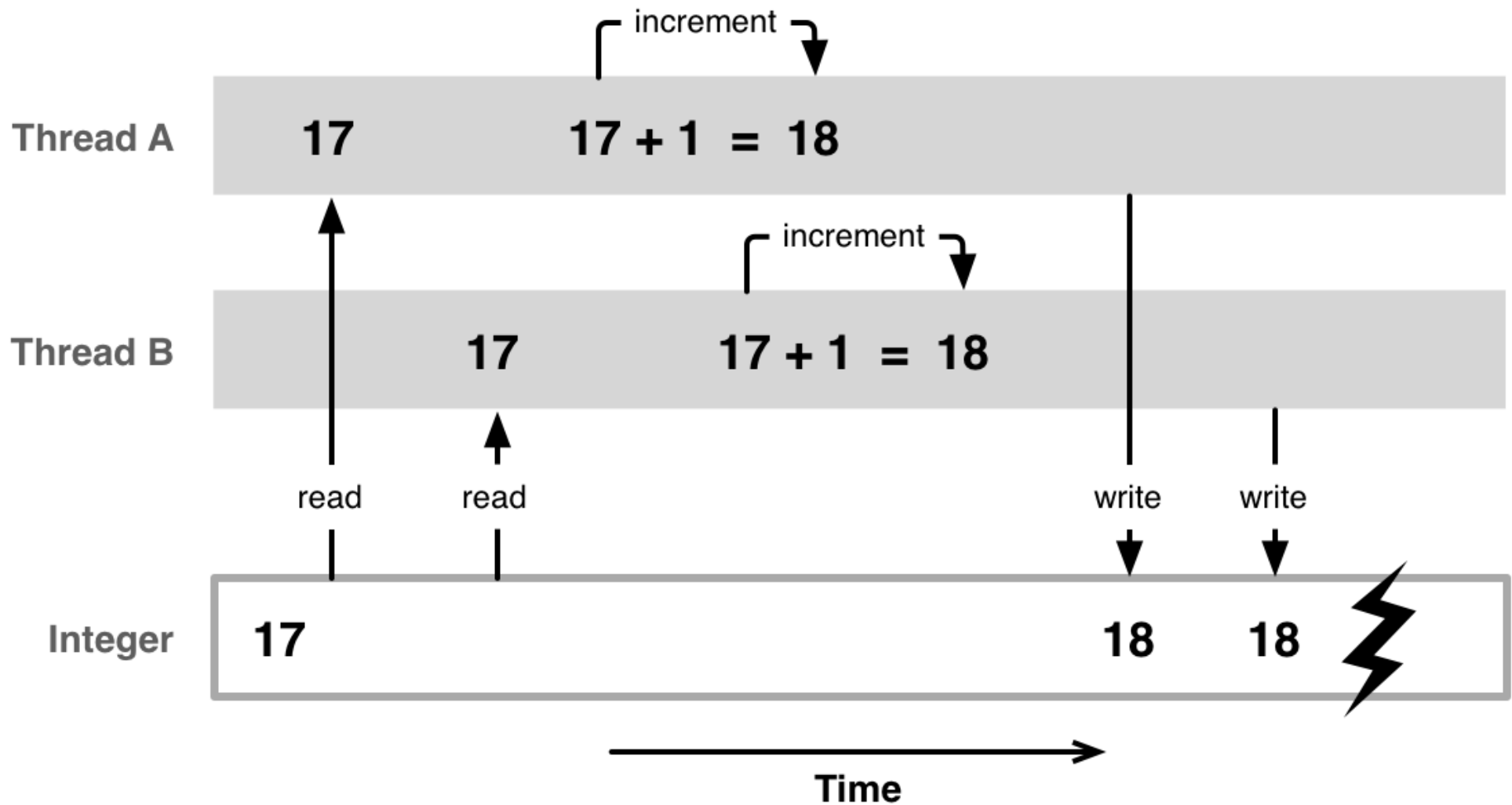
IPC

- IPC == Inter-Process Communication:

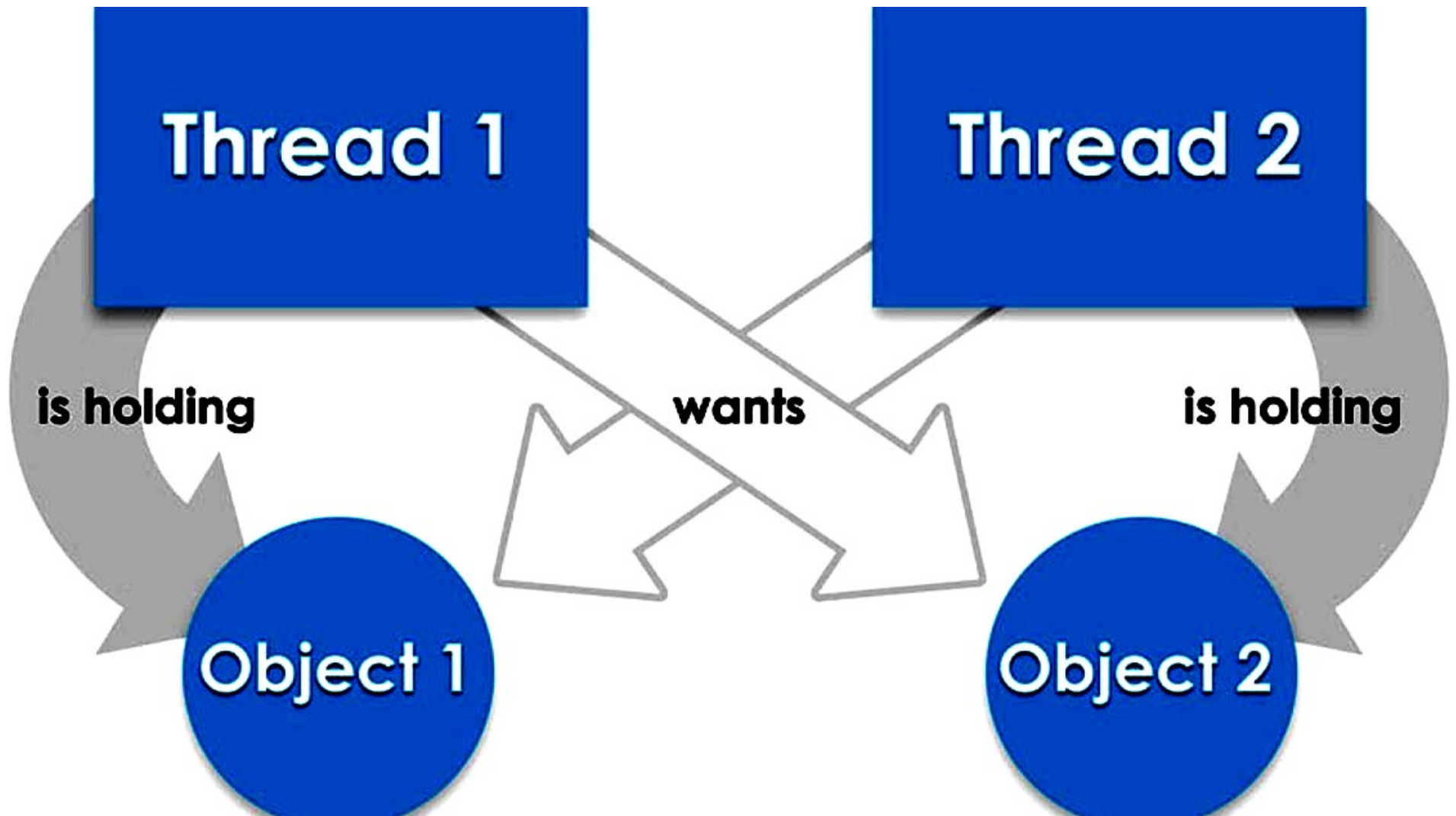
- Files
- Signals, Traps
- Sockets
- Message queue
- Pipe
- Shared memory
- Message passing
- ...



IPC: race condition



IPC: deadlock

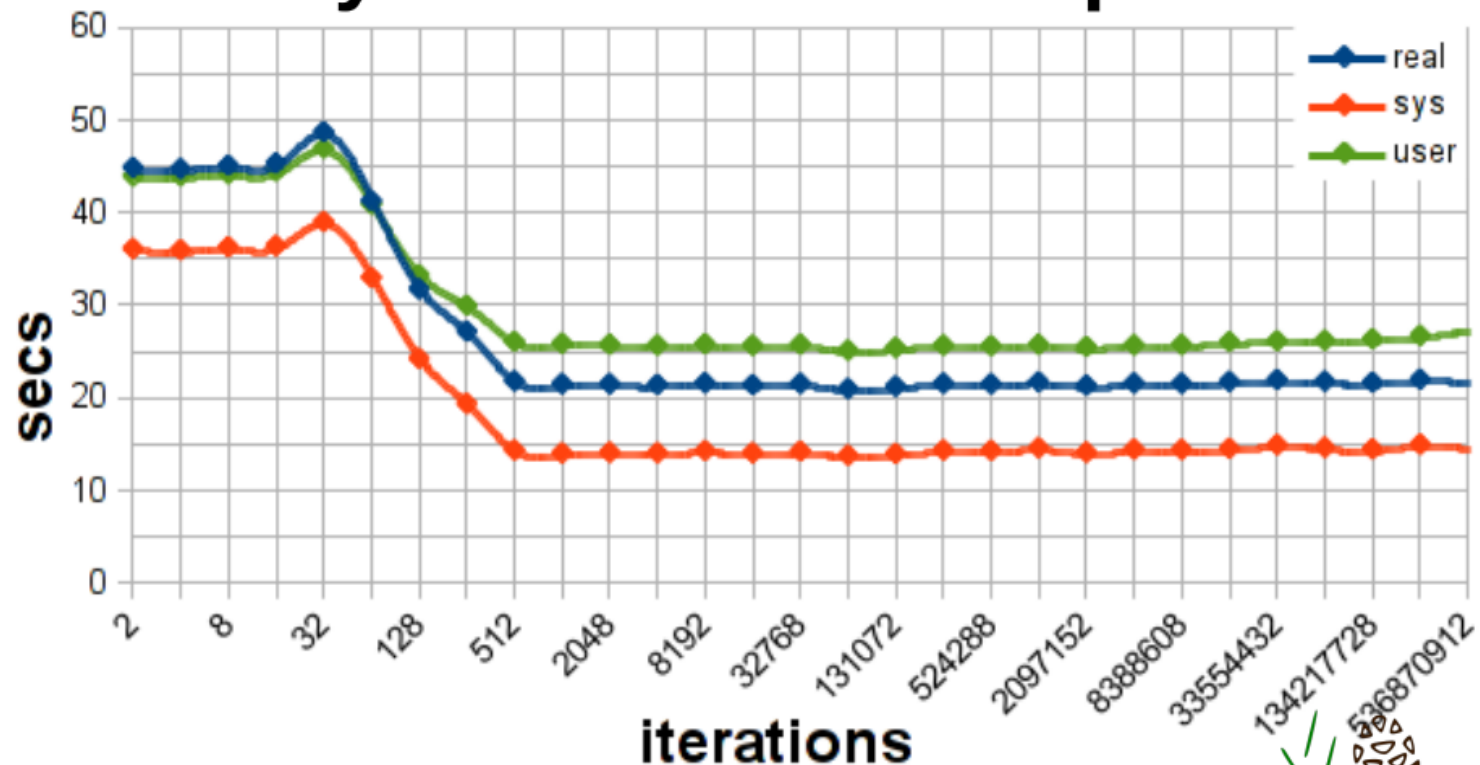


IPC: synchronization

Technique	Description	Scope
Per-CPU variables	Duplicate a data structure among CPUs	All CPUs
Atomic operation	Atomic read-modify-write instruction	All
Memory barrier	Avoid instruction re-ordering	Local CPU
Spin lock	Lock with busy wait	All
Semaphore	Lock with blocking wait (sleep)	All
Seqlocks	Lock based on access counter	All
Local interrupt disabling	Forbid interrupt on a single CPU	Local
Local softirq disabling	Forbid deferrable function on a single CPU	Local
Read-copy-update (RCU)	Lock-free access to shared data through pointers	All

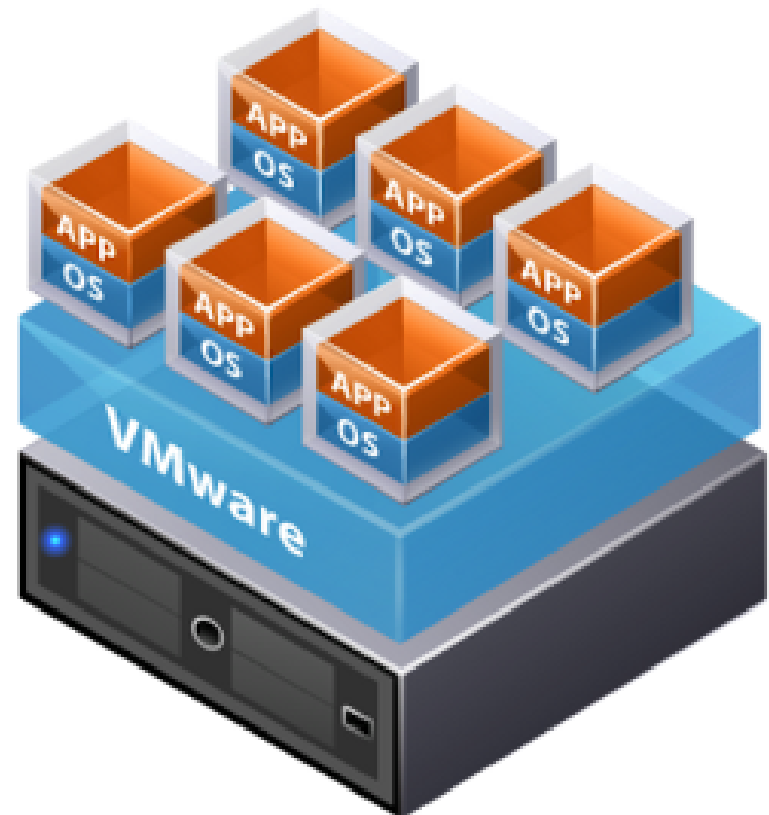
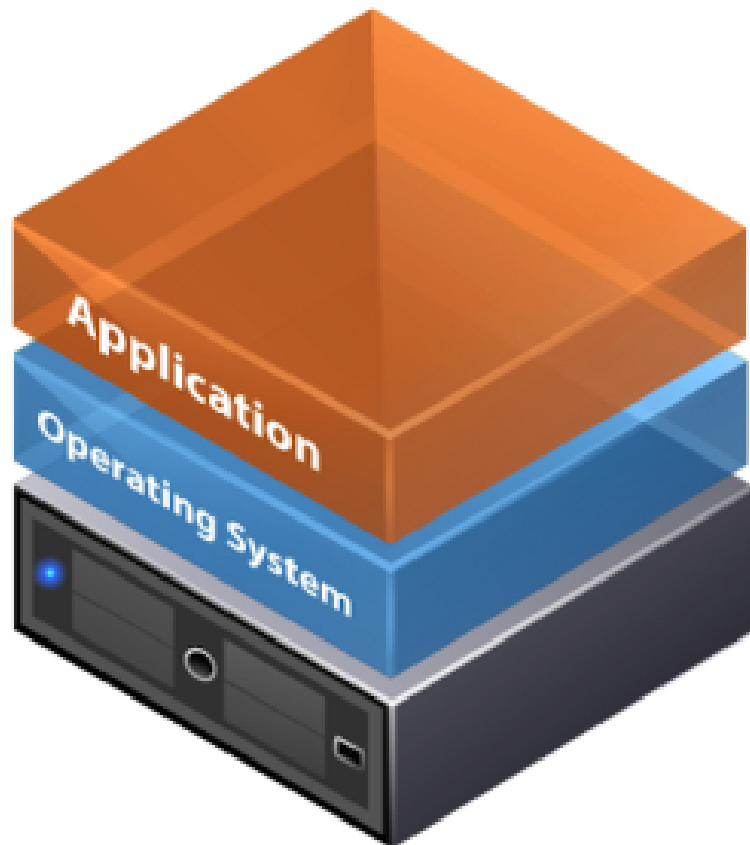
IPC: synchronization

**Зависимость времени выполнения
initial sync от timeout-a spinlock-a**



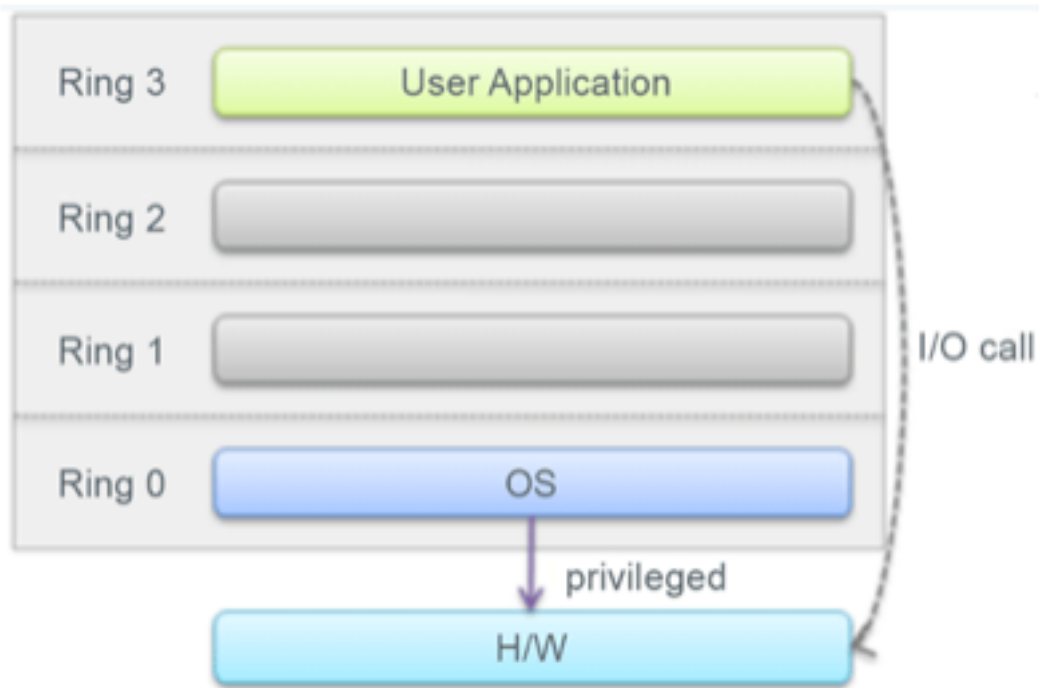
For dessert: virtualization, containers

vmx:

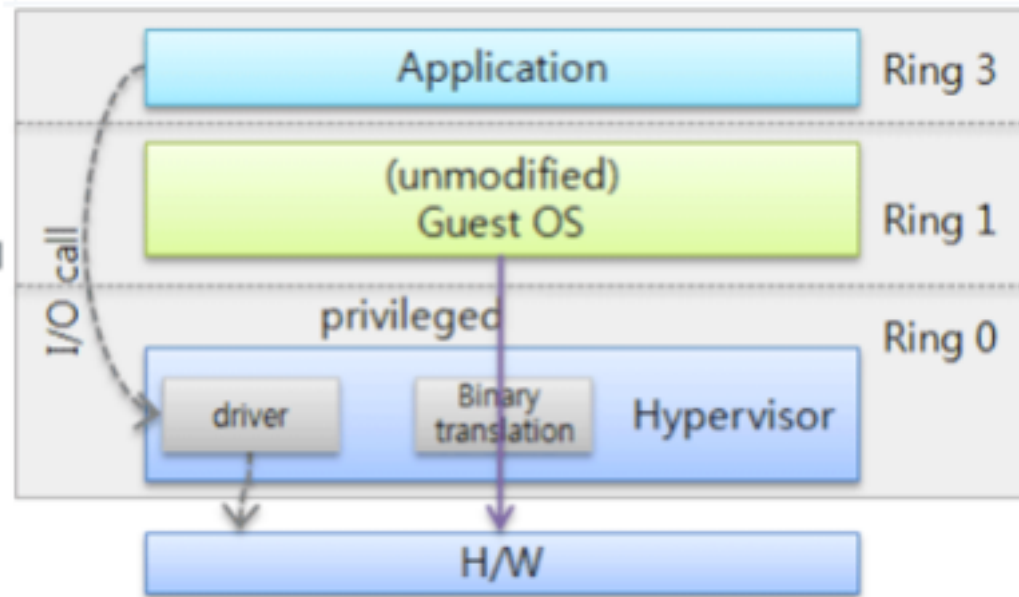


Virtualization

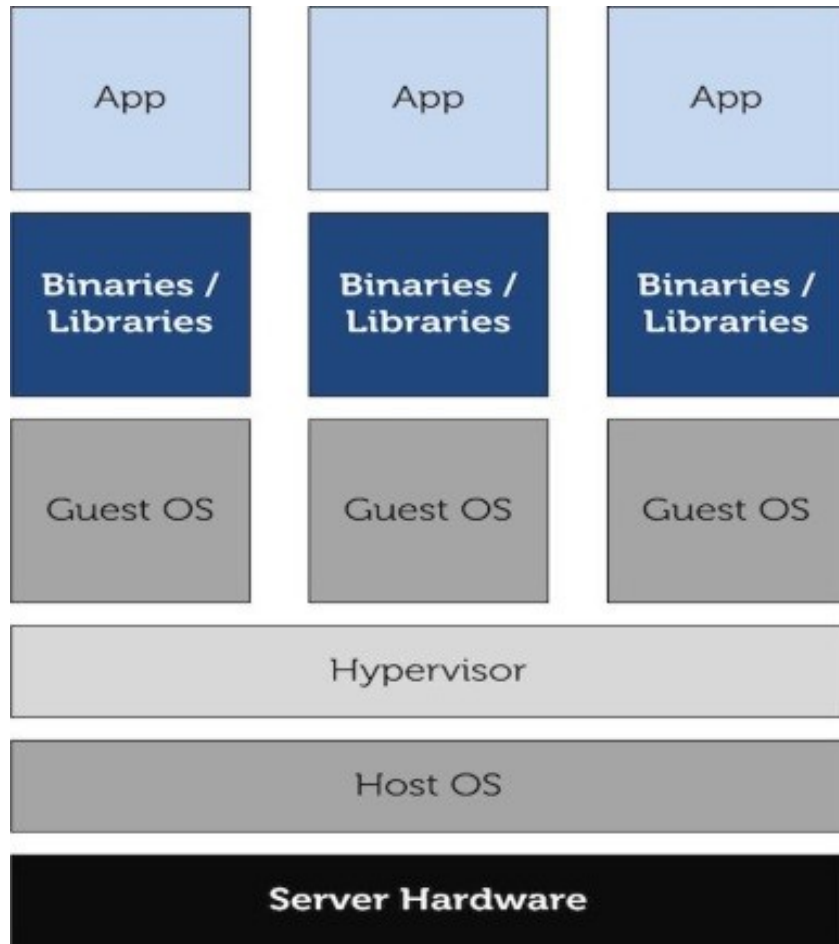
no virtualization



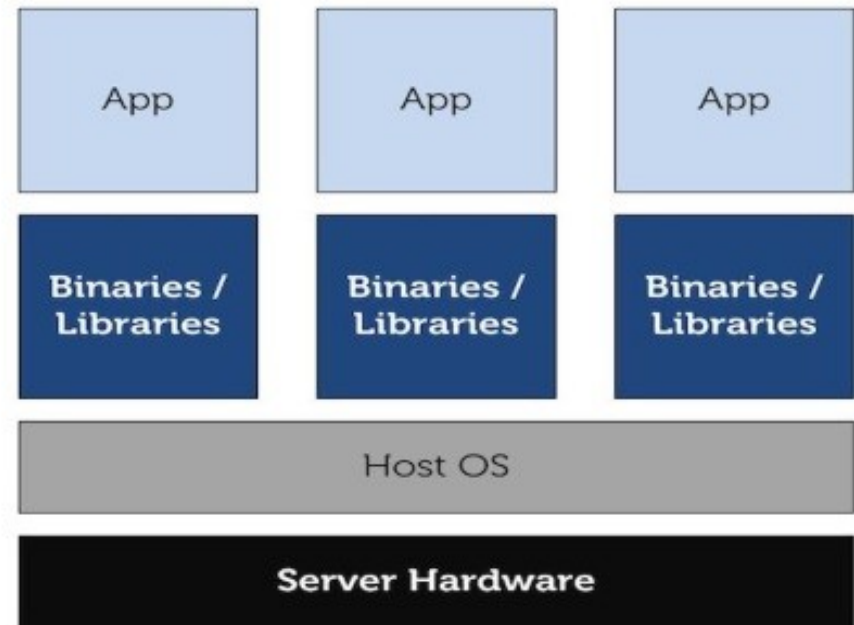
full virtualization



Containers



Virtualization



Containers

Q&A