# ICS143A: Principles of Operating Systems

# Final recap, sample questions

Anton Burtsev
March, 2017
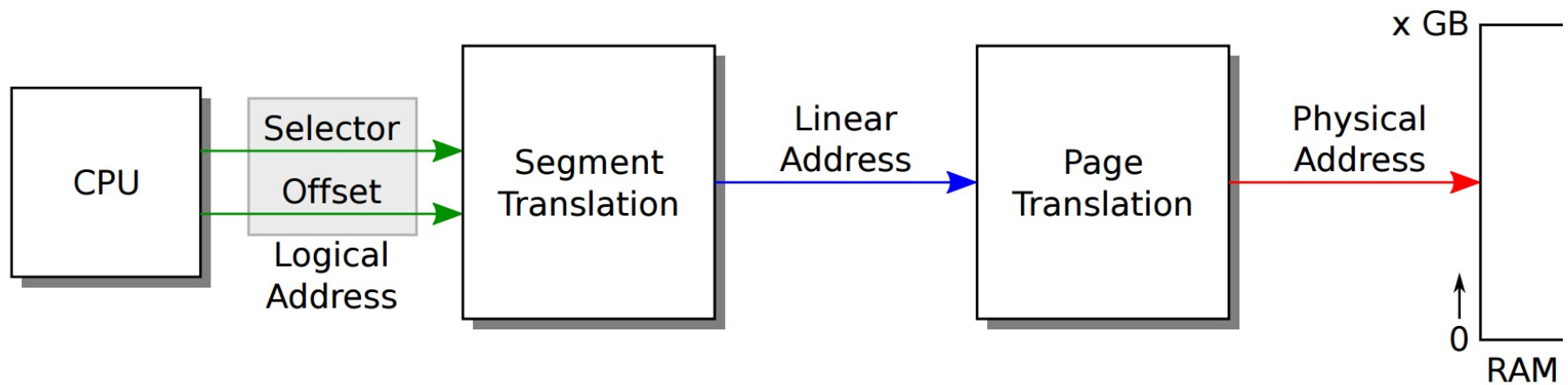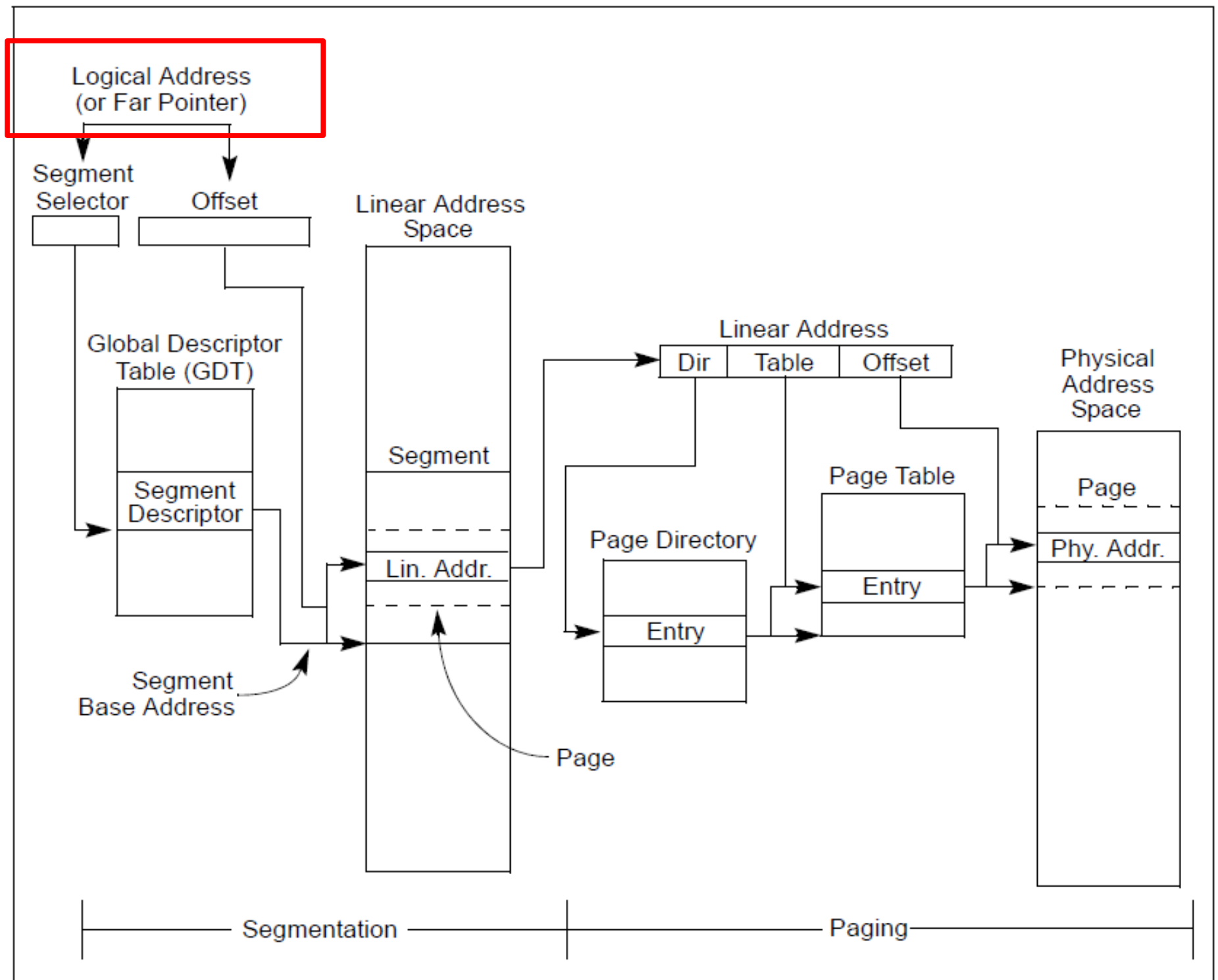
# What is operating system?
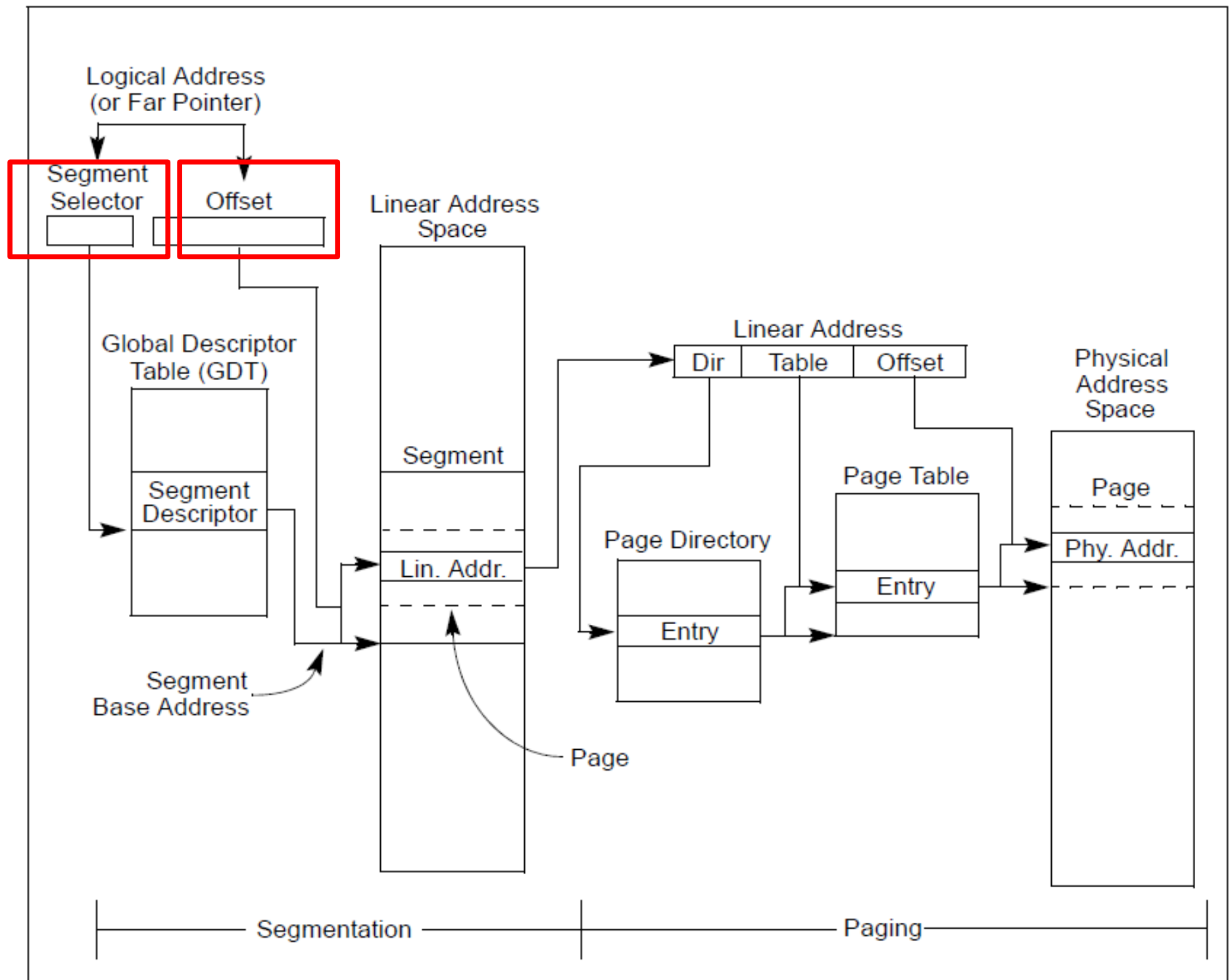
# What is a process?

# What mechanisms are involved into implementation of a process?

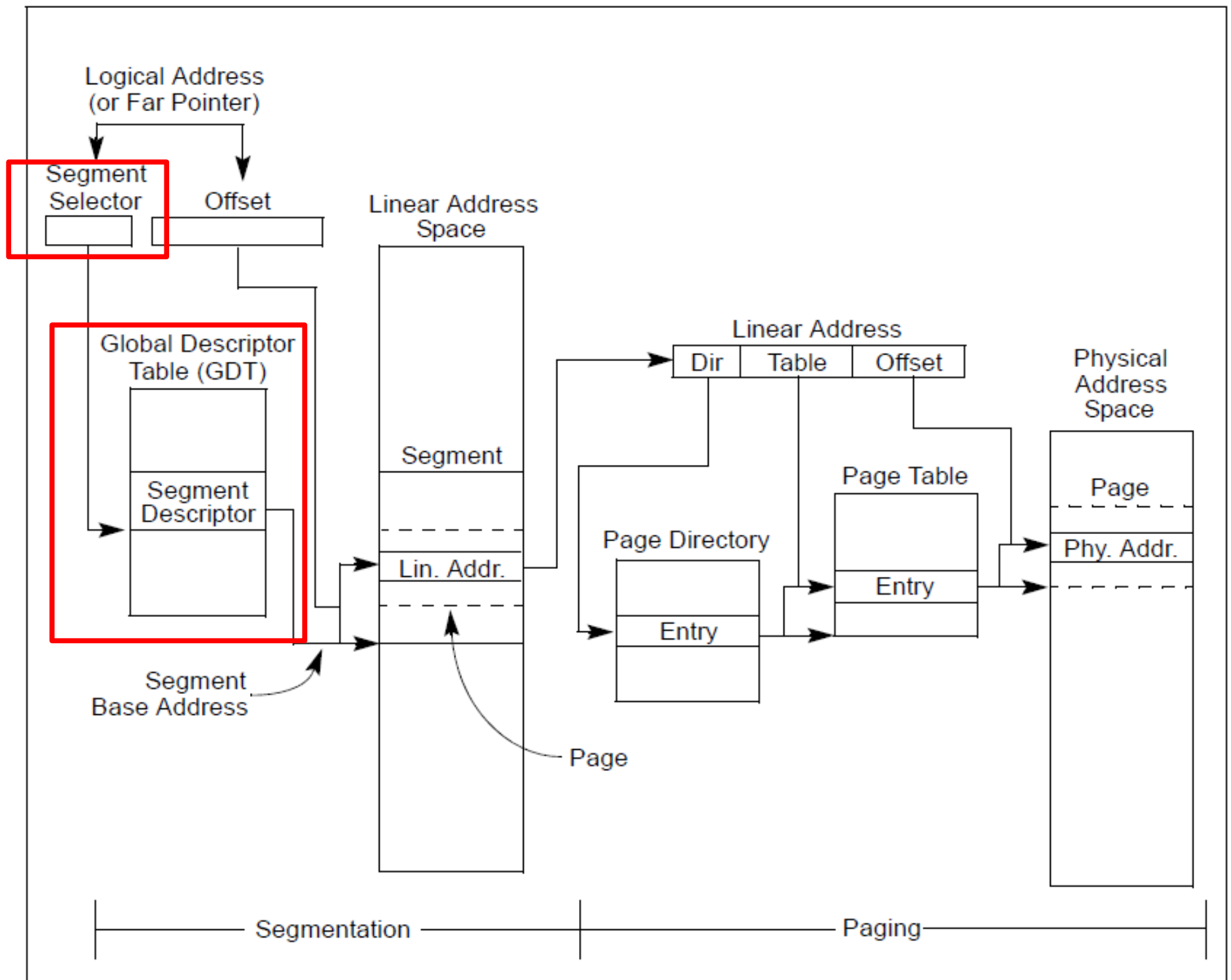Describe the x86 address translation pipeline (draw figure), explain stages.
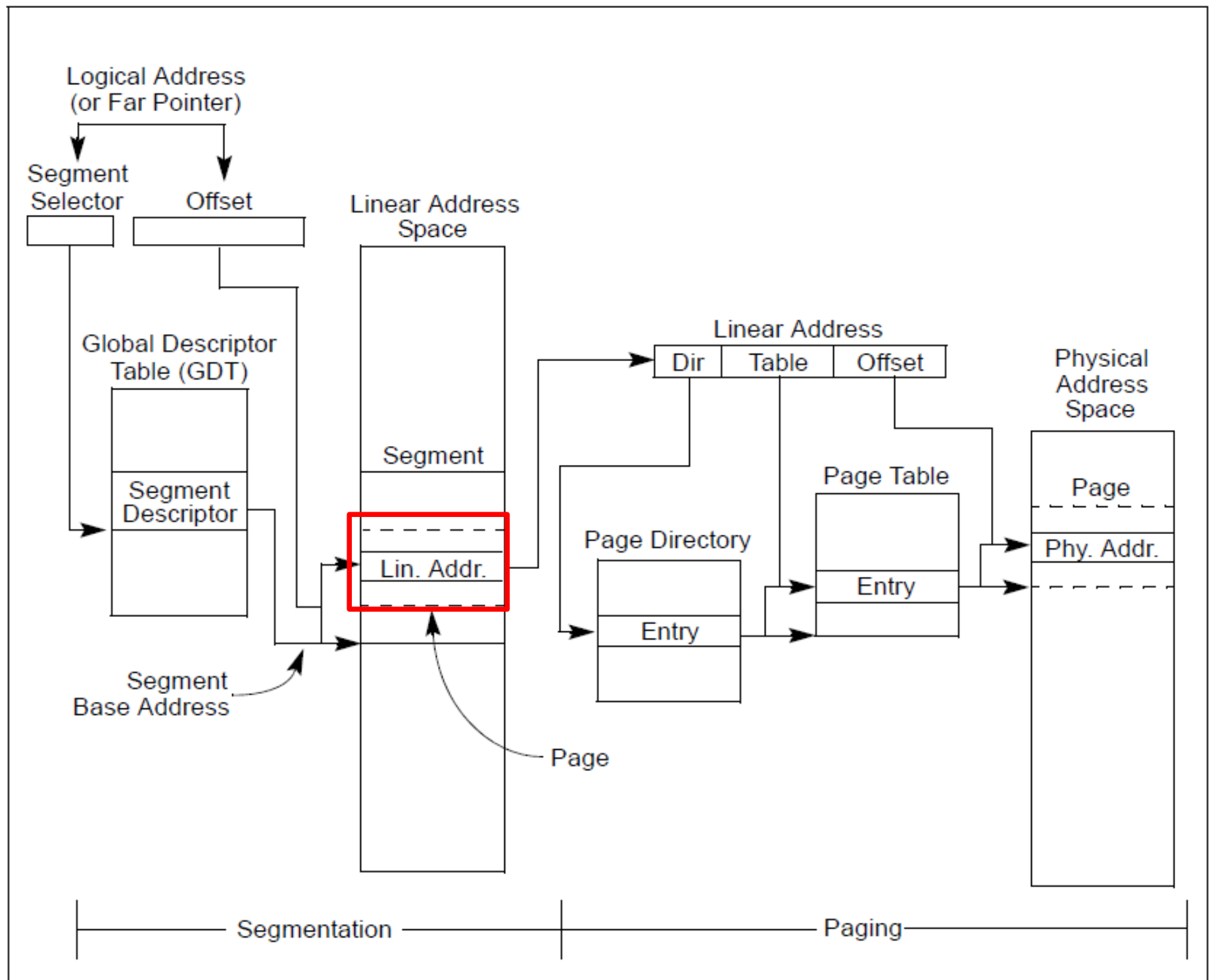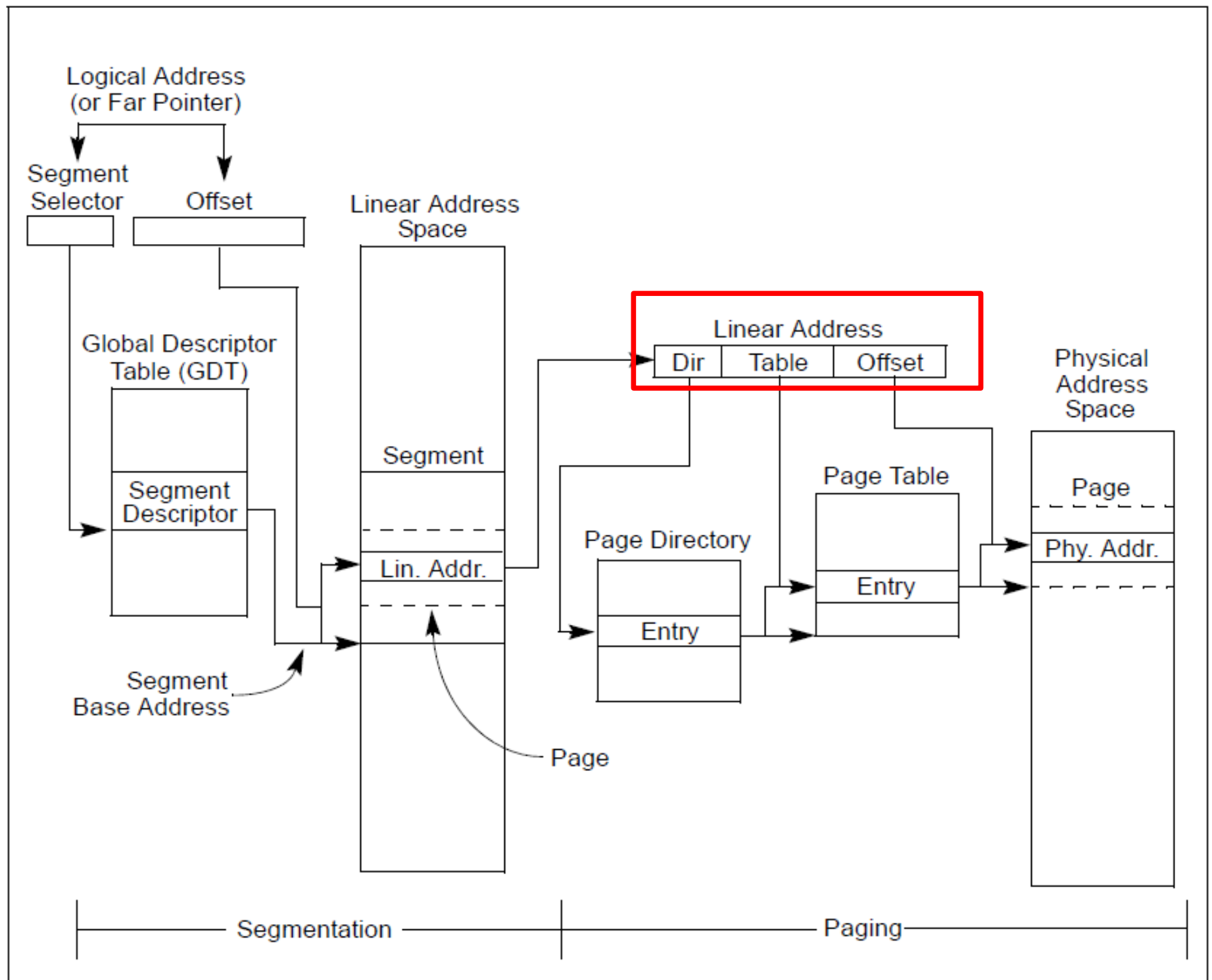
# Address translation

Logical Address
(or Far Pointer)

Segment
Selector              Offset              Linear Address
                                              Space

Global Descriptor
Table (GDT)                                                        Linear Address

                                                               Dir    Table    Offset        Physical
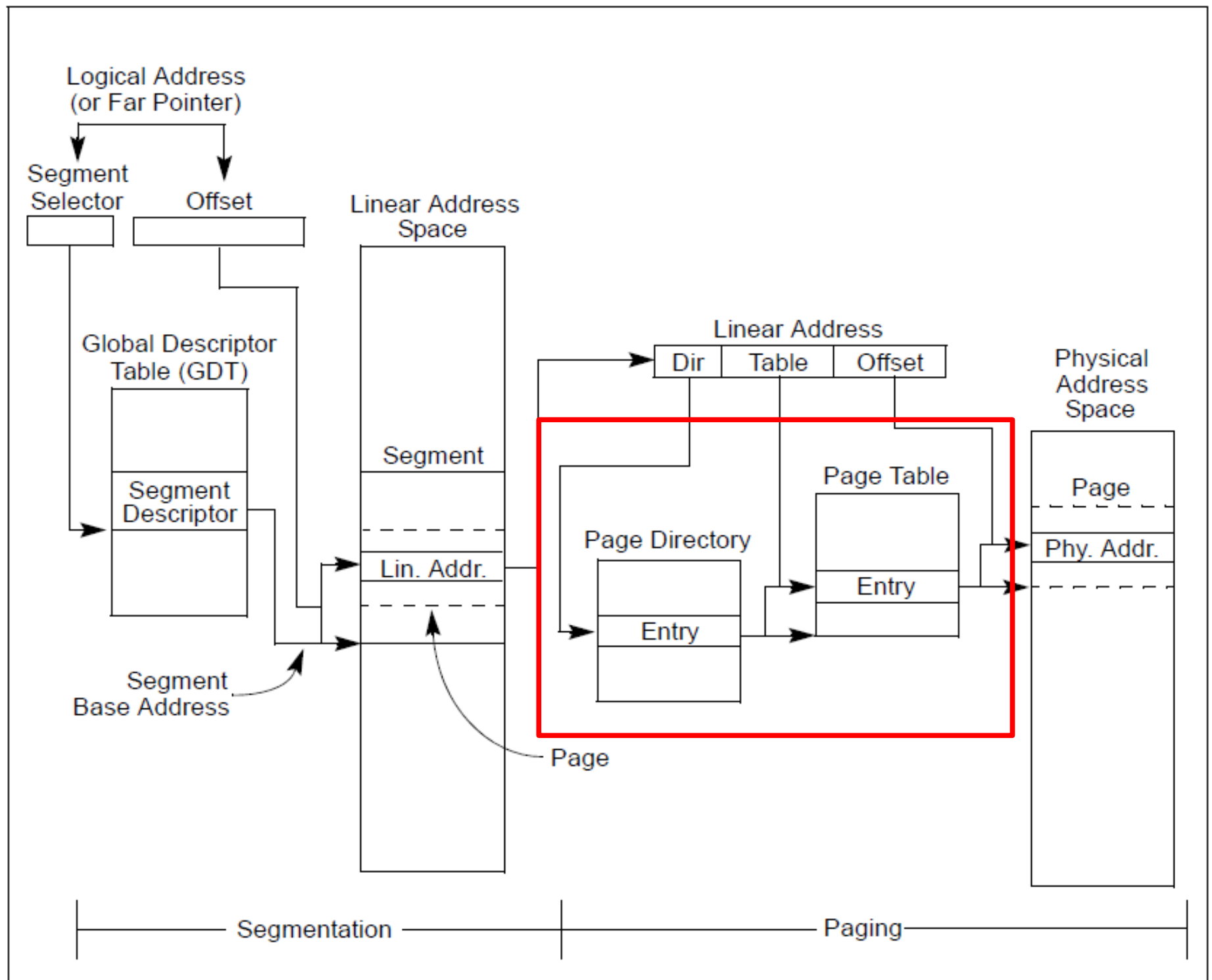                                                                                              Address
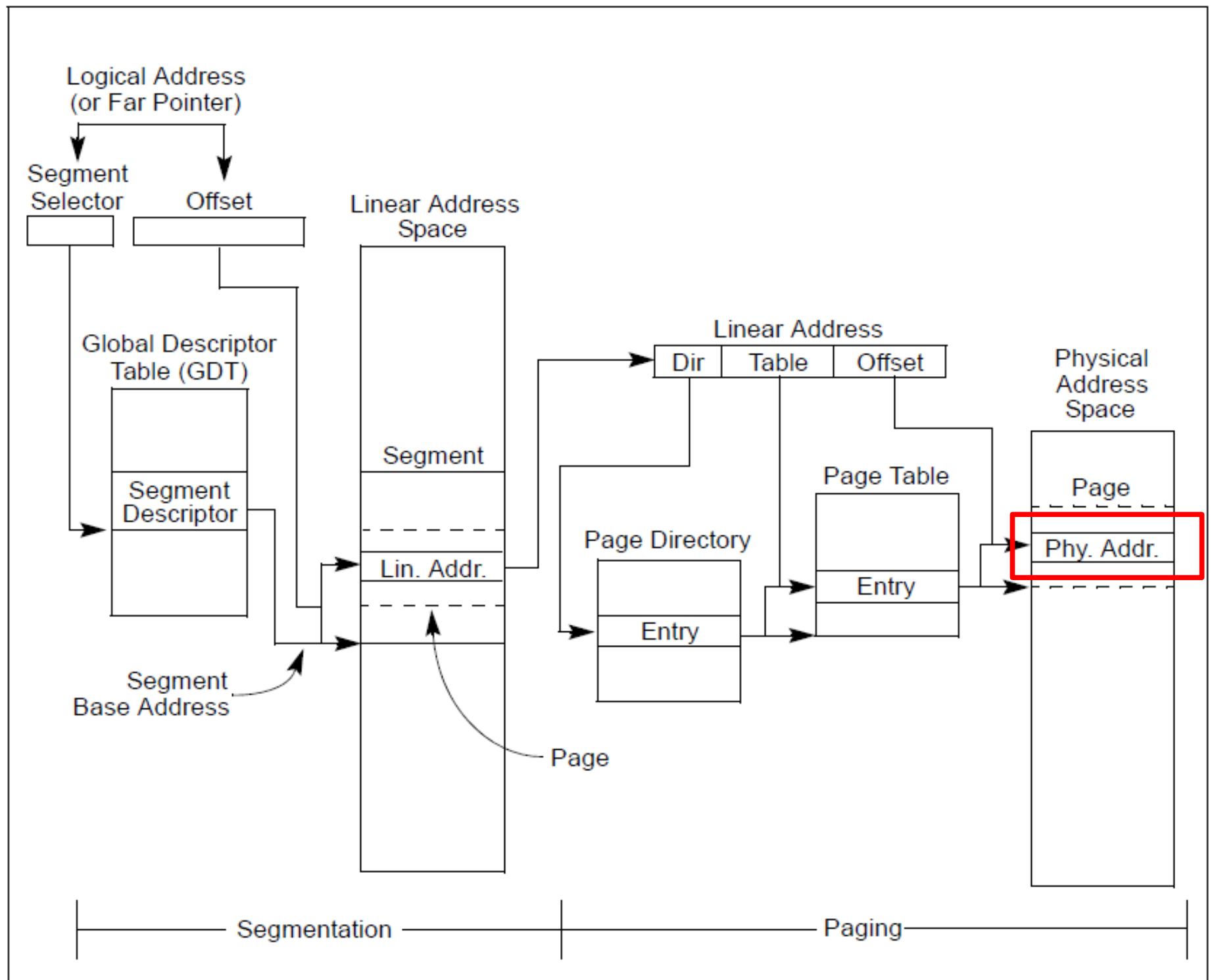                                                                                               Space

                              Segment

        Segment                                      Page Directory              Page Table      Page
        Descriptor
                                                                                                Phy. Addr.
                              Lin. Addr.                                           Entry

                                                        Entry

Segment                                                                                     Page
Base Address

                                          Page

|—————— Segmentation ——————|————————— Paging —————————|

Logical Address
(or Far Pointer)

Segment
Selector

Offset

Linear Address
Space

Global Descriptor
Table (GDT)

Linear Address

Dir | Table | Offset

Physical
Address
Space

Segment
Descriptor

Segment

Page Directory

Page Table

Page

Lin. Addr.

Entry

Phy. Addr.

Segment
Base Address

Entry

Entry

Page

Segmentation

Paging

Logical Address (or Far Pointer)

Segment Selector

Offset

Global Descriptor Table (GDT)

Segment Descriptor

Segment Base Address

Linear Address Space

Segment

Lin. Addr.

Page

Linear Address

Dir | Table | Offset

Page Directory

Entry

Page Table

Entry

Physical Address Space

Page

Phy. Addr.

Segmentation

Paging

Logical Address (or Far Pointer)

Segment Selector

Offset

Linear Address Space

Global Descriptor Table (GDT)

Segment Descriptor

Segment

Segment Base Address

Lin. Addr.

Page

Linear Address

Dir | Table | Offset

Page Directory

Entry

Page Table

Entry

Physical Address Space

Page

Phy. Addr.

Segmentation

Paging

Logical Address
(or Far Pointer)

Segment
Selector     Offset

Linear Address
Space

Global Descriptor
Table (GDT)

**Linear Address**

| Dir | Table | Offset |

Physical
Address
Space

Segment

Segment
Descriptor

Page Directory

Page Table

Page

Lin. Addr.

Entry

Entry

Phy. Addr.

Segment
Base Address

Entry

Page

Segmentation

Paging

Logical Address
(or Far Pointer)

Segment
Selector          Offset

Linear Address
Space

Global Descriptor
Table (GDT)

Linear Address

| Dir | Table | Offset |

Physical
Address
Space

Segment Descriptor

Segment

Page Table

Page

Page Directory

Phy. Addr.

Lin. Addr.

Entry

Entry

Segment
Base Address

Entry

Page

Segmentation                              Paging

Logical Address
(or Far Pointer)

Segment
Selector

Offset

Linear Address
Space

Global Descriptor
Table (GDT)

Segment
Descriptor

Segment

Lin. Addr.

Segment
Base Address

Page

Linear Address

Dir | Table | Offset

Physical
Address
Space

Page Directory

Page Table

Page

Phy. Addr.

Entry

Entry

Segmentation

Paging

What is the linear address? What address is in the registers, e.g., in %eax?

# Logical and linear addresses

- Segment selector (16 bit) + offset (32 bit)

What segments do the following instructions use? push, jump, mov

# What's on the stack? Describe layout of a stack and how it changes during function invocation?

# Example stack

```
:    :
| 10 | [ebp + 16]  (3rd function argument)
|  5 | [ebp + 12]  (2nd argument)
|  2 | [ebp + 8]   (1st argument)
| RA | [ebp + 4]   (return address)
| FP | [ebp]       (old ebp value)
|    | [ebp - 4]   (1st local variable)
:    :
:    :
|    | [ebp - X]  (esp - the current stack pointer)
```

Describe the steps and data structures involved into a user to kernel transition (draw diagrams)

# Interrupt path

**Process**

| |
|---|
| Argument 1 |
| Argument 2 |
| Calling EIP ++ |
| Old EBP |
| Local variables |
| Saved local values, e.g. push EAX, etc |

EBP →

Last stack frame

User stack of a process (can grow up to 2GBs)

Code, data, heap

Interrupt Vector #

Timer: IRQ0 -> vector 32

User state (saved by hardware)

| |
|---|
| SS |
| ESP |
| EFLAGS |
| CS |
| EIP |

ESP →

Kernel Stack of a process (4K)

**GDT**

| |
|---|
| NULL: 0x0 |
| KCODE: 0 - 4GB |
| KDATA: 0 - 4GB |
| K_CPU: 4 bytes |
| CODE: 0 - 4GB |
| DATA: 0 - 4GB |
| TSS: sizeof(ts) |

**IDT**

| |
|---|
| ... |
| CS : HANDLER_ADDR |
| ... |
| ... |

**TSS**

| |
|---|
| ... |
| SS0: |
| ESP0: |
| ... |

Page table
Level 1

| |
|---|
| 0 - 4MB |
| 4 - 8MB |
| ... |
| 2GB - 2GB + 4MB |

Level 2

| |
|---|
| 0 - 4K |
| 4K - 8K |
| ... |
| (4MB-4K) - 4MB |

Kernel code

vector32

| | |
|---|---|
| **CS : #1** | **EIP: <kernel>** |
| **SS : #2** | **ESP: <kernel>** |
| GDT: gdt | TSS: tss |
| IDT: idt | CR3: pt |

# Which stack is used for execution of an interrupt handler? How does hardware find it?

# Interrupt path

**Process**

| |
|---|
| Argument 1 |
| Argument 2 |
| Calling EIP ++ |
| Old EBP |
| Local variables |
| Saved local values, e.g. push EAX, etc |

EBP →

Last stack frame

User stack
of a process
(can grow up to 2GBs)

Code, data,
heap

**Interrupt Vector #**

Timer: IRQ0 -> vector 32

User state
(saved by
hardware)

| |
|---|
| SS |
| ESP |
| EFLAGS |
| CS |
| EIP |

ESP →

Kernel Stack
of a process (4K)

**GDT**

| NULL: 0x0 |
|---|
| KCODE: 0 - 4GB |
| KDATA: 0 - 4GB |
| K_CPU: 4 bytes |
| CODE: 0 - 4GB |
| DATA: 0 - 4GB |
| TSS: sizeof(ts) |

**IDT**

| ... |
|---|
| CS : HANDLER_ADDR |
| ... |
| ... |

**TSS**

| ... |
|---|
| SS0: |
| ESP0: |
| ... |

Kernel
code

vector32

**Page table**
**Level 1**

| 0 - 4MB |
|---|
| 4 - 8MB |
| ... |
| 2GB - 2GB + 4MB |

**Level 2**

| 0 - 4K |
|---|
| 4K - 8K |
| ... |
| (4MB-4K) - 4MB |

| **CS : #1** | **EIP: <kernel>** |
|---|---|
| **SS : #2** | **ESP: <kernel>** |
| GDT: gdt | TSS: tss |
| IDT: idt | CR3: pt |

# Describe organization of the memory allocator in xv6?

# Physical page allocator

Virtual

kmem.freelist

end (0x801126fc)

main

0x80000000
2GB

2GB +4MB
0x80400000

```
CS : 0x8        EIP: main
SS : 0x10       ESP: stack
GDT: 0x7c78     TSS: 0x0
IDT: 0x0        CR3: entrypgdir
```

Protected Mode

# Where did free memory came from?

# How do we switch between processes?

# swtch()

```
2958 swtch:
2959 movl 4(%esp), %eax
2960 movl 8(%esp), %edx
2961
2962 # Save old callee-save registers
2963 pushl %ebp
2964 pushl %ebx
2965 pushl %esi
2966 pushl %edi
2967
2968 # Switch stacksh
2969 movl %esp, (%eax)
2970 movl %edx, %esp
2971
2972 # Load new callee-save registers
2973 popl %edi
2974 popl %esi
2975 popl %ebx
2976 popl %ebp
2977 ret
```

```
2093 struct context {
2094   uint edi;
2095   uint esi;
2096   uint ebx;
2097   uint ebp;
2098   uint eip;
2099 };
```

# Stack inside swtch()

User state
(saved by
hardware)

| SS |
|---|
| ESP |
| EFLAGS |
| CS |
| CS |

vector32

| 0 |
|---|
| 32 |

alltraps

| DS |
|---|
| ES |
| FS |
| GS |
| All registers |

| ESP |
|---|

trap

| EIP (alltraps) |
|---|
| ... |
| EIP (trap) |

yield

| ... |
|---|
| EIP (yield) |

sched

| ... |
|---|
| &proc->context |
| cpu->scheduler |
| EIP (sched) |

### Kernel Stack
of a process (4K)

Trap frame

### Call stack:
vector32()
alltraps()
trap()
yield()
sched()
switch(&proc->context,
     cpu->scheduler)

### Proc

Context

### Context

| EIP (line: 2479) |
|---|
| EBP |
| EBX |
| ESI |
| EDI |

# What is the interface between the kernel and user-level processes?

```
3374 void
3375 syscall(void)
3376 {
3377   int num;
3378
3379   num = proc->tf->eax;
3380   if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
3381     proc->tf->eax = syscalls[num]();
3382   } else {
3383     cprintf("%d %s: unknown sys call %d\n",
3384     proc->pid, proc->name, num);
3385     proc->tf->eax = -1;
3386   }
3387 }
```

```
3374 void
3375 syscall(void)
3376 {
3377   int num;
3378
3379   num = proc->tf->eax;
3380   if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
3381     // proc->tf->eax = syscalls[num]();
         proc->tf->esp -= 4;
         *(int*)ptoc->tf->esp = syscalls[num]();
3382   } else {
3383     cprintf("%d %s: unknown sys call %d\n",
3384             proc->pid, proc->name, num);
3385     // proc->tf->eax = -1;
         proc->tf->esp -= 4;
         *(int*)ptoc->tf->esp = -1;
3386   }
3387 }
```

# Why do we need locks?

# Request queue (e.g. incoming network packets)

list

- Linked list, list is pointer to the first element

# CPU 1 and 2 allocate new request

list

# CPU 1 and 2 update next pointer

l->next = list

list

l->next = list

# CPUs 1 updates head pointer

list = l

list

CPU2 updates head pointer

list = l

# State after the race

list = l

list

# List implementation with locks

```
9  insert(int data)
10 {
11   struct list *l;

13   l = malloc(sizeof *l);

     acquire(&listlock);

14   l->data = data;

15   l->next = list;

16   list = l;

     release(&listlock);

17 }
```

- Critical section

# Correct implementation

```
1573 void

1574 acquire(struct spinlock *lk)

1575 {

...

1580    // The xchg is atomic.

1581    while(xchg(&lk->locked, 1) != 0)

1582      ;

...

1592 }
```

# Xchg instruction

- Swap a word in memory with a new value
  - Atomic!
  - Return old value

# Deadlocks

acquire(A)

acquire(B)

acquire(B) {
    while(xchg(&B−>locked, 1) != 0)
}

acquire(A) {
    while(xchg(&A−>locked, 1) != 0)
}

# Lock ordering

- Locks need to be acquired in the same order

# Locks and interrupts

Network
interrupt

network_packet(){

   ....

   insert() {

      acquire(A)

     ...

   }


}

network_packet(){

   ....

   insert() {

      acquire(A)

     ...

   }


}

# Locks and interrupts

- Never hold a lock with interrupts enabled

# Send/receive queue

```
201 void*
202 send(struct q *q, void *p)
203 {
204   while(q->ptr != 0)
205     ;
206   q->ptr = p;
207   wakeup(q); /*wake recv*/
208 }
```

```
210 void*
211 recv(struct q *q)
212 {
213   void *p;
214
215   while((p = q->ptr) == 0)
216     sleep(q);
217   q->ptr = 0;
218   return p;
219 }
```

# Lost wakeup problem

recv

215
test

216
sleep

wait for wakeup forever

Time

send

206
store p

207
wakeup

204
test

205
spin forever

# The role of file systems

- Sharing
  - Sharing of data across users and applications
- Persistence
  - Data is available after reboot

# Architecture

- On-disk and in-memory data structures represent

  - The tree of named files and directories

  - Record identities of disk blocks which hold data for each file

  - Record which areas of the disk are free

# Crash recovery

- File systems must support crash recovery
    - A power loss may interrupt a sequence of updates
    - Leave file system in inconsistent state
        - E.g. a block both marked free and used

# Multiple users

- Multiple users operate on a file system concurrently
  - File system must maintain invariants

# Speed

- Access to a block device is several orders of magnitude slower

  - Memory: 200 cycles

  - Disk: 20 000 000 cycles

- A file system must maintain a cache of disk blocks in memory

# Block layer

| System calls | File descriptors |
|---|---|
| Pathnames | Recursive lookup |
| Directories | Directory inodes |
| Files | Inodes and block allocator |
| Transactions | Logging |
| Blocks | Buffer cache |

- Read and write data
  - From a block device
  - Into a buffer cache
- Synchronize across multiple readers and writers

# Transactions

| | |
|---|---|
| System calls | File descriptors |
| Pathnames | Recursive lookup |
| Directories | Directory inodes |
| Files | Inodes and block allocator |
| Transactions | Logging |
| Blocks | Buffer cache |

- Group multiple writes into an atomic transaction

# Files

| | |
|---|---|
| System calls | File descriptors |
| Pathnames | Recursive lookup |
| Directories | Directory inodes |
| Files | Inodes and block allocator |
| Transactions | Logging |
| Blocks | Buffer cache |

- Unnamed files
  - Represented as inodes
  - Sequence of blocks holding file's data

# Directories

| | |
|---|---|
| System calls | File descriptors |
| Pathnames | Recursive lookup |
| Directories | Directory inodes |
| Files | Inodes and block allocator |
| Transactions | Logging |
| Blocks | Buffer cache |

- Special kind of inode
  - Sequence of directory entries
  - Each contains name and a pointer to an unnamed inode

# Pathnames

| | |
|---|---|
| System calls | File descriptors |
| Pathnames | Recursive lookup |
| Directories | Directory inodes |
| Files | Inodes and block allocator |
| Transactions | Logging |
| Blocks | Buffer cache |

- Hierarchical path names
  - /usr/bin/sh
  - Recursive lookup

# System call

| System calls | File descriptors |
|---|---|
| Pathnames | Recursive lookup |
| Directories | Directory inodes |
| Files | Inodes and block allocator |
| Transactions | Logging |
| Blocks | Buffer cache |

- Abstract UNIX resources as files
  - Files, sockets, devices, pipes, etc.
- Unified programming interface

# File system layout on disk

| boot | super | inodes | bit map | data | .... | data | log |
|------|-------|--------|---------|------|------|------|-----|

0      1      2

- Block #0: Boot code
- Block #1: Metadata about the file system
  - Size (number of blocks)
  - Number of data blocks
  - Number of inodes
  - Number of blocks in log

# File system layout on disk

| boot | super | inodes | bit map | data | .... | data | log |
|------|-------|--------|---------|------|------|------|-----|

0     1     2

- Block #2 (inode area)

- Bit map area: track which blocks are in use

- Data area: actual file data

- Log area: maintaining consistency in case of a power outage or system crash

```
begin_op();
...
bp = bread(...);
bp->data[...] = ...;
log_write(bp);
...
end_op();
```

# Typical use of transactions

# Strawman scheduler (xv6)

- Organize all processes as a simple list

- In schedule():

  - Pick first one on list to run next

  - Put suspended task at the end of the list

- Problem?

# Xv6 scheduler

```
2458 scheduler(void)
2459 {

2462   for(;;){

2468     for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){
2469       if(p->state != RUNNABLE)
2470         continue;

2475       proc = p;
2476       switchuvm(p);
2477       p->state = RUNNING;
2478       swtch(&cpu->scheduler, proc->context);
2479       switchkvm();

2483       proc = 0;
2484     }
2487   }
2488 }
```

# Strawman scheduler (xv6)

- Organize all processes as a simple list

- In schedule():

  - Pick first one on list to run next

  - Put suspended task at the end of the list

- Problem?

# Strawman scheduler

- Organize all processes as a simple list

- In schedule():

  - Pick first one on list to run next

  - Put suspended task at the end of the list

- Problem?

  - Only allows round-robin scheduling

  - Can't prioritize tasks

# Priority based scheduling

- Higher-priority processes run first
- Processes within the same priority are round-robin

# O(1) scheduler (Linux 2.6 – 2.6.22)

- Priority based scheduling
- Goal: decide who to run next, independent of number of processes in system
  - Still maintain ability to prioritize tasks, handle partially unused quanta, etc

# O(1) data structures

- runqueue: a list of runnable processes
  - Blocked processes are not on any runqueue
  - A runqueue belongs to a specific CPU
  - Each task is on exactly one runqueue
  - Task only scheduled on runqueue's CPU unless migrated
- 2 *40 * #CPUs runqueues
  - 40 dynamic priority levels (more later)
  - 2 sets of runqueues – one active and one expired

# O(1) data structures (contd)

# O(1) intuition

- Take the first task off the lowest-numbered runqueue on active set

  - Confusingly: a lower priority value means higher priority

- When done, put it on appropriate runqueue on expired set

- Once active is completely empty, swap which set of runqueues is active and expired

- Constant time, since fixed number of queues to check; only take first item from non-empty queue

# O(1) example

# Blocked tasks

- What if a program blocks on I/O, say for the disk?
    - It still has part of its quantum left
    - Not runnable, so don't waste time putting it on the active or expired runqueues
- We need a "wait queue" associated with each blockable event
    - Disk, lock, pipe, network socket, etc.

# Blocking example

# Buddy memory allocator

# Buddy allocator

# What's wrong with buddy?

# What's wrong with buddy?

- Buddy allocator is ok for large allocations
    - E.g. 1 page or more
- But what about small allocations?
    - Buddy uses the whole page for a 4 bytes allocation
        - Wasteful
    - Buddy is still slow for short-lived objects
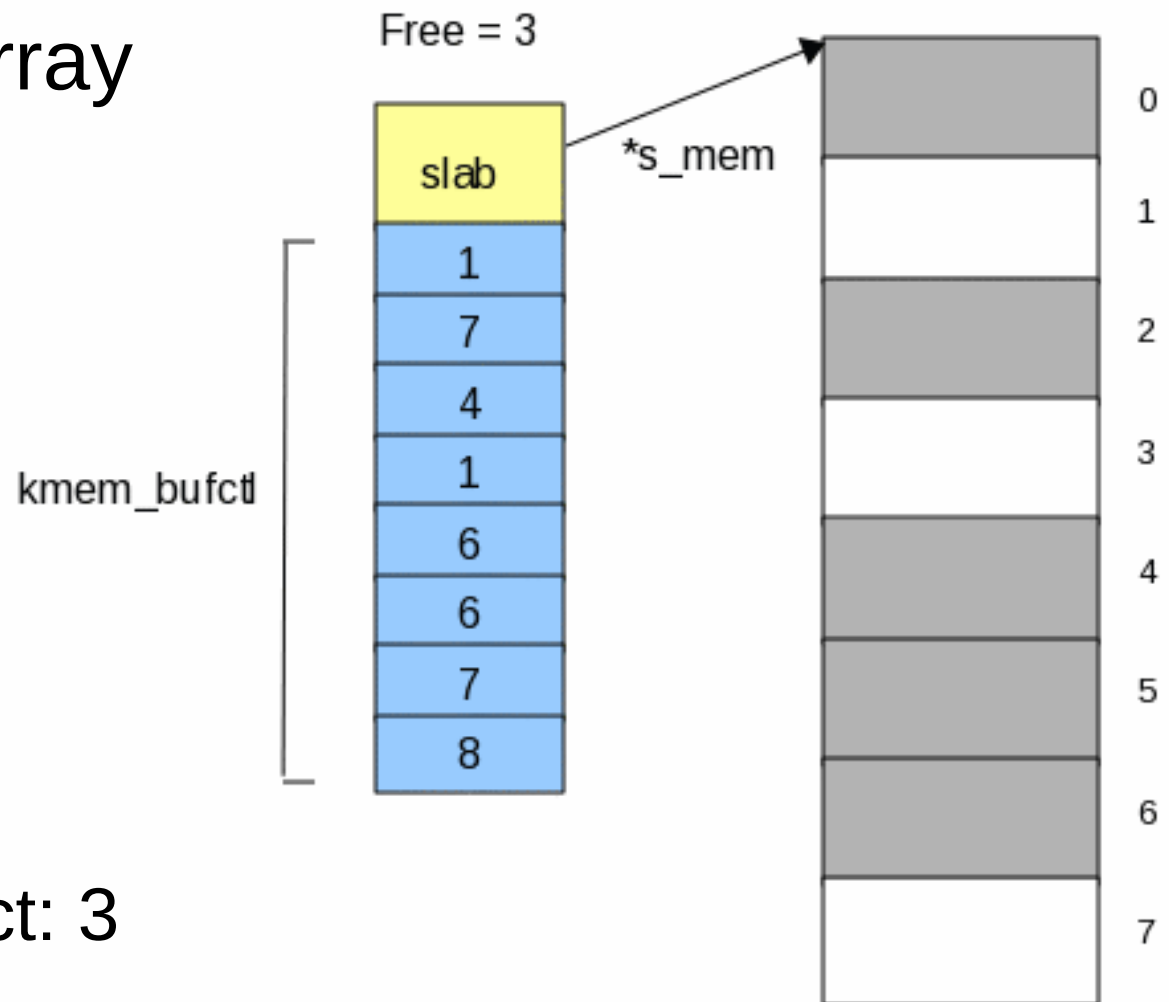
# Slab:
## Allocator for object of a fixed size
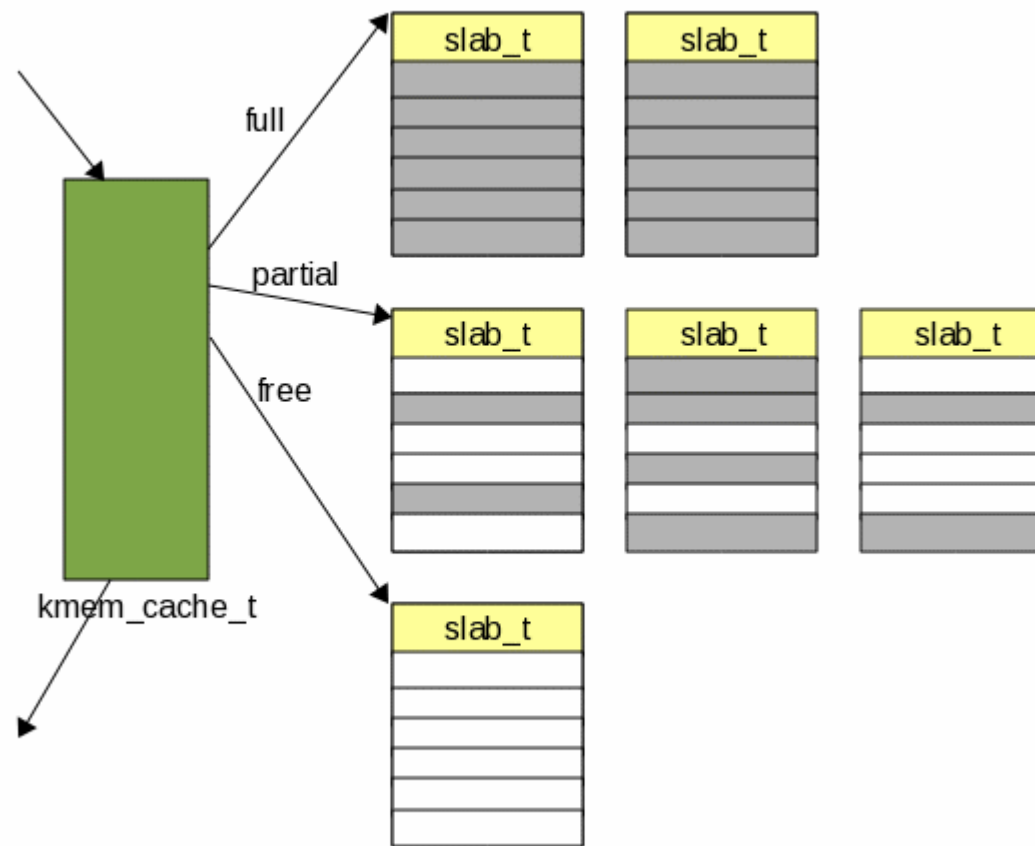
# Slab

- A 2 page slab with 6 objects

# Keeping track of free objects

- kmem_bufctl array is effectively a linked list

Free = 3

slab    *s_mem

kmem_bufctl

| 1 |
| 7 |
| 4 |
| 1 |
| 6 |
| 6 |
| 7 |
| 8 |

0
1
2
3
4
5
6
7

- First free object: 3
- Next free object: 1
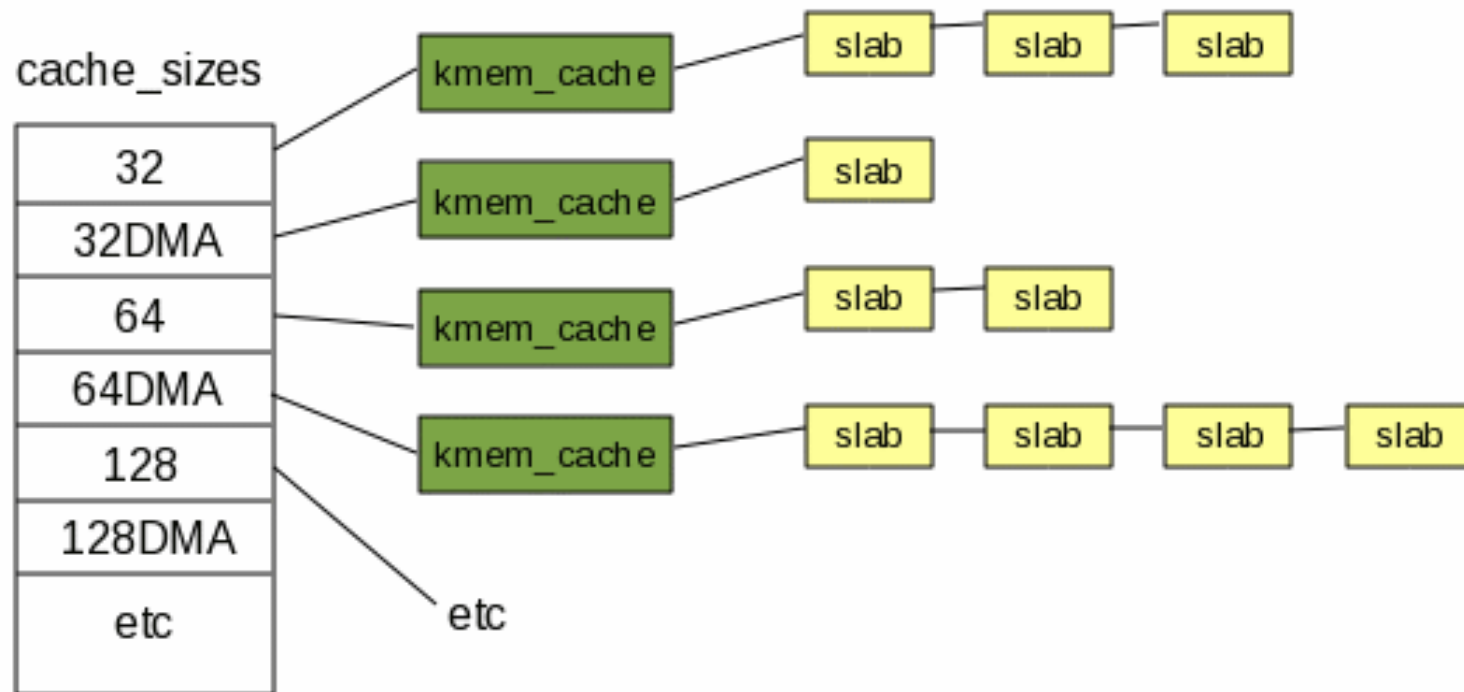
# A cache is formed out of slabs

# Slab is fine, but what's wrong?

# Slab is fine, but what's wrong?

- We can only allocate objects of one size

# Kmalloc(): variable size objects

- A table of caches
  - Size: 32, 64, 128, etc.

# Thank you!