# STOS - Pagination

Gabriel Laskar <gabriel@lse.epita.fr>

# Vocabulary

- MMU
- Page
- Frame
- Address Space
- Logical Address
- Linear Address
- Physical Address

# Virtual Memory?

- Separate Address Space (per-process)
- "Map" a virtual address to a physical address
- Fine grained allocation (Page granularity)
- Per-page permissions (R/W, U/S)

# Module : pagination.ko

```
MODINFO {
    module_name("pagination"),
    module_init_once(pagination_init),
    module_type(M_PAGINATION | M_PAGE_ALLOCATOR),
    module_deps(M_INTERRUPTS | M_FRAME_ALLOCATOR)
};

EXPORT_SYMBOL(alloc_pages);
EXPORT_SYMBOL(map_pages);
EXPORT_SYMBOL(map_io);
EXPORT_SYMBOL(unmap_pages);
```

# Pagination API (<kernel/page.h>)

```c
void* alloc_pages(struct frame* pdbr, size_t n);

int map_pages(struct frame* pdbr, void* vaddr,
              struct frame* frames[],
              size_t n, int flags);

void unmap_pages(struct frame* pdbr, void* vaddr, int n);

void* map_io(phys_t ioaddr, size_t len);
```

# Code already in STOS

- Macros and Flags are defined inside `<kernel/arch/page.h>`
- Generic page flags are consumed by the api
  - P_KERNEL
  - P_USER_RO
  - P_USER_RW
- "`frame_allocator.ko`" is here to allocate physical frames.
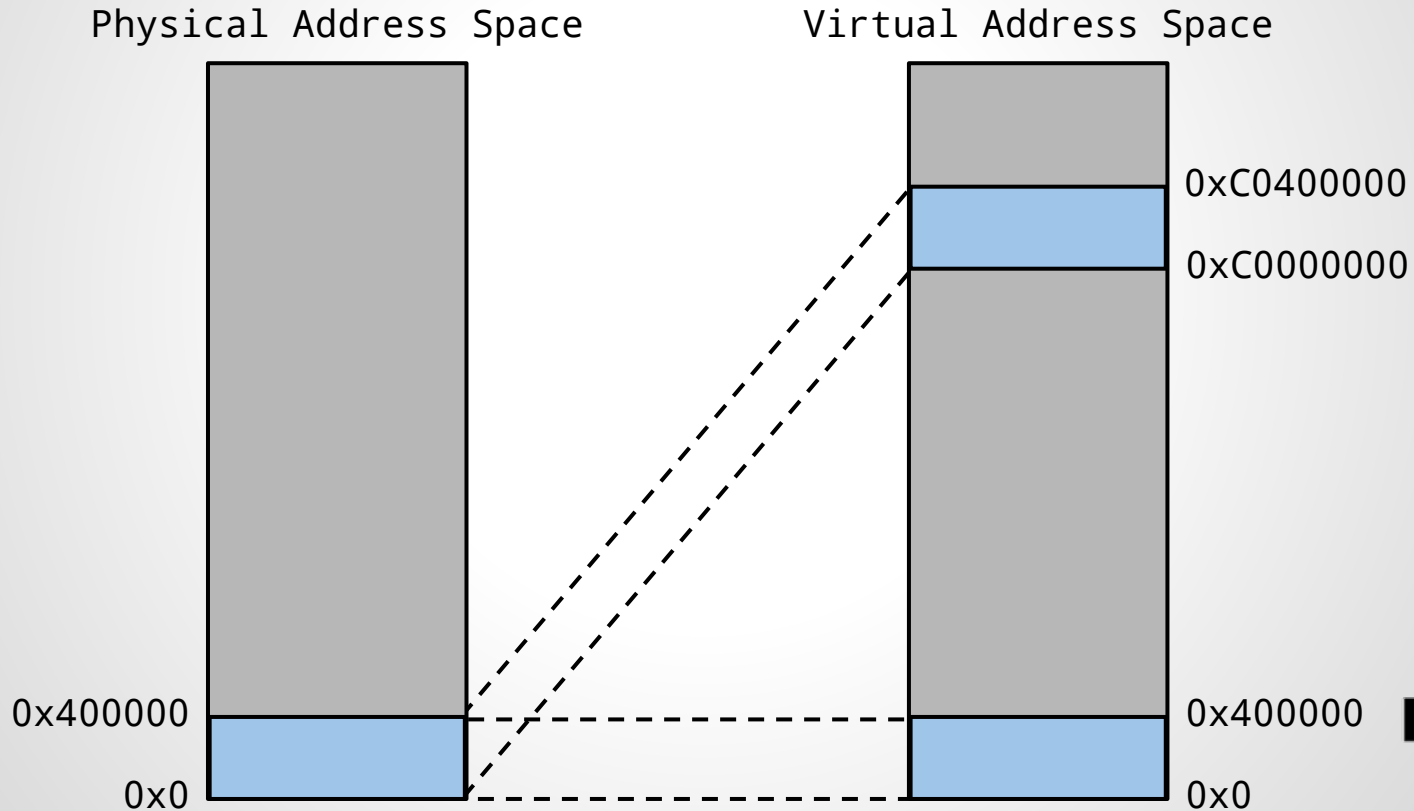- `<kernel/memory.h>` contains physical memory informations
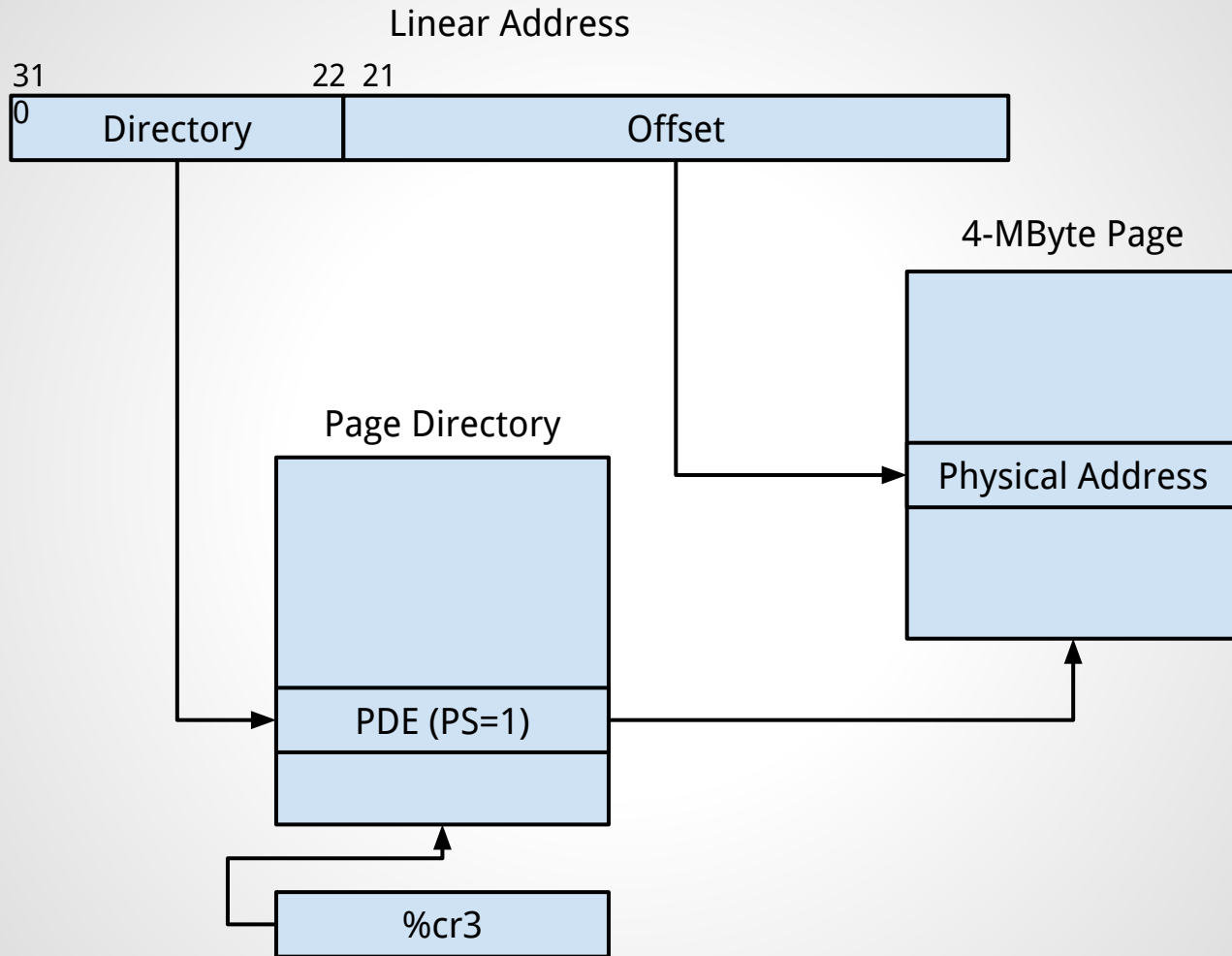
# Frame Allocator (<kernel/frame.h>)

```c
struct frame* alloc_frame(void);
void free_frame(struct frame* frame);

int alloc_frames(u8 n, struct frame** frames);

static inline phys_t frame_to_phys(struct frame* f);
static inline struct frame* phys_to_frame(phys_t
addr);
```

# More on frame allocator

- frames are stored in an array (from `framestart` through `frameend`)
- index of a frame is the frame number
- inside a frame, you must maintain the `vaddr`
- kernel maintain only a free list
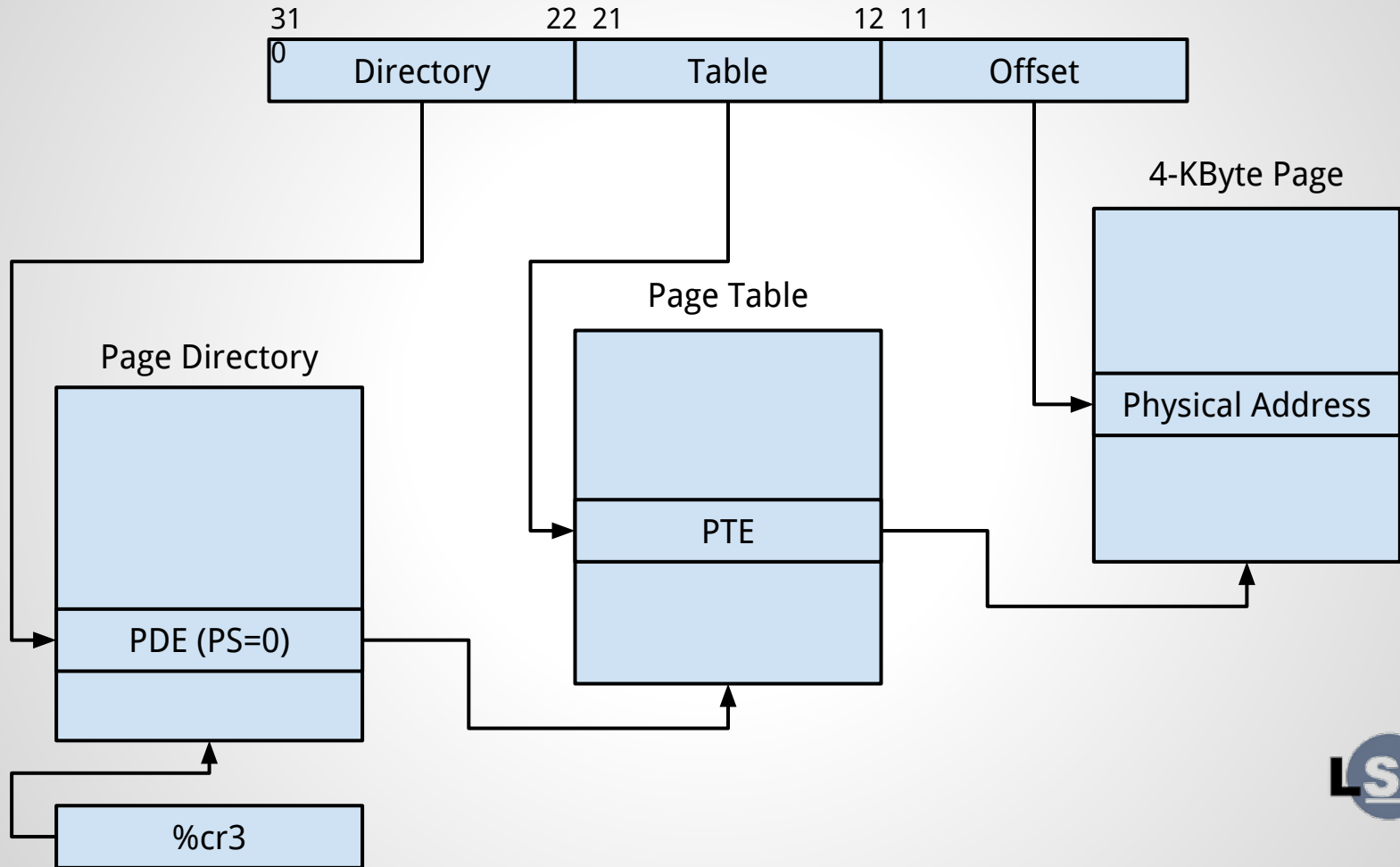
# Reminder: Current paging state

Physical Address Space

Virtual Address Space

0xC0400000

0xC0000000

0x400000

0x400000

0x0

0x0

# Linear Address

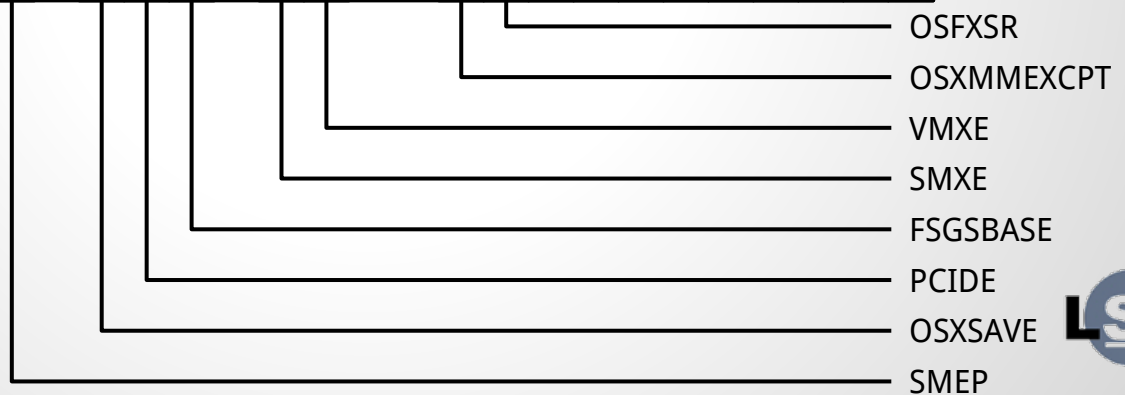| 31 | 22 | 21 | 0 |
|---|---|---|---|
| Directory | | Offset | |

## 4-MByte Page

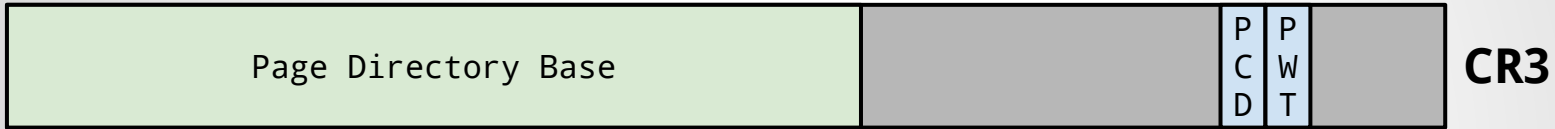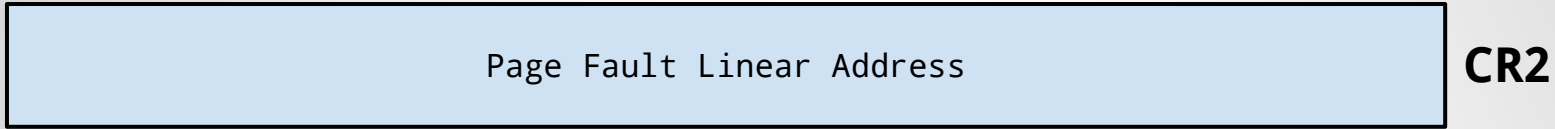## Page Directory

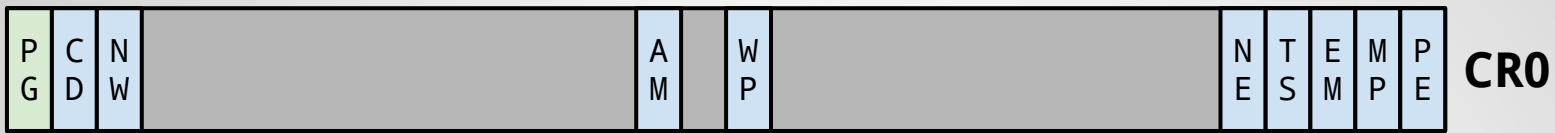| |
|---|
| PDE (PS=1) |
| |

Physical Address

%cr3

# Pagination on x86_32 (vol 3a, chapter 4)

- 3 paging modes for x86
- multiple page sizes
- control registers used:
    - %cr0: activation (**PG**, WP)
    - %cr2: Page Fault Linear Address
    - %cr3: Page Directory Base
    - %cr4: extra features (PAE, **PSE**, **PGE**, PCIDE, SMEP)
    - IA32_EFER MSR: (LME, NXE)

# Linear Address

| | Directory | Table | Offset |
|---|---|---|---|

31        22 21        12 11

0

## 4-KByte Page

Physical Address

## Page Table

PTE

## Page Directory

PDE (PS=0)

%cr3

# PDE and PTE

| addr [21:22] | 0 | addr [39:32] | P A T | | | G | 1 | D | A | P C D | P W T | U / S | R / W | P | PDE 4MB Page |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| Address of Page Table | | | | 0 | | A | P C D | P W T | U / S | R / W | P | PDE Page table |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| Address of 4KB Page Frame | | G | P A T | D | A | P C D | P W T | U / S | R / W | P | PTE 4KB Page |
|---|---|---|---|---|---|---|---|---|---|---|---|

- R/W: Read/Write
- U/S: User/System
- PWT: Page Level write-through
- PCD: Page Level Cache disable

- A: Accessed
- D: Dirty
- G: Global (if `%cr4.pge = 1`)
- PAT: Reserved

# Page Fault Handling

| | I/D | RSVD | U/S | W/R | P |
|---|---|---|---|---|---|

- Which address? Content of `%cr2`
- Error Code:
  - P: non-present (clear), page-level protection violation (set)
  - W/R: read (clear) or write (set) error
  - U/S: supervisor (clear) or user-mode (set)
  - RSVD: reserved bit violation (set)
  - I/D: data (clear) or instruction (set)

# Translation Lookaside Buffer

- Page Directory Walk costs
- TLB cache the page walks
- TLB is automatically managed on x86
- Flush on :
    - mov %cr3
    - on change of some %cr registers
    - `invlpg addr:` instruction