# Operating Systems Practice
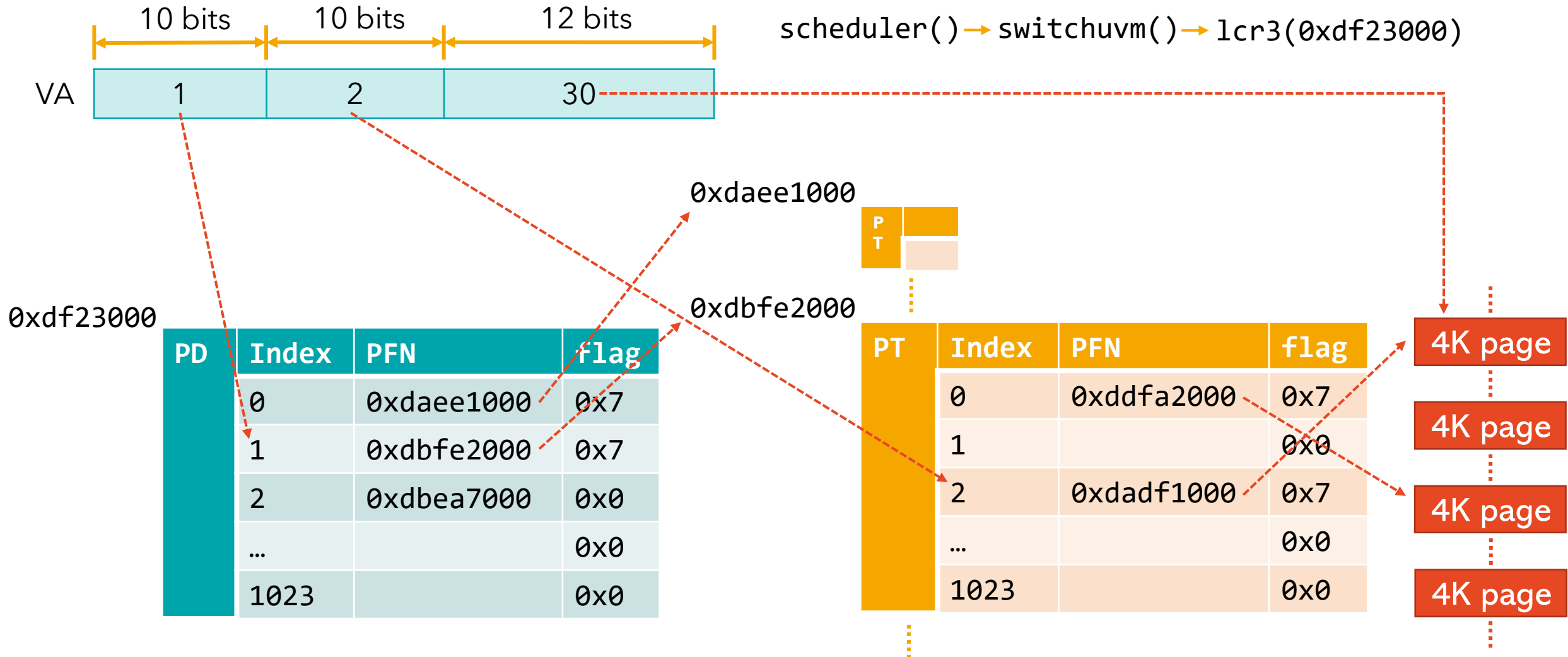
Project #2 - 3-level paging (Part 1)

Juhyung Park
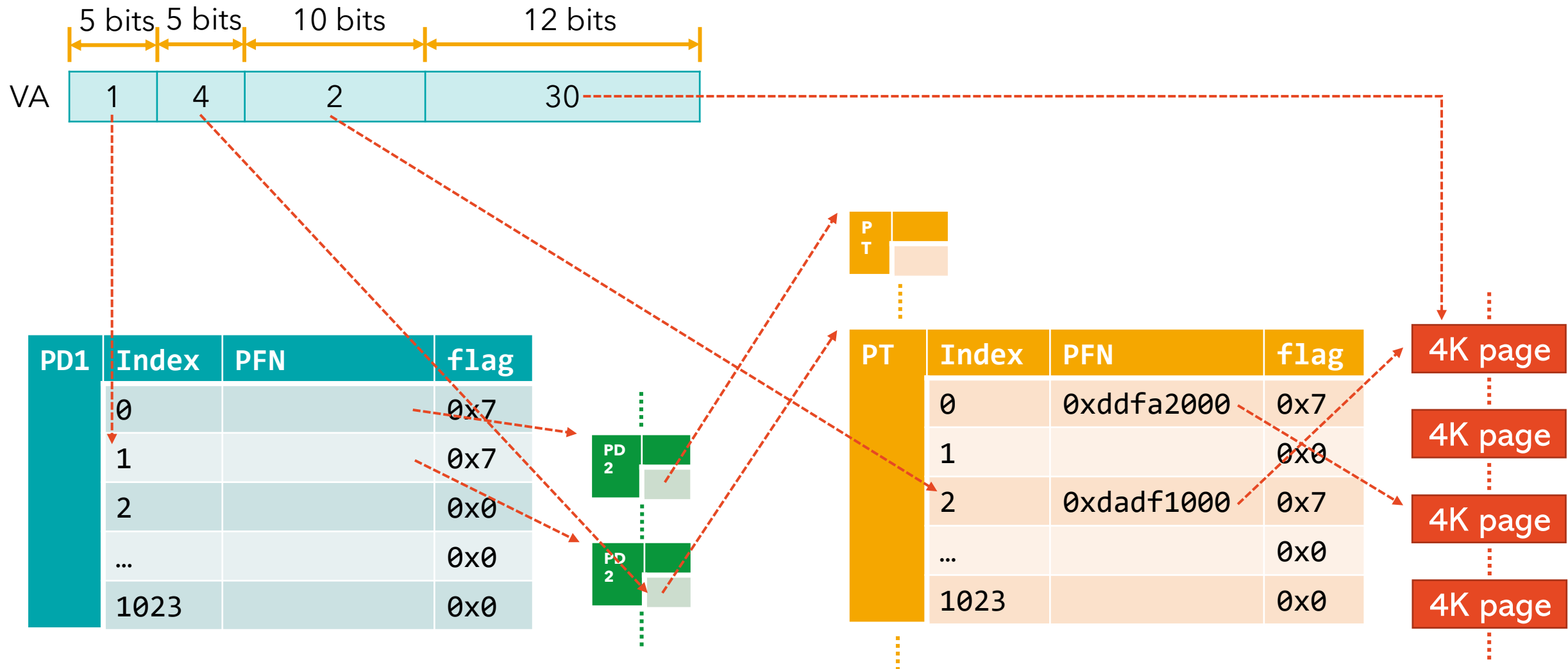
arter97@dgist.ac.kr
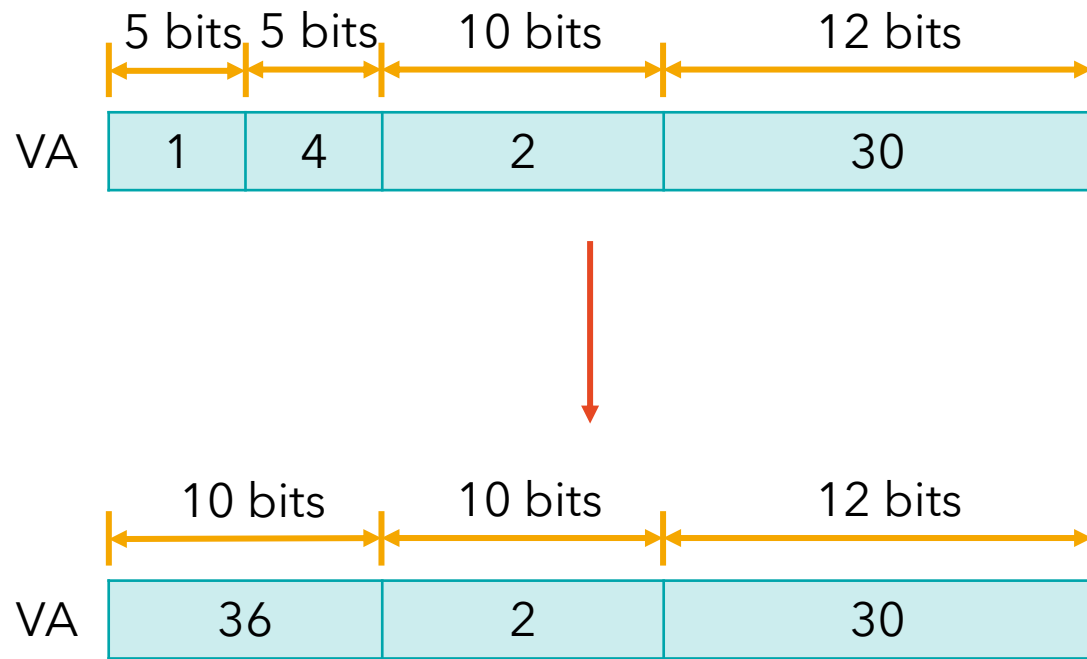
# Recall how VA is translated to PA (2-level page table)

| 10 bits | 10 bits | 12 bits |
|---------|---------|---------|

scheduler() → switchuvm() → lcr3(0xdf23000)

VA

| 1 | 2 | 30 |
|---|---|----|

0xdaee1000

0xdbfe2000

0xdf23000

| PD | Index | PFN | flag |
|----|-------|-----|------|
| | 0 | 0xdaee1000 | 0x7 |
| | 1 | 0xdbfe2000 | 0x7 |
| | 2 | 0xdbea7000 | 0x0 |
| | … | | 0x0 |
| | 1023 | | 0x0 |

| PT | Index | PFN | flag |
|----|-------|-----|------|
| | 0 | 0xddfa2000 | 0x7 |
| | 1 | | 0x0 |
| | 2 | 0xdadf1000 | 0x7 |
| | … | | 0x0 |
| | 1023 | | 0x0 |

4K page

4K page

4K page

4K page

# How can we customize VA to PA translation? (e.g., 3-level)

# x86-32 MMU still expects 2-level VA



VA diagram (top): 5 bits | 5 bits | 10 bits | 12 bits → values 1 | 4 | 2 | 30

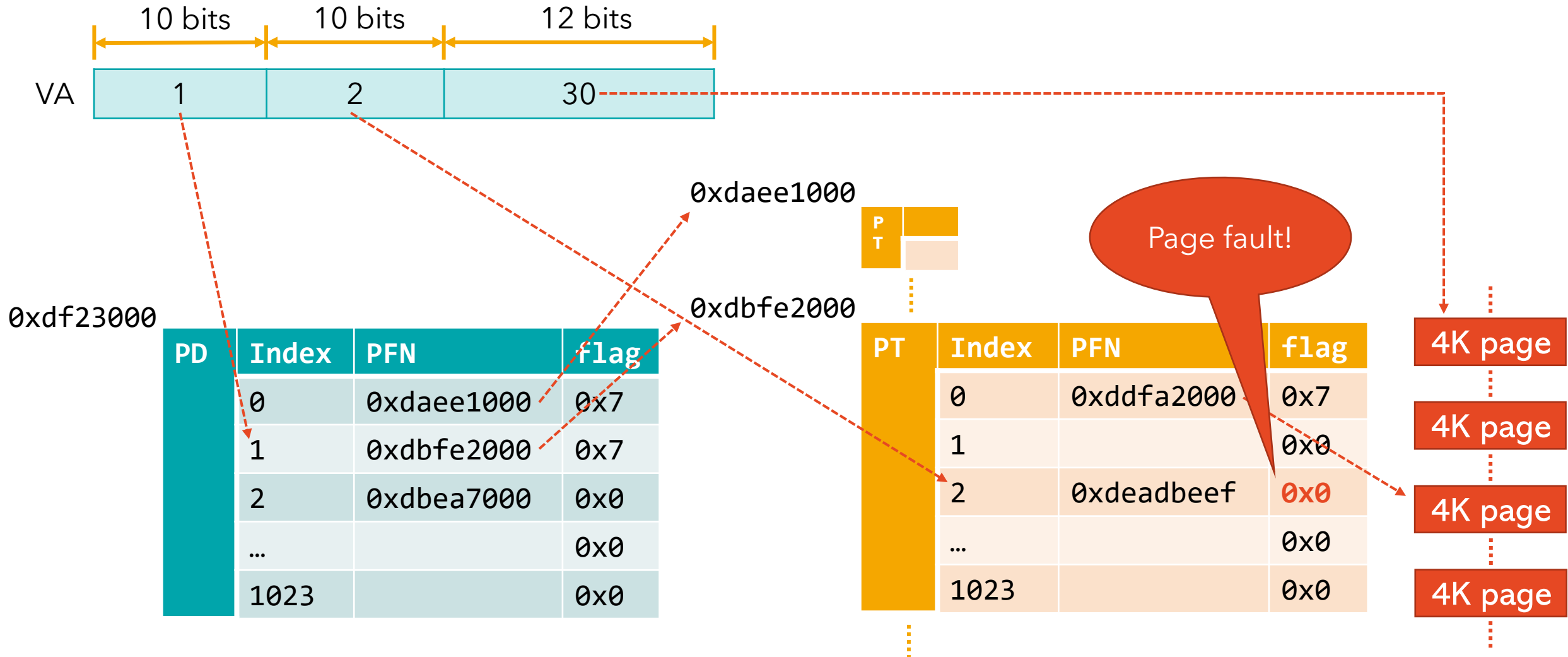VA diagram (bottom): 10 bits | 10 bits | 12 bits → values 36 | 2 | 30
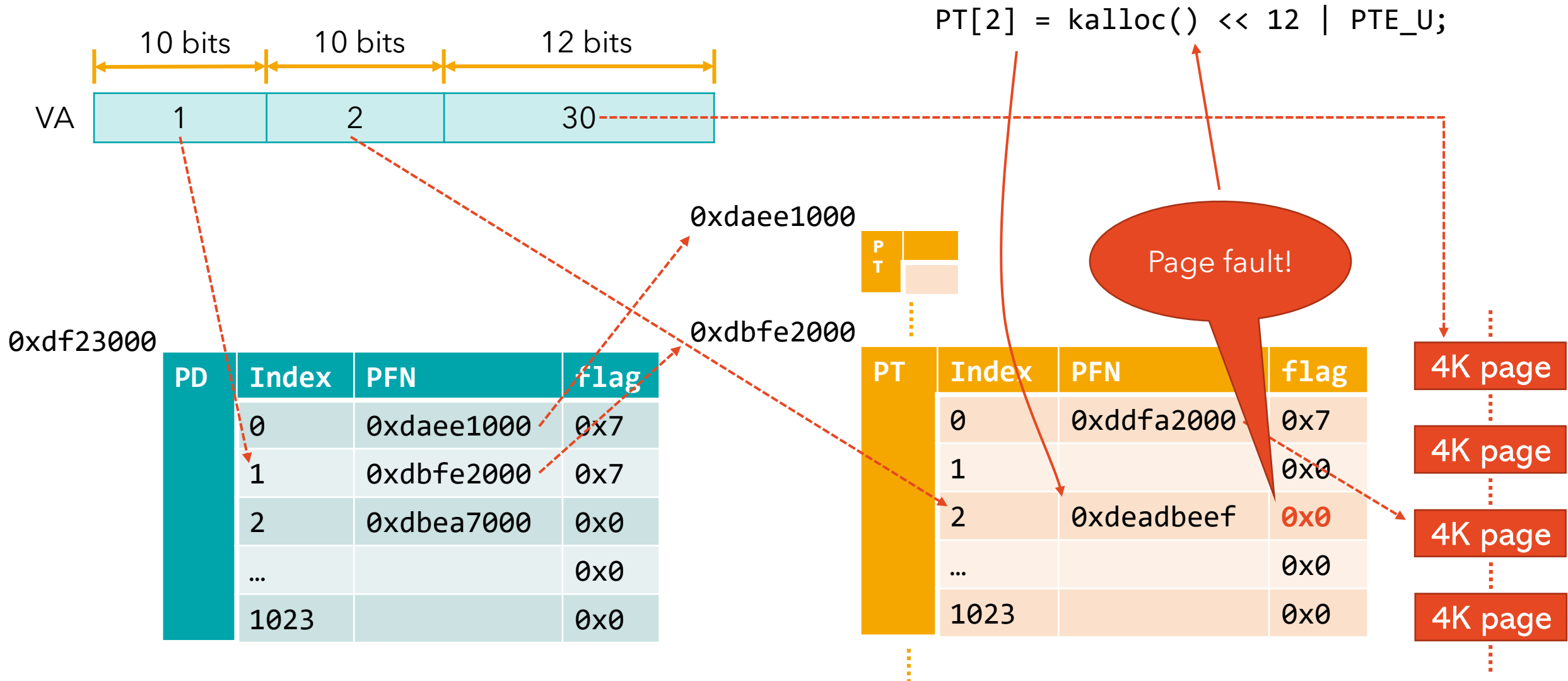
**Error!**

How can we implement custom page addressing within software?

# Page fault (1 - occurs when PTE_P flag is not set)

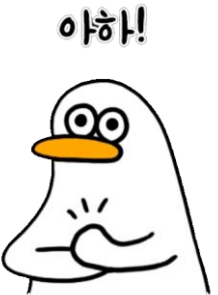# Page fault (2 - kernel's page fault handler is executed)

```
PT[2] = kalloc() << 12 | PTE_U;
```

10 bits    10 bits    12 bits

VA | 1 | 2 | 30 |

Page fault!

0xdaee1000

0xdbfe2000

0xdf23000

| PD | Index | PFN | flag |
|----|-------|-----|------|
|  | 0 | 0xdaee1000 | 0x7 |
|  | 1 | 0xdbfe2000 | 0x7 |
|  | 2 | 0xdbea7000 | 0x0 |
|  | … |  | 0x0 |
|  | 1023 |  | 0x0 |

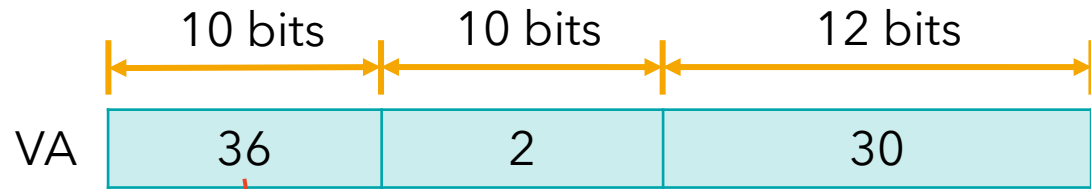| PT | Index | PFN | flag |
|----|-------|-----|------|
|  | 0 | 0xddfa2000 | 0x7 |
|  | 1 |  | 0x0 |
|  | 2 | 0xdeadbeef | 0x0 |
|  | … |  | 0x0 |
|  | 1023 |  | 0x0 |

4K page

4K page

4K page

4K page

# Page fault (3 - code execution is resumed from page fault handler)

Triggering page fault on every page access allows the software to fully redirect VA!
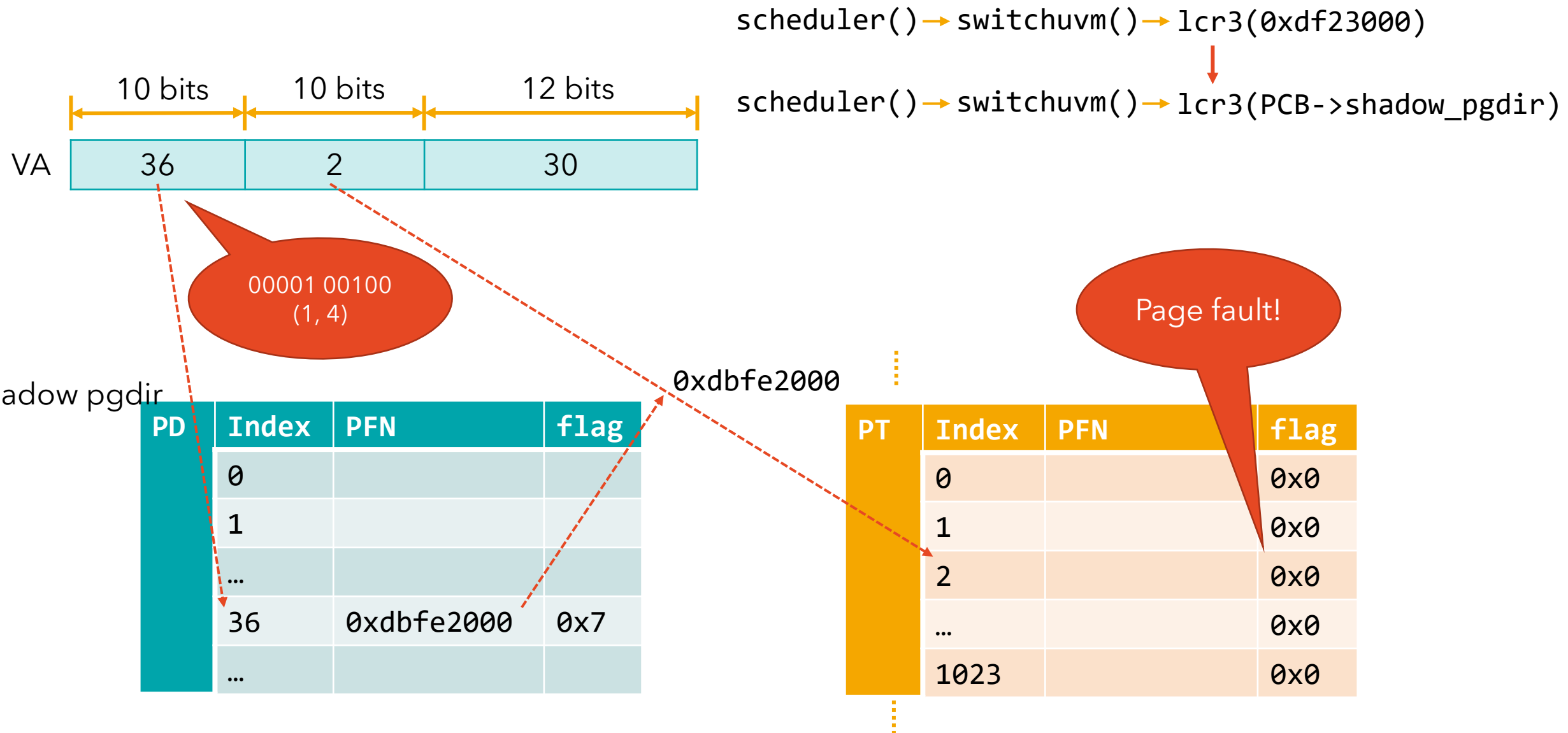
아하!

# How?

scheduler() → switchuvm() → lcr3(0xdf23000)

scheduler() → switchuvm() → lcr3(PCB->shadow_pgdir)

|        | 10 bits | 10 bits | 12 bits |
|--------|---------|---------|---------|
| VA     | 36      | 2       | 30      |

00001 00100
(1, 4)

Page fault!

Shadow pgdir

| PD | Index | PFN | flag |
|----|-------|-----|------|
|    | 0     |     |      |
|    | 1     |     |      |
|    | …     |     |      |
|    | 36    |     | 0x0  |
|    | …     |     |      |

# How?

# How?

scheduler() → switchuvm() → lcr3(0xdf23000)

scheduler() → switchuvm() → lcr3(PCB->shadow_pgdir)

| 10 bits | 10 bits | 12 bits |
|---------|---------|---------|

VA | 36 | 2 | 30 |

00001 00100
(1, 4)

Shadow pgdir

| PD | Index | PFN | flag |
|----|-------|-----|------|
| | 0 | | |
| | 1 | | |
| | … | | |
| | 36 | 0xdbfe2000 | 0x7 |
| | … | | |

0xdbfe2000

Insert the page address you want

Page fault!

| PT | Index | PFN | flag |
|----|-------|-----|------|
| | 0 | | 0x0 |
| | 1 | | 0x0 |
| | 2 | 0xda3b0000 | 0x7 |
| | … | | 0x0 |
| | 1023 | | 0x0 |

# Software TLB

- Trigger a page fault on every new page access

- Reference the vanilla pgdir and make a copy to shadow_pgdir


- We maintain 2-level paging in this project but the CPU's MMU must never access PCB->pgdir
  - PCB->shadow_pgdir must return the appropriate PA instead

# Trigger a page fault on every new page access (1/2)

- Setting the CR3 register tells the hardware to use it as the current page directory address
  - When is switchuvm() called?
  - Can we abuse this to force page faults?
  - Hint: a new empty pgdir can be allocated with setupkvm();

```
vm.c:

// Switch TSS and h/w page table to correspond to process p.
void
switchuvm(struct proc *p)
{

  ...

  ltr(SEG_TSS << 3);
  lcr3(V2P(p->pgdir));  // switch to process's address space
  popcli();
}
```

# Reference the vanilla pgdir and make a copy to shadow_pgdir (2/2)

- The CPU's MMU must never access PCB->pgdir
  - It should use PCB->shadow_pgdir

- You must copy PA from PCB->pgdir and set it up to PCB->shadow_pgdir
  - Recall what walkpgdir() function did in vm.c

```c
void pagefault(void)
{

    ...

    // Map pgdir's page address to shadow_pgdir's page table
    // XXX

    ...
}
```

# Counting page accesses

- Once shadow_pgdir is set, page fault handler no longer runs if it was already handled before

- To count page accesses, shadow_pgdir must be cleared

- shadow_pgdir may contain addresses to page table itself

- Clearing all of shadow_pgdir will cause recursion loop during page fault

# Handling sbrk'ed memory separately

- We use a reserved bit available on PDE

## Page Directory Entry (4 MB)

| 31 ... 22 | 21 | 20 ... 13 | 12 | 11 ... 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Bits 31-22 of address | RSVD (0) | Bits 39-32 of address | PAT | AVL | G | PS (1) | D | A | PCD | PWT | U/S | R/W | P |

## Page Directory Entry

| 31 ... 12 | 11 ... 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|
| Bits 31-12 of address | AVL | PS (0) | AVL | A | PCD | PWT | U/S | R/W | P |

| | |
|---|---|
| **P:** Present | **D:** Dirty |
| **R/W:** Read/Write | **PS:** Page Size |
| **U/S:** User/Supervisor | **G:** Global |
| **PWT:** Write-Through | **AVL:** Available |
| **PCD:** Cache Disable | **PAT:** Page Attribute |
| **A:** Accessed | Table |

```
#define PTE_SBRK        (1 << 10)  // sbrk'ed memory (malloc)
```

# Handling sbrk'ed memory separately

- We use a reserved bit available on PDE

- malloc => sbrk => growproc => PTE_SBRK set!

- pagefault => proc->last_va => if PTE_SBRK is set, clear!

- This allows page fault handler to run again on the same memory address
  - Which in turn allows us to "count" its occurrence

- https://github.com/dgist-datalab/xv6/commit/27ce2ad

# Tips

- Try to understand 2-level paging fully before writing actual code

- Be careful of pointer accesses

- Reference mmu.h for helper macros
  - P2V, PDX, PTX, PTE_ADDR, …

- Have a look at skeleton code's commit history
  - https://github.com/dgist-datalab/xv6/commits/page-counter


- cprintf() is buggy when called from page fault handler
  - You can set #define LOG 1 from vm.c to forcefully enable clprintf()
  - If a page fault is triggered during cprintf()…
  - It'll cause a crash during normal operations

```
pagefault++
Page fault by process "init" (pid: 1) at 0x38a
virt_to_phys: translated "init"(1)'s VA 0x38a to PA 0xdf3838a (pgdir)
Clearing last_pde_entry (0x8dffe000)
virt_to_phys: translated "init"(1)'s VA 0x38a to PA 0xdf3838a (shadow_pgdir)
pagefault--
pagefault++
Page fault by process "init" (pid: 1) at 0x38a
virt_to_phys: translated "init"(1)'s VA 0x38a to PA 0xdf3838a (pgdir)
Clearing last_pde_entry (0x8dffe000)
virt_to_phys: translated "init"(1)'s VA 0x38a to PA 0xdf3838a (shadow_pgdir)
pagefault--
lapicid 0: panic: acquire
 801046aa 801007ed 80107203 80105b65 8010588f 80101170 80104e29 80104ad9 80105ab5 8010588f
```

# In action

- `#define LOG 1 (vm.c)`

# Extra points

- If you found a mistake in the skeleton code

- If you found a different/better method

- Kernel memory leak test has been disabled in usertests
  - https://github.com/dgist-datalab/xv6/commit/63776c9cc
  - If you manage to pass this test (after reverting the commit)


- Write a report for extra points

# Finally …

## **<span style="color:red">Do NOT hesitate</span>** to ask questions!

# Thank You!