# CS162
# Operating Systems and
# Systems Programming
# Lecture 15

# Virtual Memory (2)

Professor Natacha Crooks

https://cs162.org/

# Recall: Memory Management Wishlist

Memory Protection

Memory Sharing

Flexible Memory Placement

Support for Sparse Addresses

Runtime Lookup Efficiency

Compact Translation Table

# Recall: Increasingly powerful mechanisms

**No protection. Living life on the edge**

**Can access all memory**

**Base & Bound**

**Absolute memory addressing. Hard to relocate**

**Base & Bound with Relocation**

**Internal fragmentation when address space is sparse**

**Segmentation**

**External fragmentation as assigning variably sized chunks**

**Paging**

# Paging

Divide logical address space of process into fixed sized chunks called pages

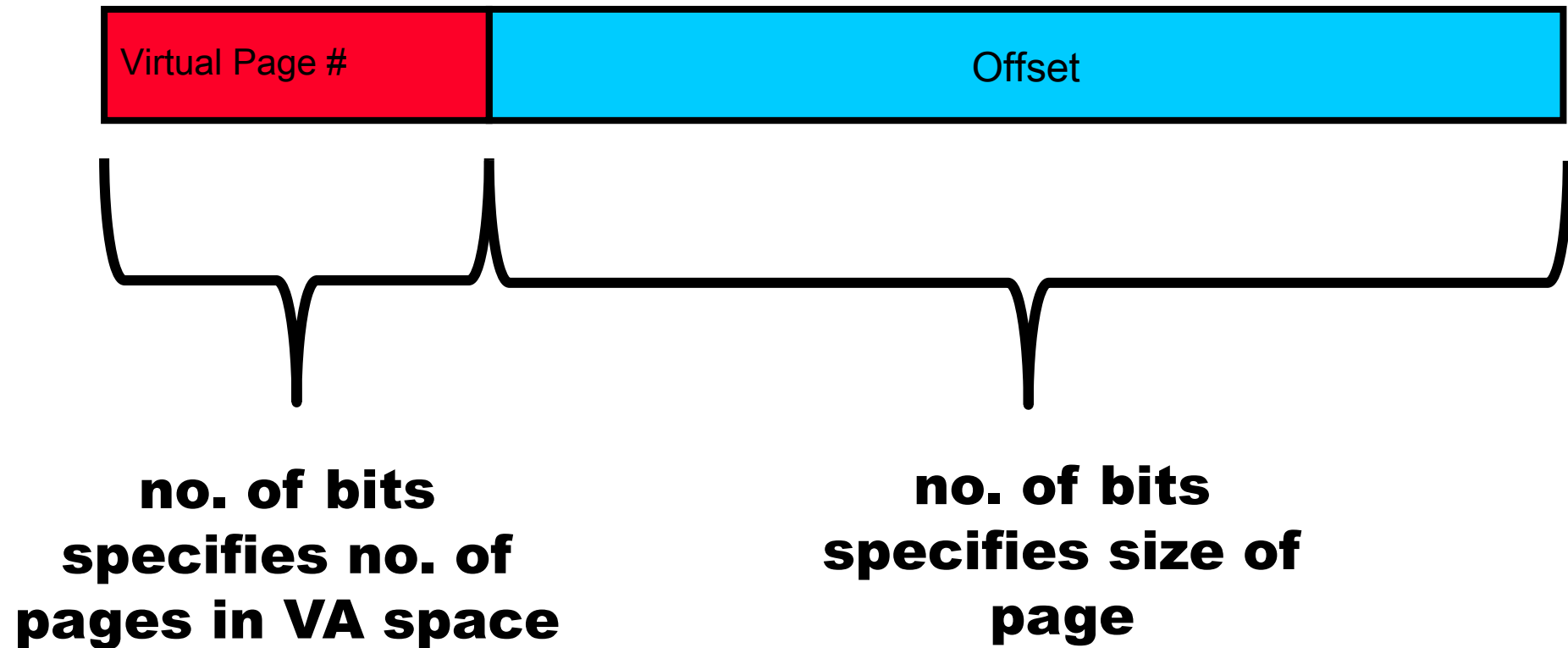View physical memory as an array of fixed-sized slots called **page frames**

Each page frame can contain a

single virtual-memory page

Pages should be small to minimise internal fragmentation (1K-16k)
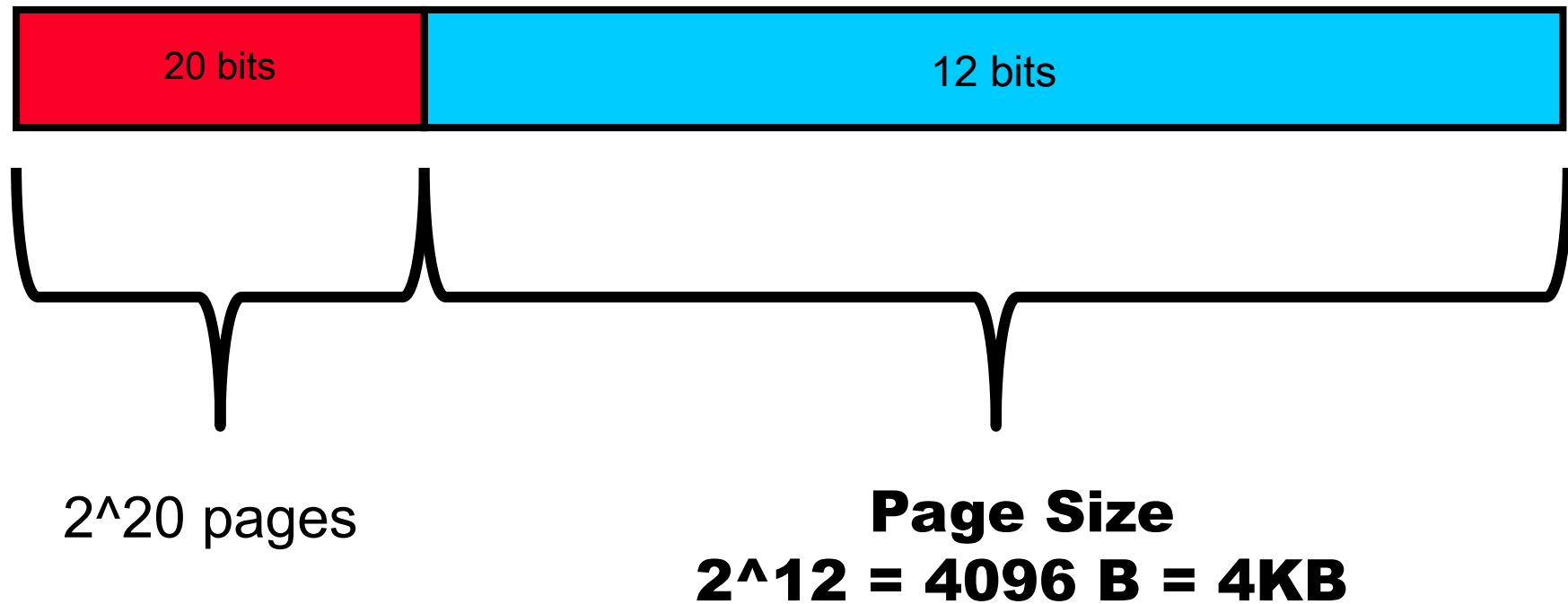
# How to Implement Simple Paging?

Interpret virtual address as two components

| Virtual Page # | Offset |
|---|---|

**no. of bits specifies no. of pages in VA space**

**no. of bits specifies size of page**

# How to Implement Simple Paging?

Interpret virtual address as two components

| 20 bits | 12 bits |
|:---:|:---:|

$2^{20}$ pages

**Page Size**
**$2^{12}$ = 4096 B = 4KB**

# A (Simplified) Page Table

A page table stores

virtual-to-physical address translations

One page table per process. Lives in memory.

Address stored in the in the Page Table Base Register

PTBR value saved/restored in PCB on context switch

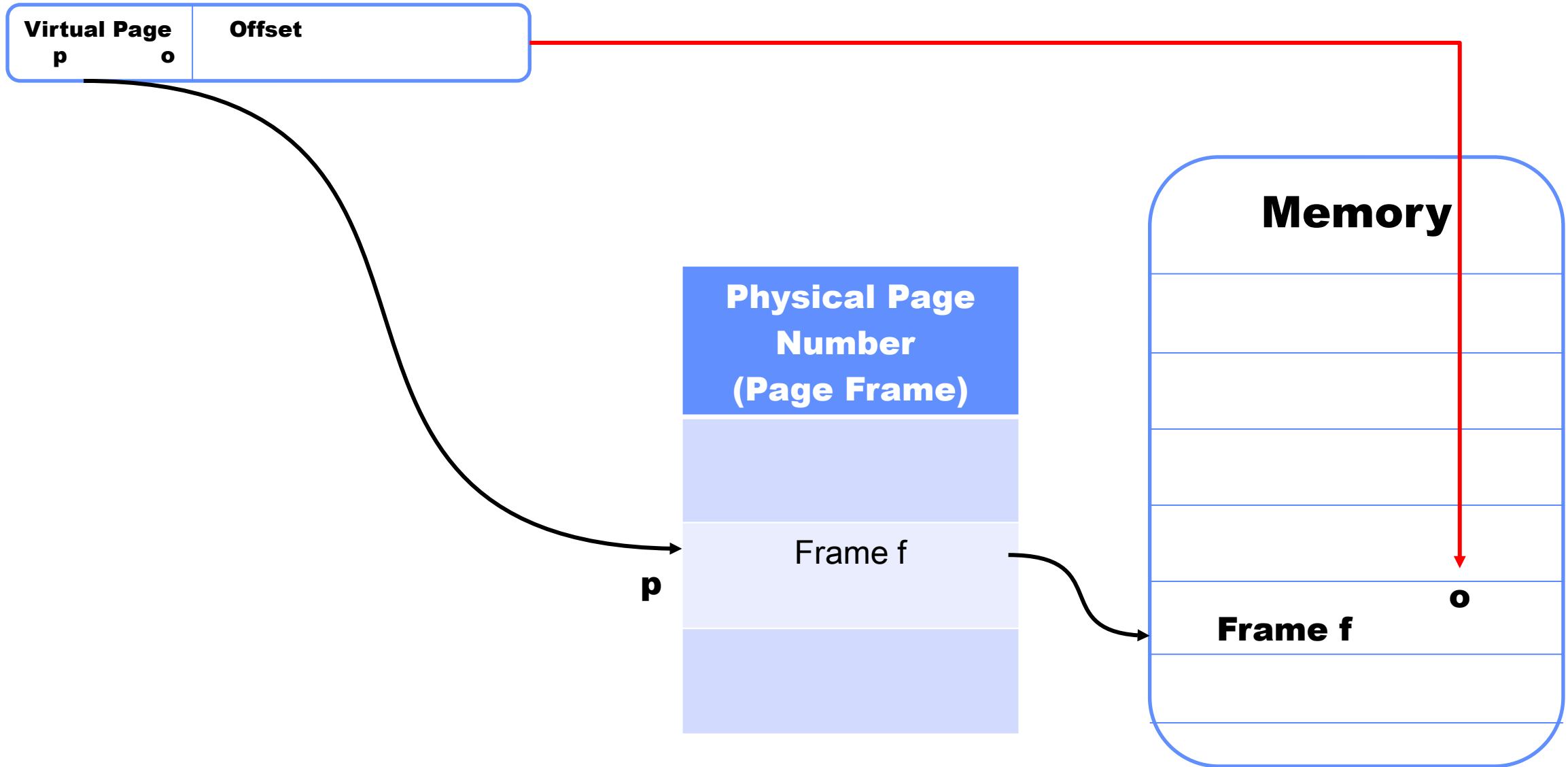# How to access a byte?

Extract page number (first p bits)

Map virtual page number into a frame number
(also called physical page number) using a page table

Extract offset (last o bits)

Convert to physical memory location: access byte at offset in frame

# A (Simplified) Page Table

| Virtual Page p | Offset o |
| --- | --- |

**Physical Page Number (Page Frame)**

p → Frame f

**Memory**

Frame f

o

# Example: A Mini Page Table

Assume we have a 64 bytes ($2^6$) of physical memory

Assume we want pages of 4 bytes ($2^2$)

How long should our addresses be?
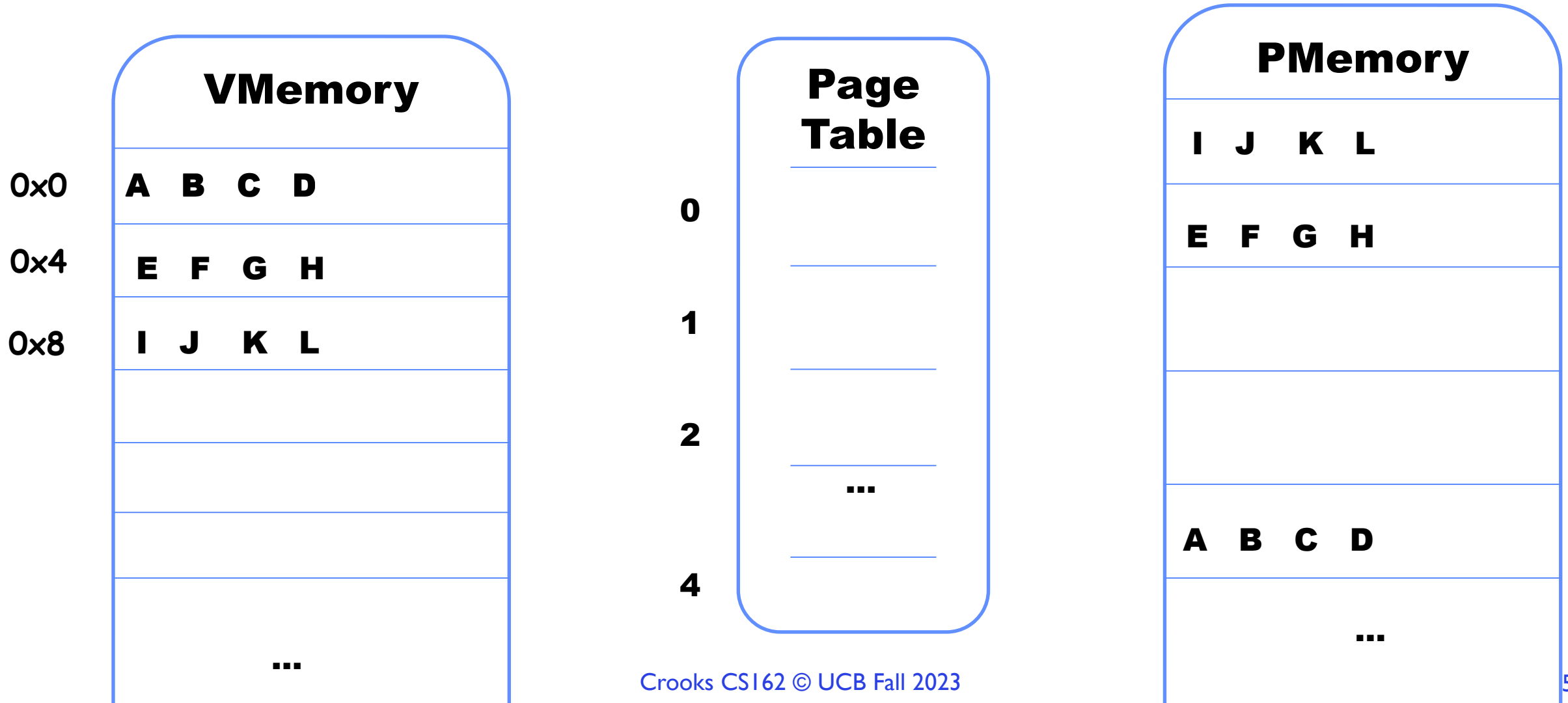
6 bits

How many offset bits should we assign?
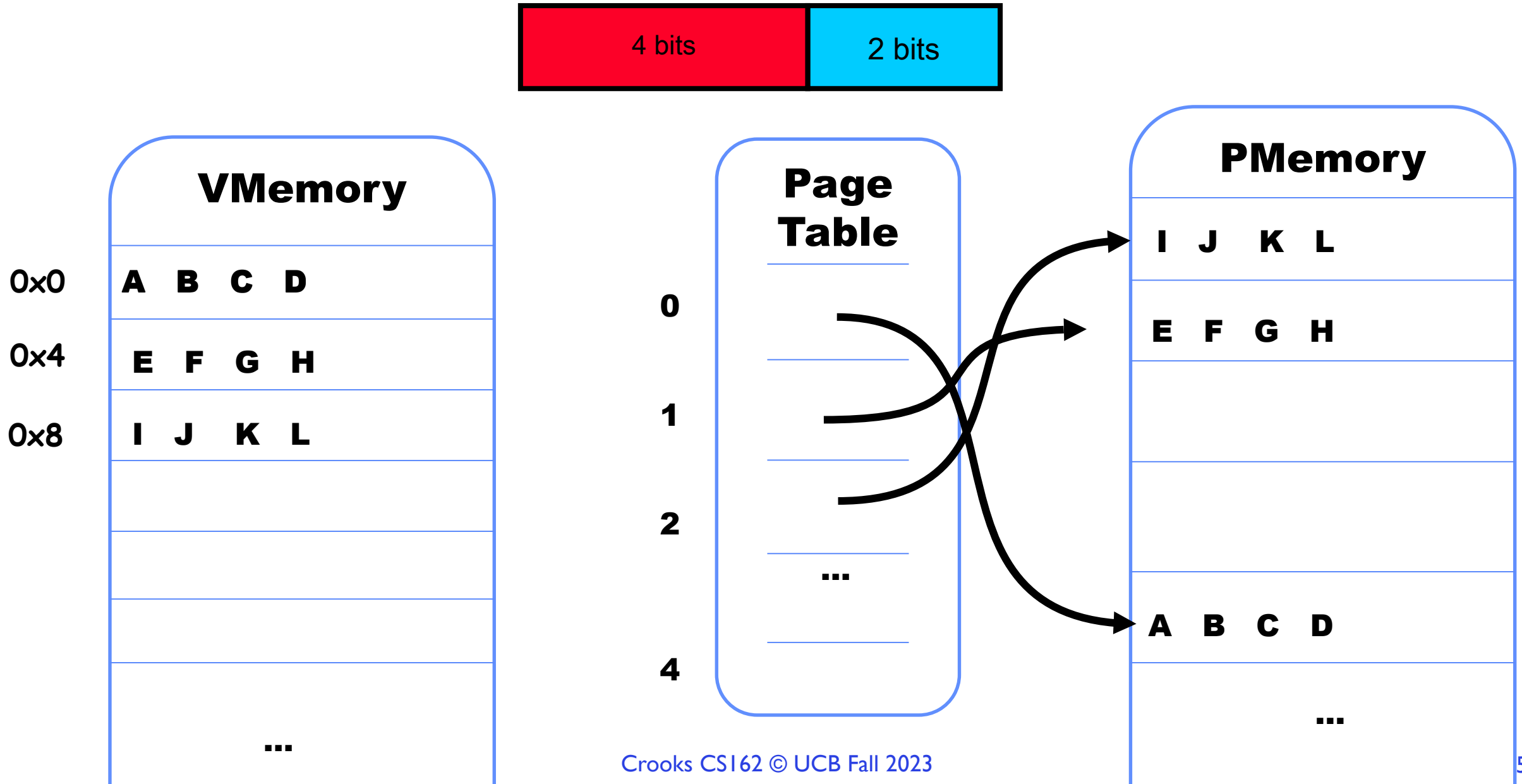
2 bits

How many virtual pages can we have?

6 bit addresses, 2 bit for offsets, 4 bits for VPN.

$2^4 = 16$ pages

# Example: A Mini Page Table
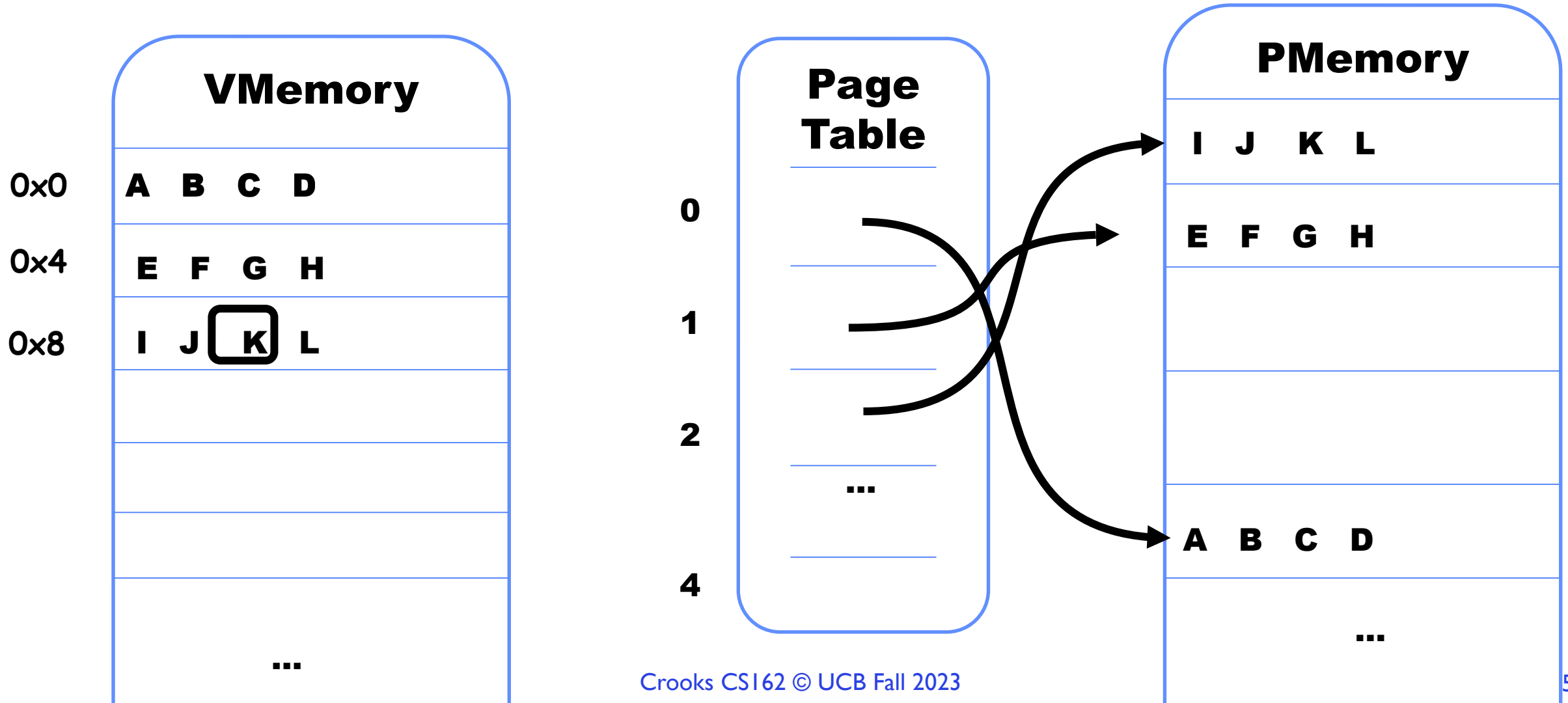
| 4 bits | 2 bits |
|---|---|

**VMemory**

0x0    A   B   C   D

0x4     E   F   G   H

0x8    I   J   K   L

...

**Page Table**

0

1

2

...

4

**PMemory**

I   J   K   L

E   F   G   H

A   B   C   D

...

# Example: A Mini Page Table

# Example: A Mini Page Table

0x9 =

| 0010 | 01 |
|------|-----|

### VMemory

| 0x0 | A | B | C | D |
| 0x4 | E | F | G | H |
| 0x8 | I | J | K | L |
| | ... | | | |

### Page Table

0

1

2

...

4

### PMemory

| I | J | K | L |
| E | F | G | H |
| | | | |
| A | B | C | D |
| | ... | | |

$$0x9 =$$

| 0010 | 01 |
|------|----|

**VMemory**

**Page Table**

**PMemory**

0x0    A    B

0x4    E    F

0x8    I    J

**Virtual Page Number:**

**0010 => 2.**

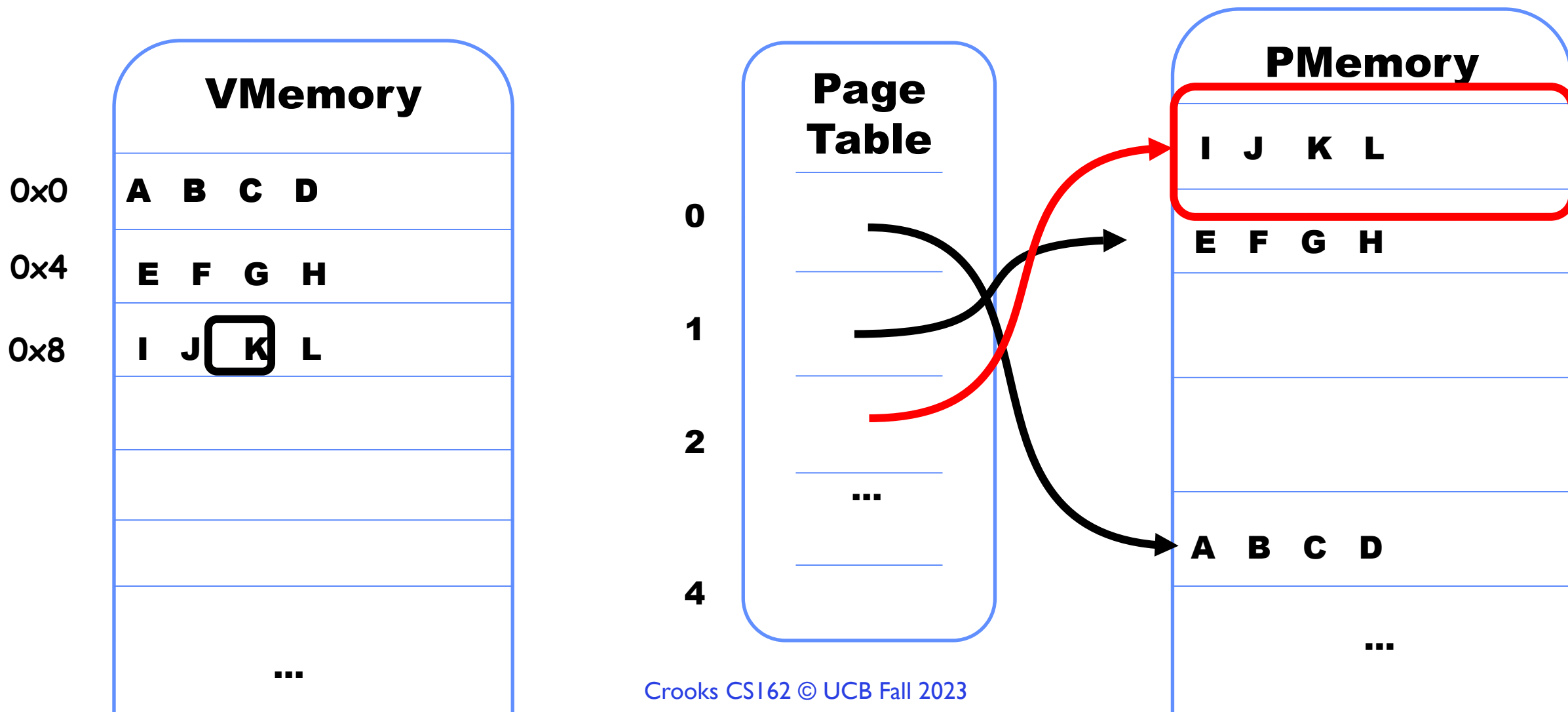**Access Index 2 of Page Table**

2

...

4

A    B    C    D

...

...

0x9 =

| 0010 | 01 |
|------|----|

**VMemory**

0x0  A  B  C  D

0x4  E  F  G  H

0x8  I  J  K  L

...

**Page Table**

0

1

2

...

4

**PMemory**

I  J  K  L

E  F  G  H

A  B  C  D

...

# Step 3: Extract Frame Offset

0x9 = | 0010 | 01 |

## VMemory

| | | | |
|---|---|---|---|
0x0 | A | B | C | D
0x4 | E | F | G | H
0x8 | I | J | K | L

...

## Page Table

0

1

2

...

4

## PMemory

| I | J | K | L |
|---|---|---|---|
| E | F | G | H |

| A | B | C | D |

...

# Step 3: Extract Frame Offset

0x9 =

| 0010 | 01 |
|------|-----|

**VMemory**

**PMemory**

**Page Table**

0x0 | A  B

0x4 | E  F

0x8 | I  J

**Offset:**

**01 => 1.**

**Access Byte 1 of Frame**

2

...

4

A   B   C   D

...

# Step 3: Extract Frame Offset

0x9 =

| 0010 | 01 |
|------|-----|

**VMemory**

0x0  A  B  C  D

0x4  E  F  G  H

0x8  I  J  K  L

...

**Page Table**

0

1

2

...

4

**PMemory**

I  J  K  L

E  F  G  H

A  B  C  D

...

# Step 4: Convert to Physical Address

0x9 =    | 0010 | 01 |

**VMemory**

| | | | |
|---|---|---|---|
| 0x0 | A B C D |
| 0x4 | E F G H |
| 0x8 | I J K L |

...

**Page Table**

0

1

**PMemory**

| | | | |
|---|---|---|---|
| I J K L |
| E F G H |

...

**Physical Page Number * Page Size
+ Offset
=
0 * 4 + 1 = 1**

# What is a page table entry? (32 bits)

| Page Frame Number (Physical Page Number) | Free (OS) | 0 | PS | D | A | PCD | PWT | U | W | P |
|---|---|---|---|---|---|---|---|---|---|---|
| 31-12 | 11-9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |

P:      Present (same as "valid" bit in other architectures)

W:      Writeable

U:      User accessible

PWT:    Page write transparent: external cache write-through

PCD:    Page cache disabled (page cannot be cached)

A:      Accessed: page has been accessed recently

D:      Dirty: page has been modified recently

PS:     Page Size

> **Size of page table entry:**
> **PFN (20 bits) + 12 bits for access control/caching**
>
> **4 bytes**

# The Great Power of the PTE

### Demand Paging

Keep only active pages in memory
Place others on disk and mark their PTEs invalid

### Copy-on-Write

UNIX fork gives *copy* of parent address space to child. Use combination of page sharing + marking pages as read-only

### Zero Fill On Demand

New data pages must carry no information

Mark PTEs as invalid; page fault on use gets zeroed page
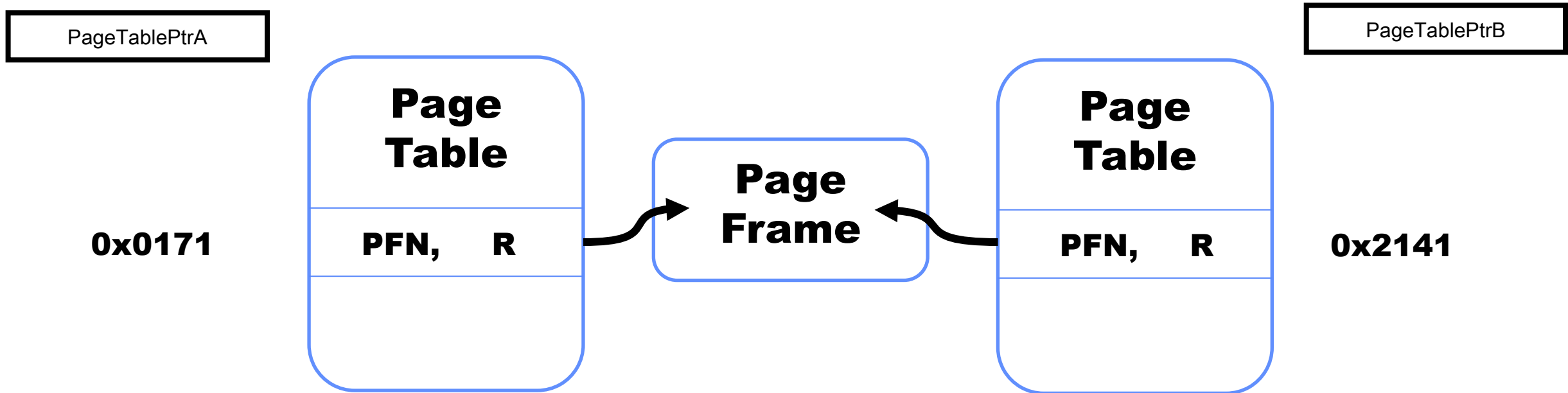
### Data Breakpoints

For debugger, mark instruction page as read-only. Will trigger page-fault when try to execute

# Paging & Sharing

Processes share a page by each mapping a page of their own virtual address space to the same frame

Use protection bits for fine-sharing

| PageTablePtrA |
| --- |

| PageTablePtrB |
| --- |

**0x0171**

**Page Table**

PFN,    R

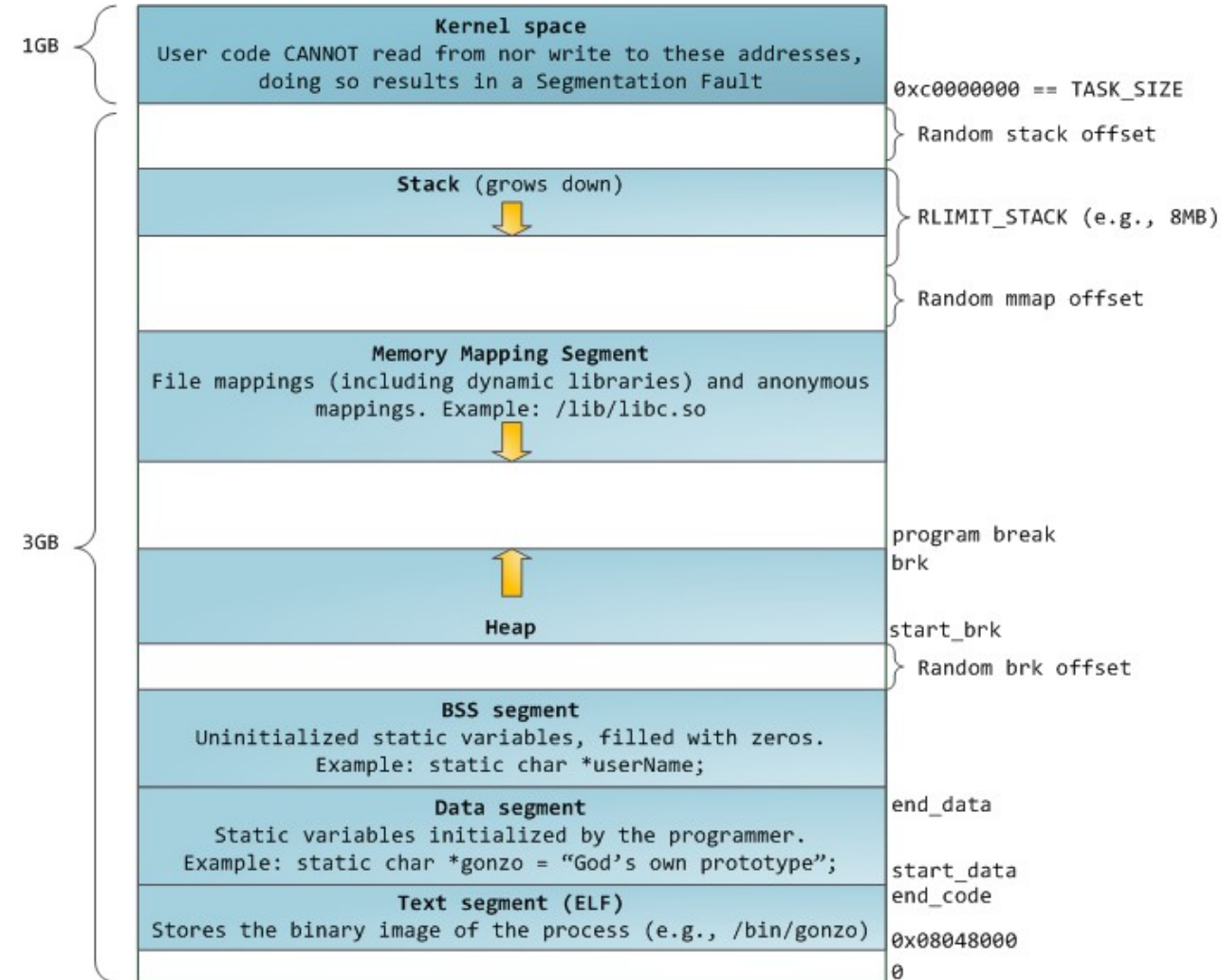**Page Frame**

**Page Table**

PFN,    R

**0x2141**

# Where is page sharing used ?

Kernel region of every process has the same page table entries

Different processes running same binary!
Do not need to duplicate code segments

Shared-memory segments between different processes

# Memory Layout for Linux 32-bit

# An aside: Meltdown

From the paper:

Meltdown is a novel attack that allows overcoming memory isolation completely by providing a simple way for any user process to read the entire kernel memory of the machine it executes on, including all physical memory mapped in the kernel region. Meltdown does not exploit any software vulnerability, i.e., it works on all major operating systems.

```
1. raise_exception();
2. // the line below is never reached
3. access(probe_array[data * 4096]);]
```

# Are we done?

How big can a page table get on x86 (32 bits)?

4KB page => 2^12
2^32/2^12 => 2^20 pages
2^20 * 4 bytes = 4 MB (approx.)
That's a lot per process!!

How big can a page table get on x86 (64 bits)?

4KB page => 2^12
2^64/2^12 => 2^52 pages
2^20 * 8 bytes = 36 petabytes (approx.)
That's a lot per process!!

# Limitations of paging

### Space overhead
With a 64-bit address space, size of page table can be huge

### Time overhead
Accessing data now requires two memory accesses must also access page table, to find mapped frame

### Internal Fragmentation
4KB pages

# The Secret to the Whole of CS

Batching

Caching

Indirection

Specialised Hardware

# Sparsity

Address space is sparse, i.e. has holes that are not mapped to physical memory
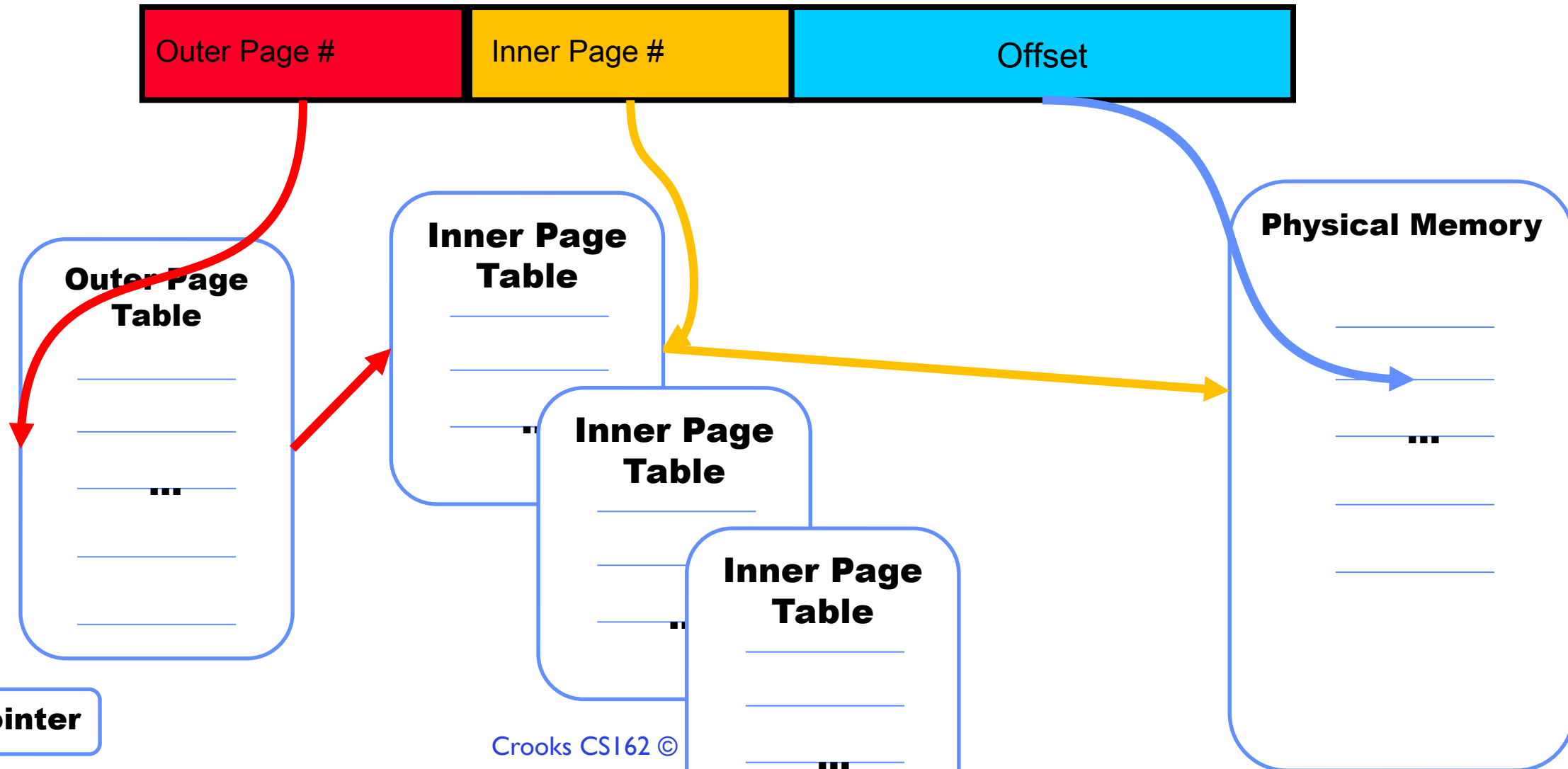
Most this space is taken up by page tables mapped to nothing

Process has access to full 2^64 bytes (virtually)

Physically, that would be 17,179,869,184 gigabytes

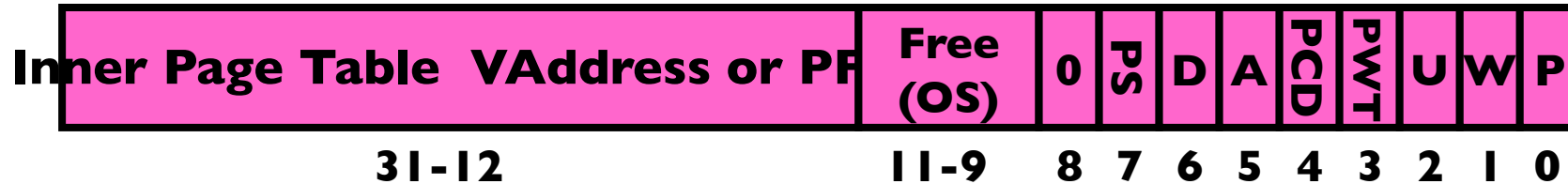# Paging the page table: 2-level paging

## Tree of Page Tables



| Outer Page # | Inner Page # | Offset |

**Outer Page Table**

**Inner Page Table**

**Inner Page Table**

**Inner Page Table**

**Physical Memory**

**PageTablePointer**

Crooks CS162 ©

# V2: What is a page table entry? (32 bits)

| Inner Page Table VAddress or PP | Free (OS) | 0 | PS | D | A | PCD | PWT | U | W | P |
|---|---|---|---|---|---|---|---|---|---|---|
| 31-12 | 11-9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |

P:      Present (same as "valid" bit in other architectures)

W:      Writeable

U:      User accessible

PWT:    Page write transparent: external cache write-through

PCD:    Page cache disabled (page cannot be cached)

A:      Accessed: page has been accessed recently

D:      Dirty: page has been modified recently

PS:     Page Size

# Paging the page table: 2-level paging

**Tree of Page Tables**

| Outer Page # | Inner Page # | Offset |
|---|---|---|

| Number of top-level pages | Ensure that fits on a single page | Defines size of a page |

# Paging the page table: 2-level paging

## Tree of Page Tables

| Outer Page # | Inner Page # | Offset |
|:---:|:---:|:---:|

| 10 bits | 10 bits | 4 KB<br>12 bits |

Want to make sure that
inner page table fits in a
page!
$2^{12}/2^2 = 2^{10}$
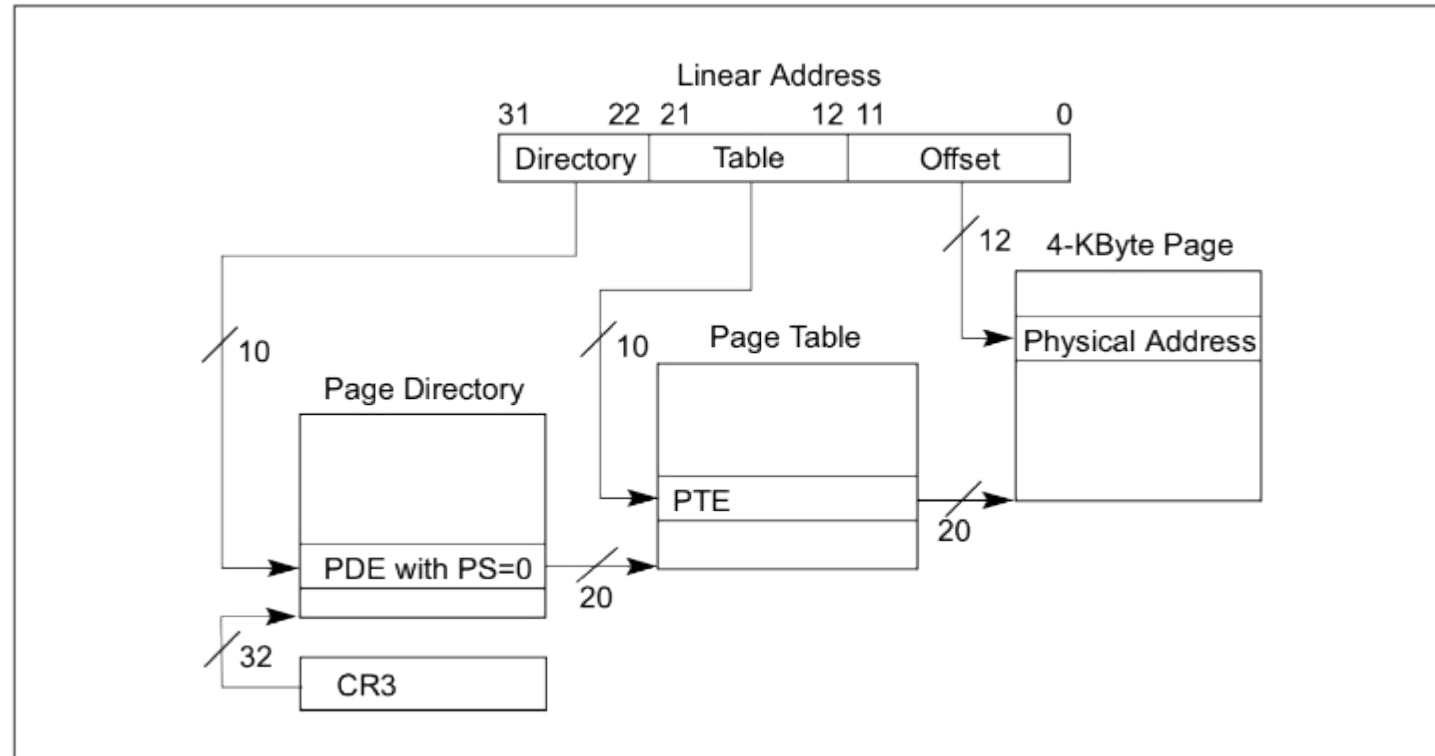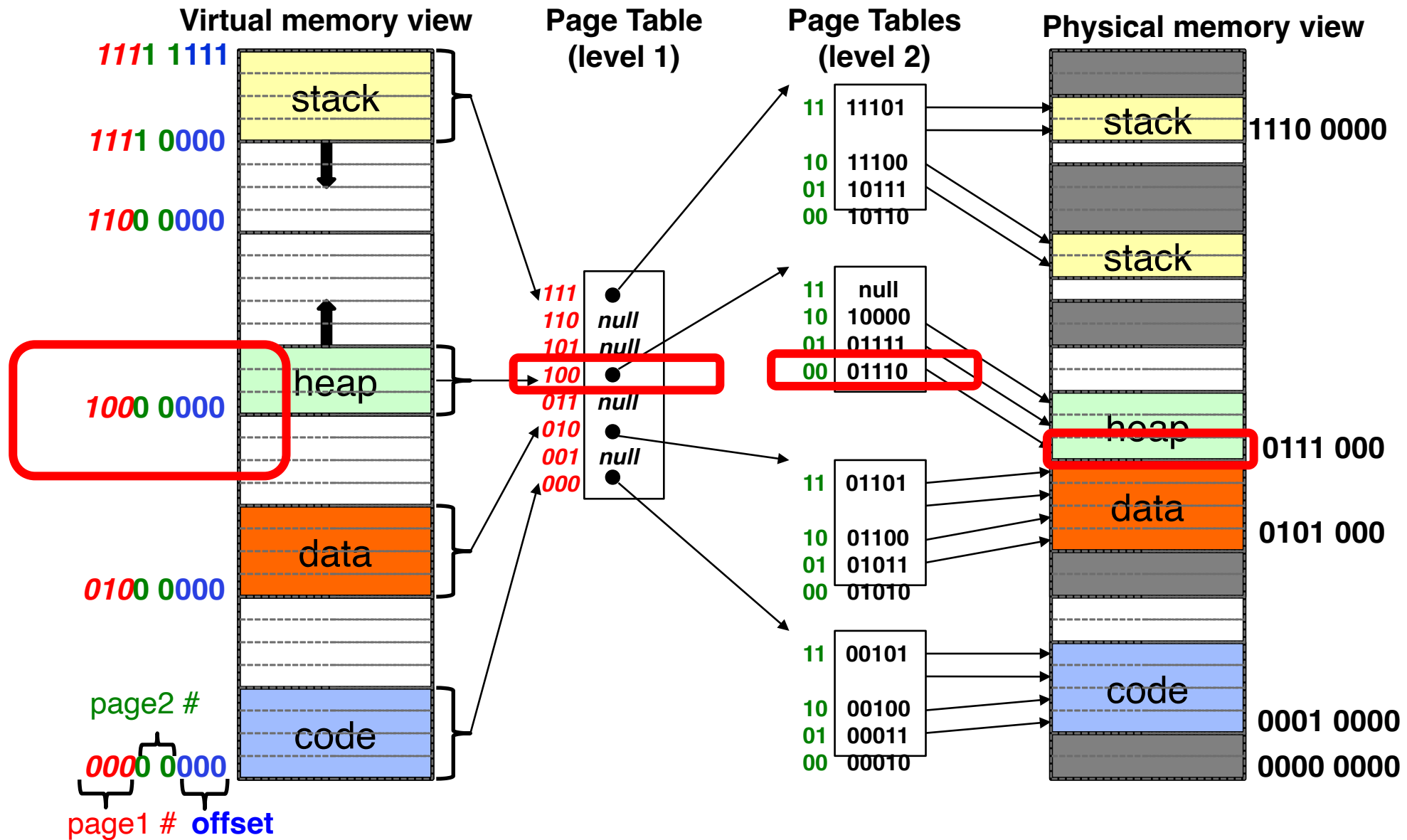
# Example: x86 classic 32-bit address translation



Figure 4-2. Linear-Address Translation to a 4-KByte Page using 32-Bit Paging
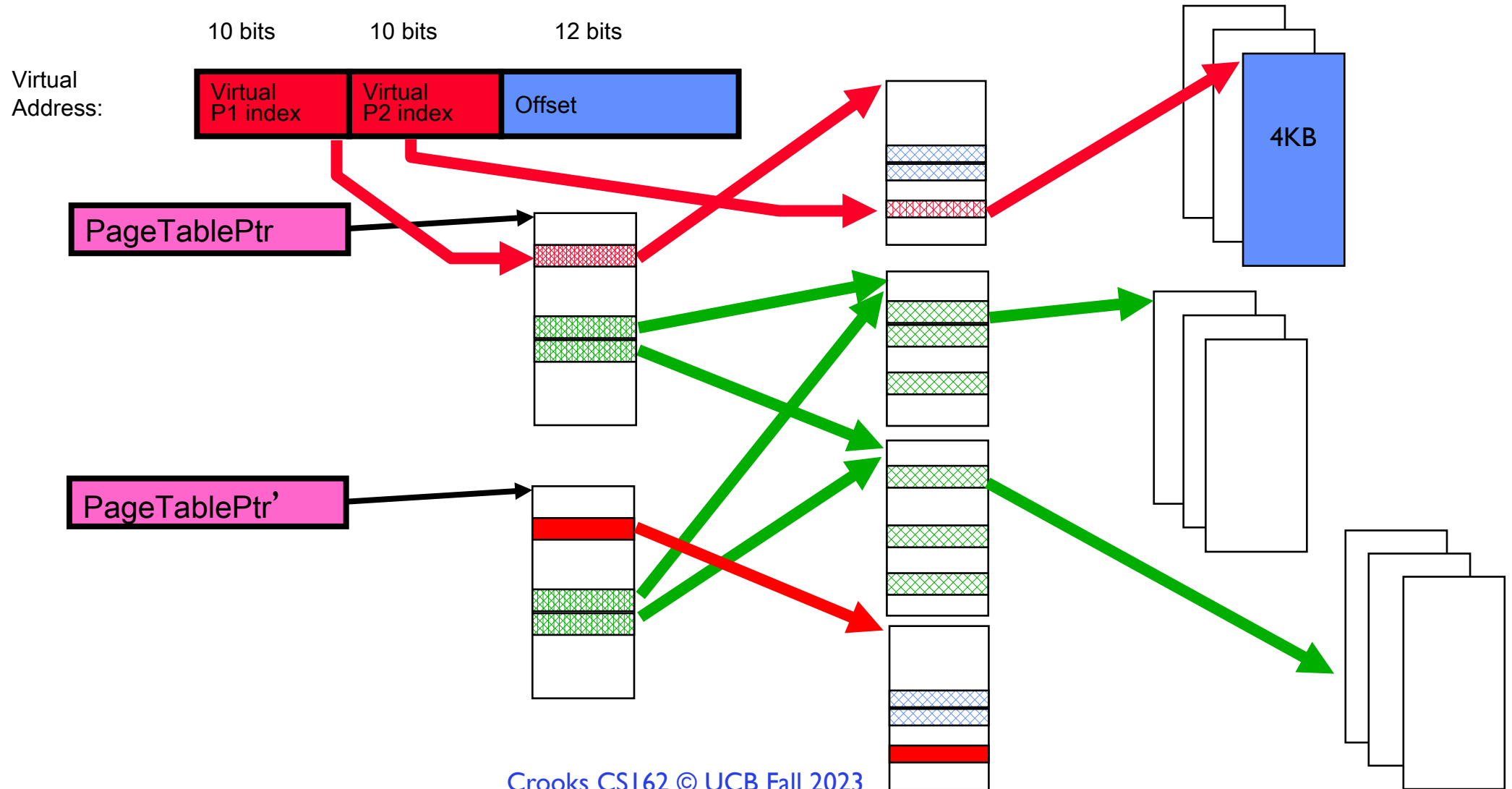
Top-level page-table: Page Directory

Inner page-table:  Page Directory Entries

# Example Address Space View



**Virtual memory view**

1111 1111
stack
1111 0000

1100 0000

1000 0000
heap

0100 0000
data

0000 0000

page2 #
page1 # offset

**Page Table (level 1)**

111
110 null
101 null
100
011 null
010
001 null
000

**Page Tables (level 2)**

11 11101
10 11100
01 10111
00 10110

11 null
10 10000
01 01111
00 01110

11 01101
10 01100
01 01011
00 01010

11 00101
10 00100
01 00011
00 00010

**Physical memory view**

stack 1110 0000

stack

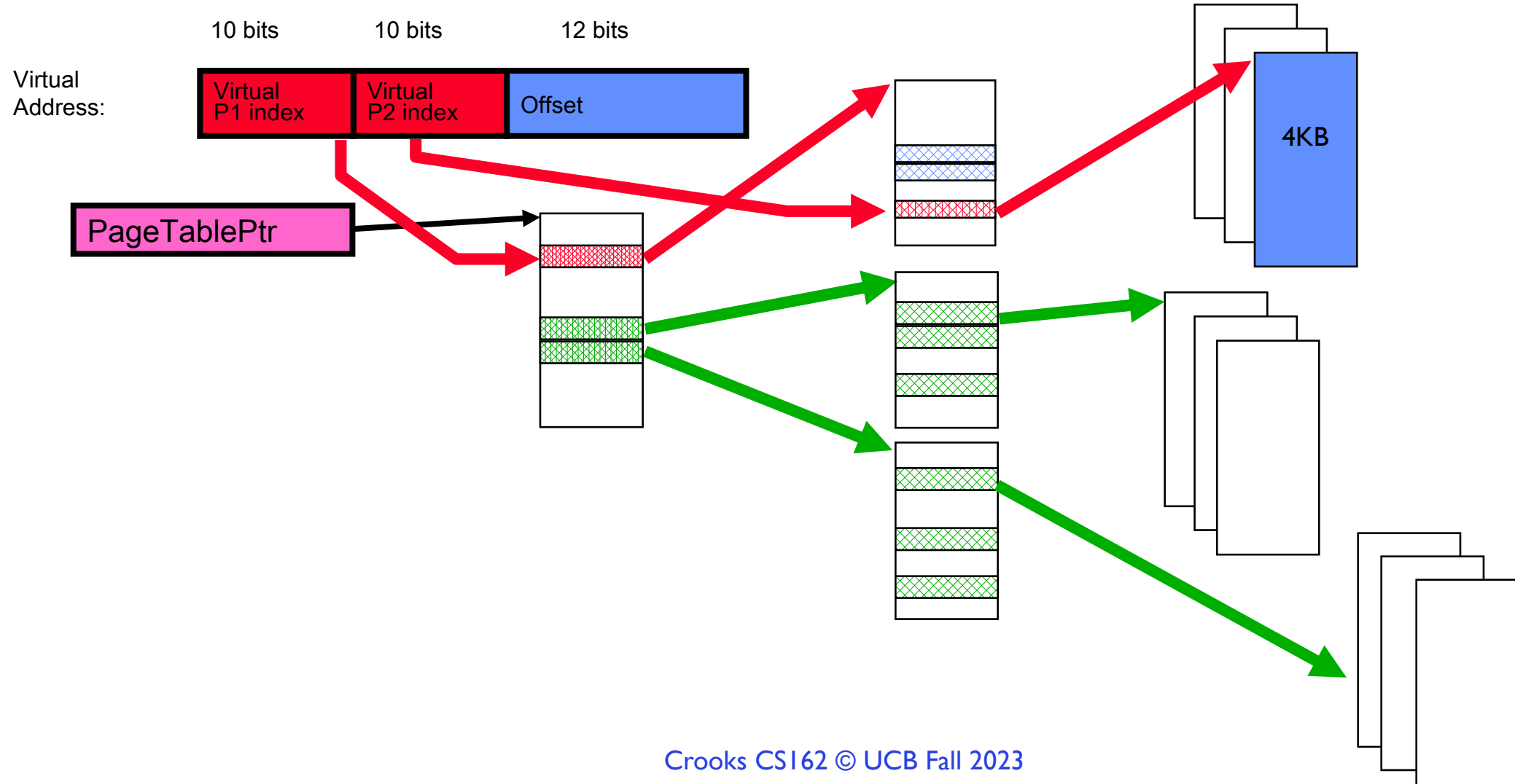heap 0111 000

data 0101 000

code 0001 0000
0000 0000

# Sharing with multilevel page tables

Entire regions of the address space can be efficiently shared
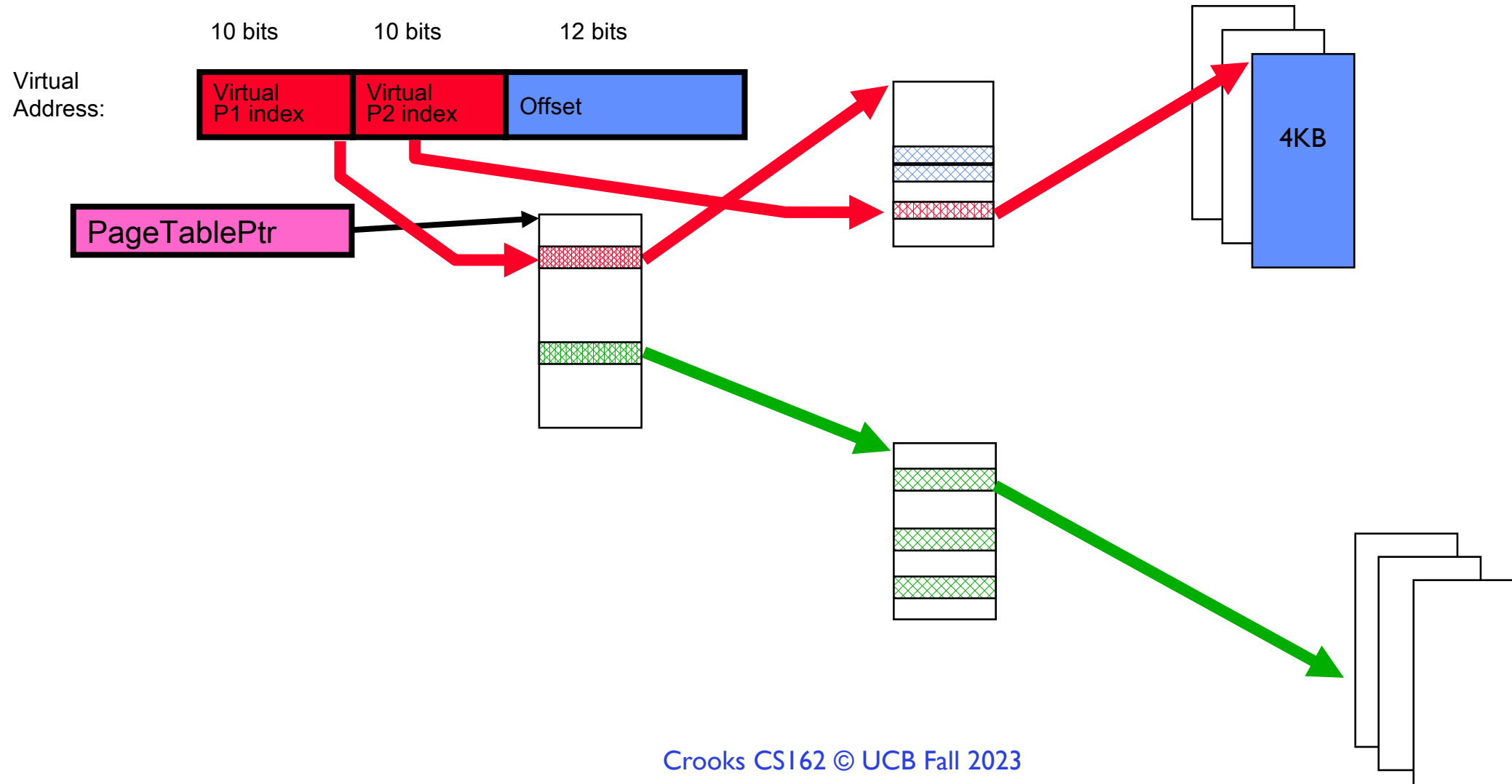
# Marking entire regions as invalid!

If region of address space unused, can mark entire inner region as invalid

# Marking entire regions as invalid!

If region of address space unused, can mark entire inner region as invalid

# Has this helped?

Assuming 10/10/12 split:

Size of Page Table

Outer: ($2^{10}$ * 4 bytes) +
Inner: $2^{10}$ * ($2^{10}$ * 4 bytes)

Overhead of indirection! BUT Marking inner pages as invalid helps when address spaces are sparse
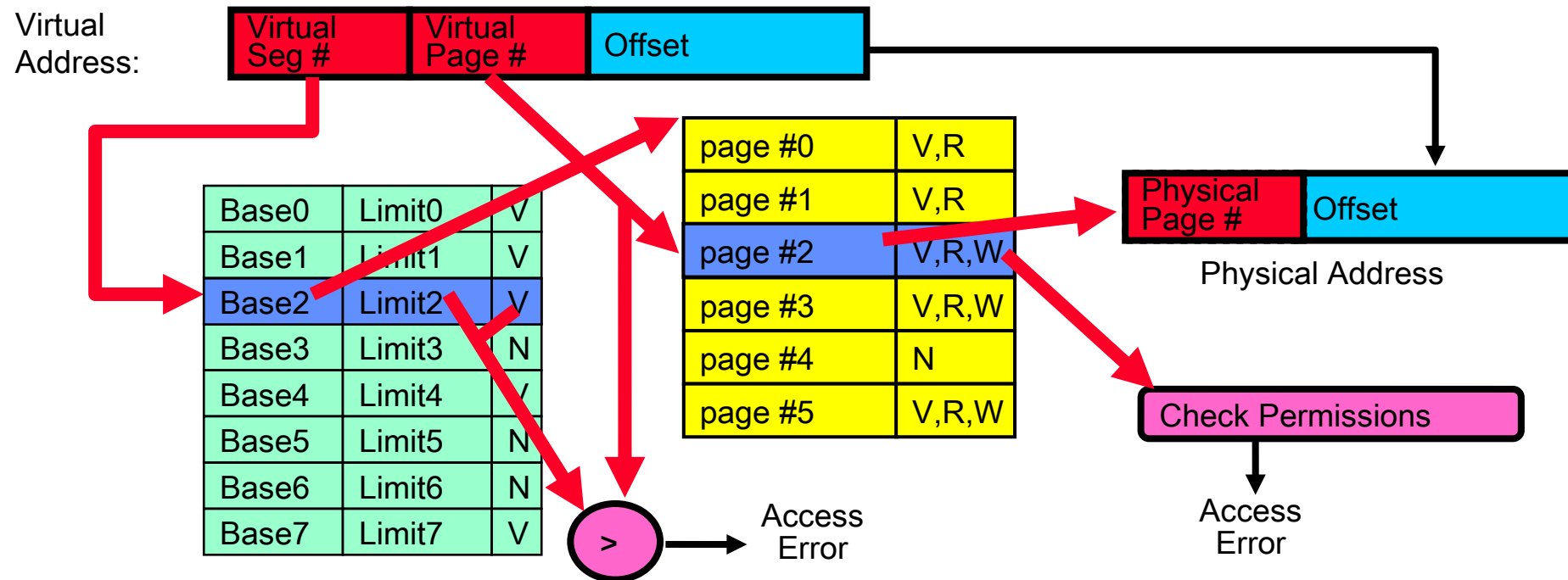
Downside: now have to do

two memory accesses for translation
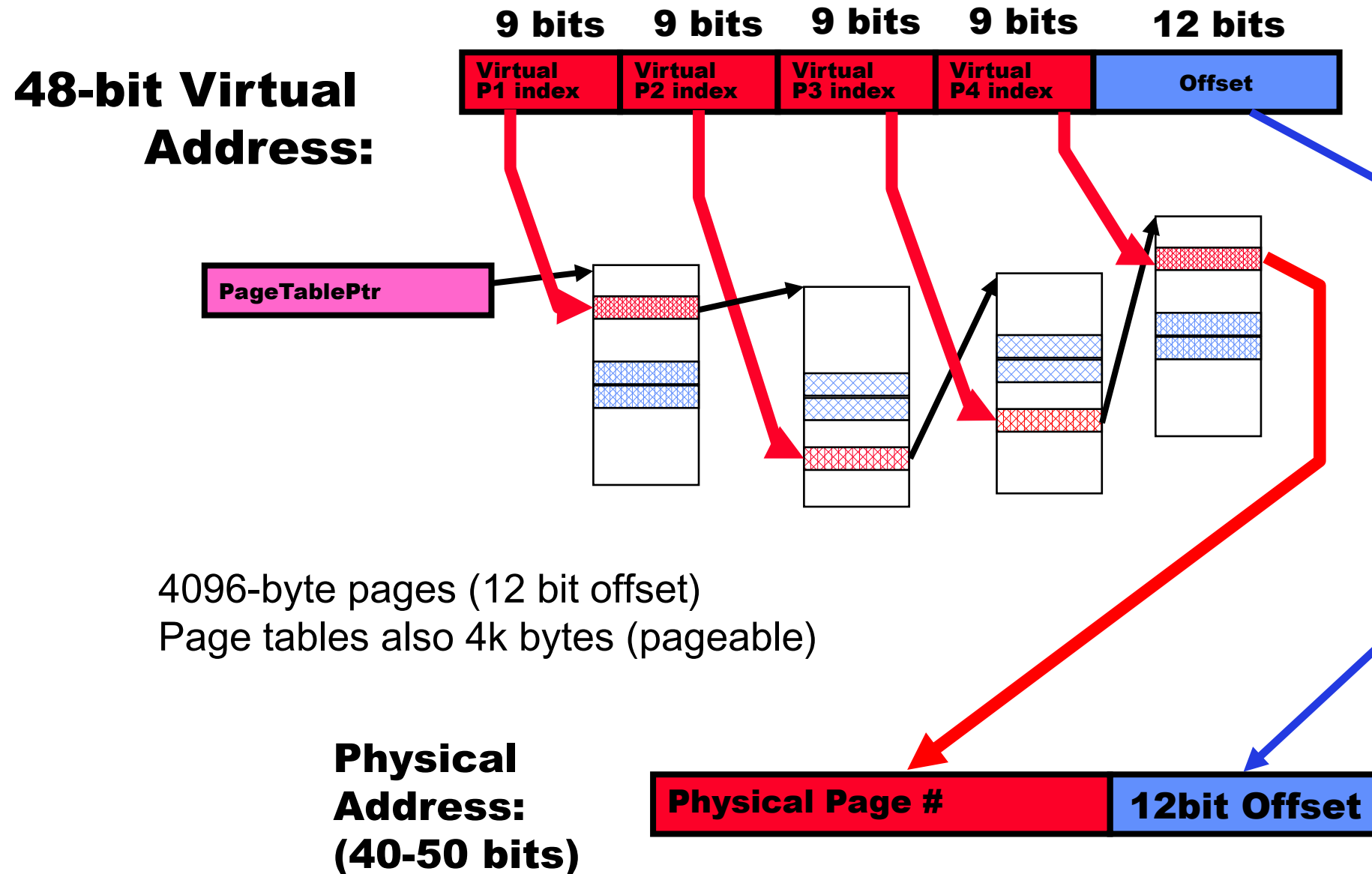
# Paged Segmentation

Use segments for top level. Paging within each segment.

Used in x86 (32 bit).
Code Segment, Data Segment, etc.

# X86 64 bits has a four-level page table!

**48-bit Virtual Address:**

| 9 bits | 9 bits | 9 bits | 9 bits | 12 bits |
|---|---|---|---|---|
| Virtual P1 index | Virtual P2 index | Virtual P3 index | Virtual P4 index | Offset |

PageTablePtr

4096-byte pages (12 bit offset)
Page tables also 4k bytes (pageable)

**Physical Address: (40-50 bits)**

| Physical Page # | 12bit Offset |
|---|---|

# Inverted Page Table

A single page table that
has an entry for each physical page of the
system

Each entry contains process ID + which
virtual page maps to physical page
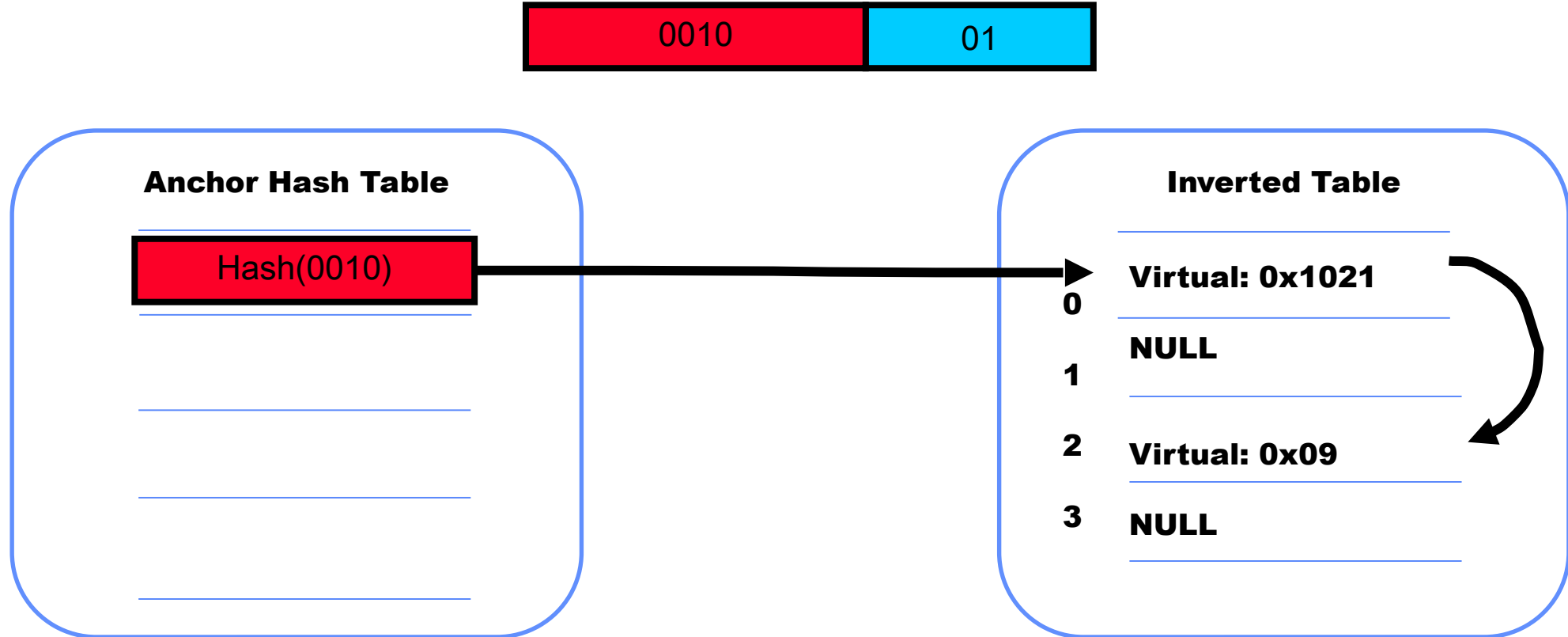
Physical memory much smaller than virtual
memory

Size proportional to size of physical
memory

**Inverted Table**

| | |
|---|---|
| **0** | **Virtual: 0x1021** |
| **1** | **NULL** |
| **2** | **Virtual: 0x0123** |
| **3** | **NULL** |

# Inverted Page Table

Don't we have it backwards?

Add a hash table. Virtual memory can only map to specific physical frames

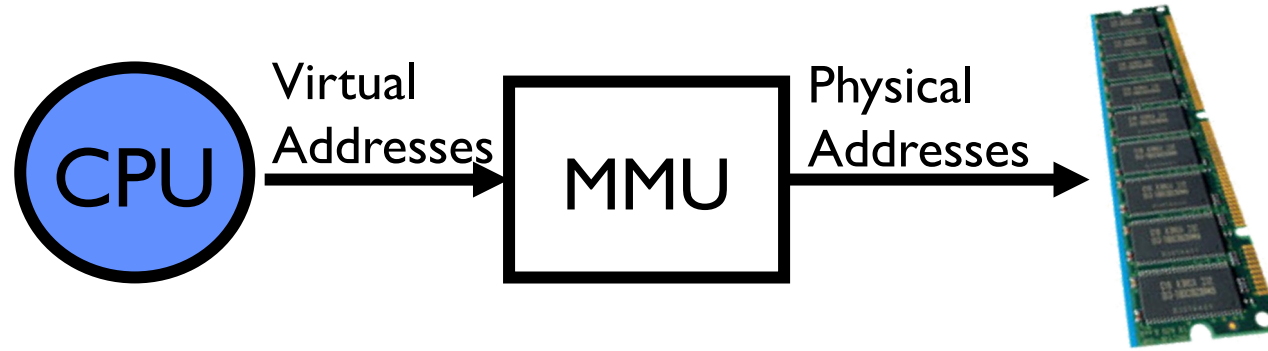| 0010 | 01 |
|------|-----|

**Anchor Hash Table**

Hash(0010)

**Inverted Table**

0    **Virtual: 0x1021**

1    **NULL**

2    **Virtual: 0x09**

3    **NULL**

# Address Translation Comparison

|  | Advantages | Disadvantages |
|---|---|---|
| Simple Segmentation | Fast context switching (segment map maintained by CPU) | External fragmentation |
| Paging (Single-Level) | No external fragmentation<br>Fast and easy allocation | Large table size (~ virtual memory)<br>Internal fragmentation |
| Paged Segmentation | Table size ~ # of pages in virtual memory<br>Fast and easy allocation | Multiple memory references per page access |
| Multi-Level Paging | | |
| Inverted Page Table | Table size ~ # of pages in physical memory | Hash function more complex<br>No cache locality of page table |

# How is the Translation Accomplished?



MMU must translate virtual address to physical address on every instruction fetch, load or store

What does the MMU need to do to translate an address?
Read, check, and update PTE
(set accessed bit/dirty bit on write)

# How can we speedup translation?

MMU must make at least 2 memory reads to walk page table. Slow!

Use specialized hardware to
cache virtual-physical memory translations!

Introducing the Translation Lookaside Buffer (TLB)

# Recall: CS61c Caching Concept

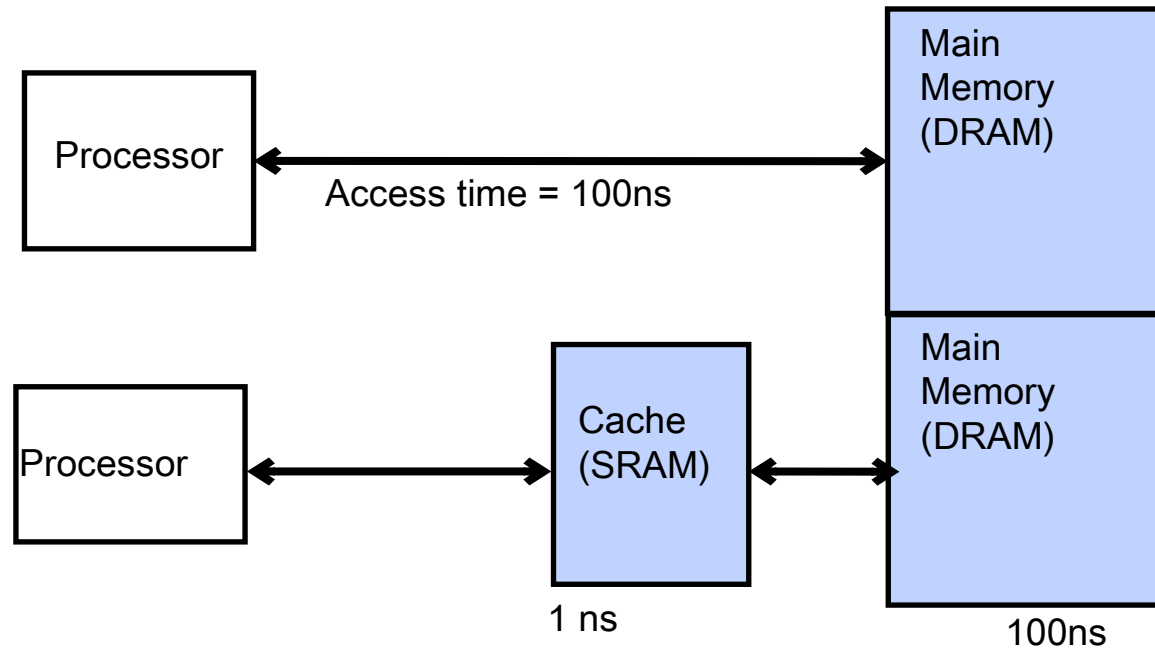Cache: a repository for copies that can be accessed more quickly than the original

Only good if:
Frequent case frequent enough and
Infrequent case not too expensive

Important measure
Average Access time =
(Hit Rate x Hit Time) + (Miss Rate x Miss Time)

# Recall: In Machine Structures (eg. 61C) ...

Caching is the key to memory system performance

Processor ←→ Main Memory (DRAM)

Access time = 100ns

Processor ←→ Cache (SRAM) ←→ Main Memory (DRAM)

1 ns          100ns

Average Memory Access Time (AMAT)

= (Hit Rate x HitTime) + (Miss Rate x MissTime)

Where HitRate + MissRate = 1

HitRate = 90% => AMAT = (0.9 x 1) + (0.1 x 101)=11 ns

HitRate = 99% => AMAT = (0.99 x 1) + (0.01 x 101)=2.01 ns

$MissTime_{L1}$ includes

$HitTime_{L1}+MissPenalty_{L1} \equiv HitTime_{L1} +AMAT_{L2}$

# Why Does Caching Help? Locality!

## Temporal Locality (Locality in Time):

Keep recently accessed data items closer to processor

## Spatial Locality (Locality in Space):

Move contiguous blocks to the upper levels
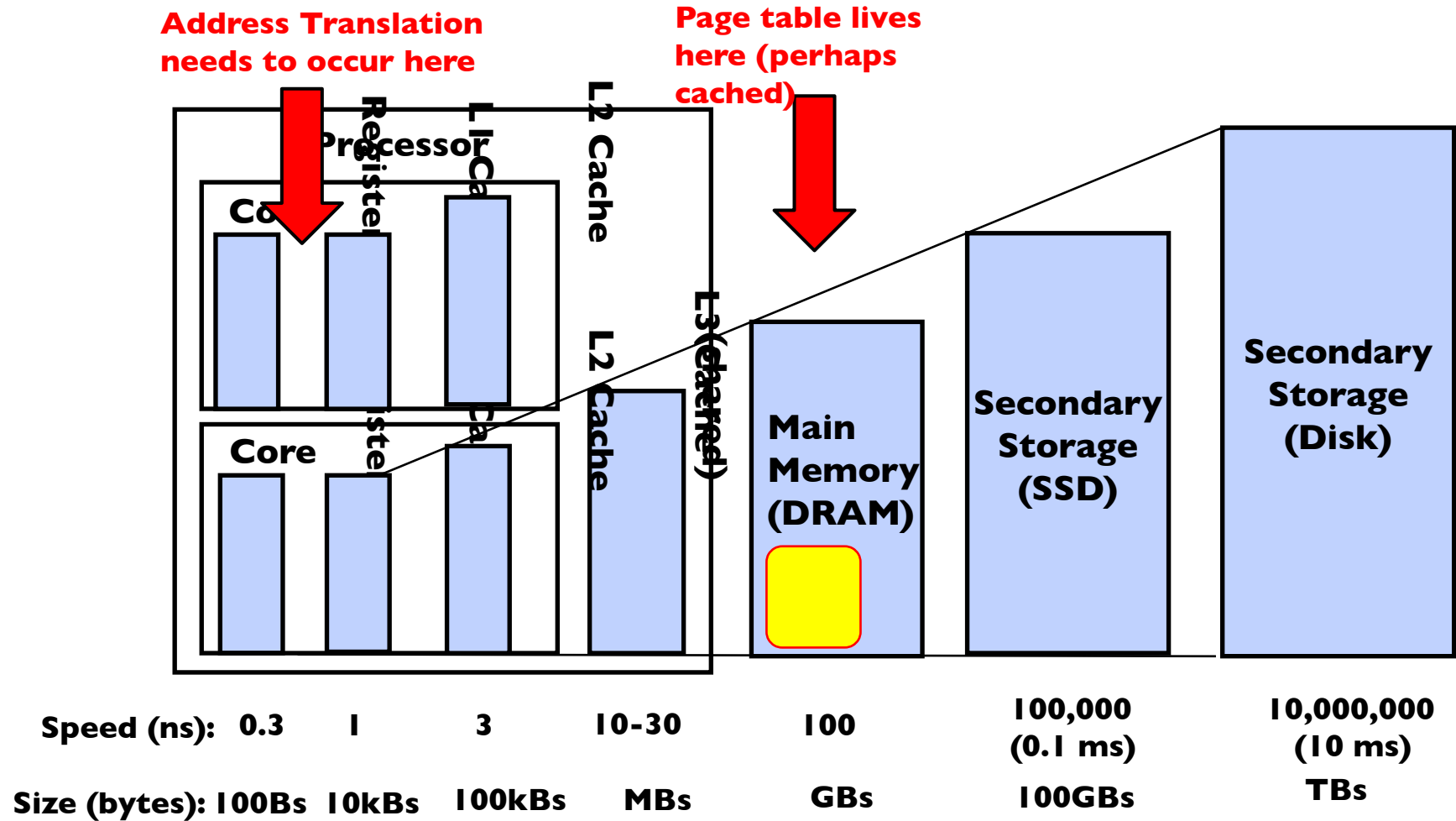
# Recall: Memory Hierarchy

Take advantage of the principle of locality to:

1) Present the illusion of having as much memory as in the cheapest technology

2) Provide average speed similar to that offered by the fastest technology
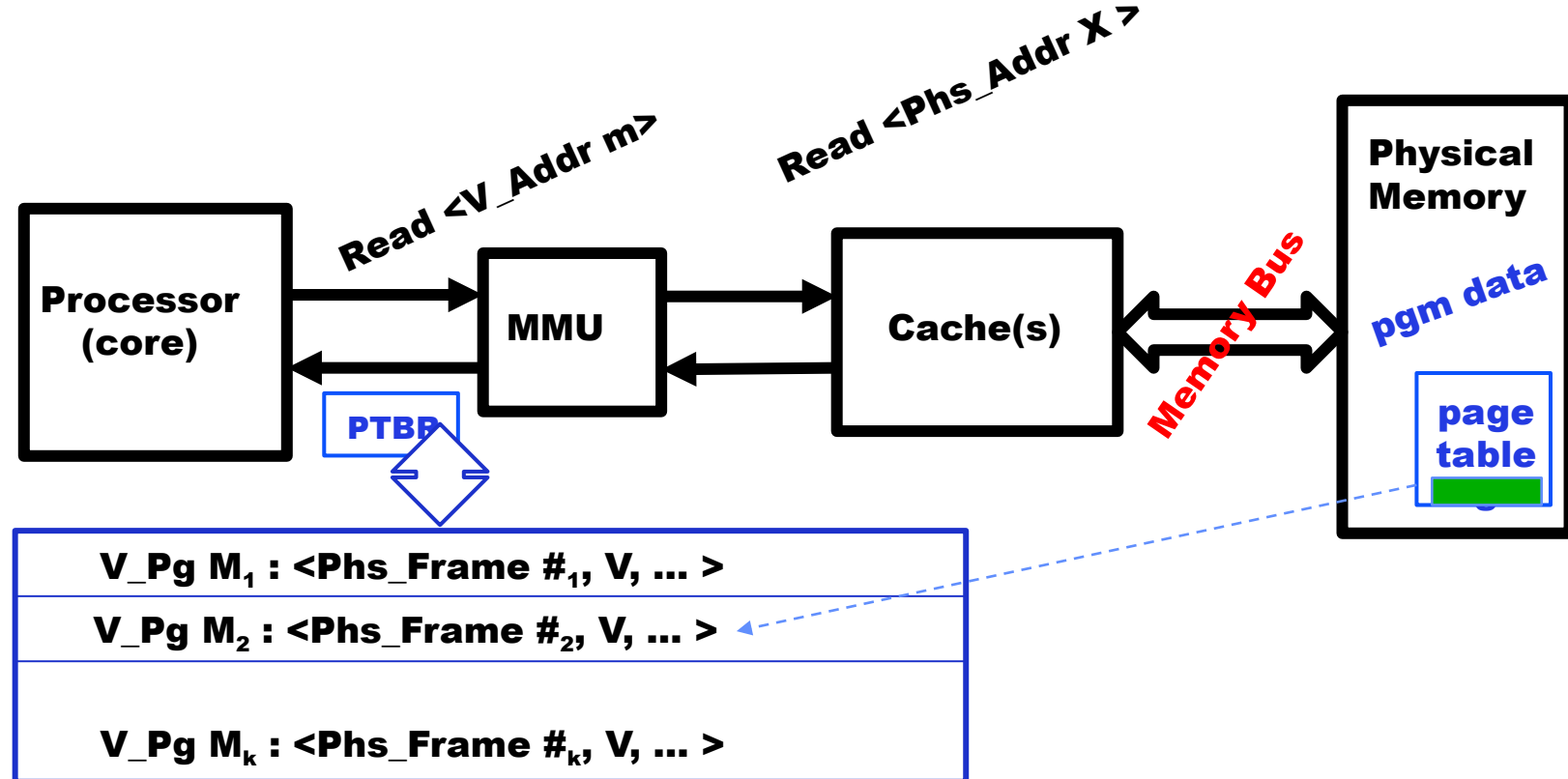
Recall: fast but small/expensive. Slow but large!

# Recall: Memory Hierarchy



| Speed (ns): | 0.3 | 1 | 3 | 10-30 | 100 | 100,000 (0.1 ms) | 10,000,000 (10 ms) |
|---|---|---|---|---|---|---|---|
| Size (bytes): | 100Bs | 10kBs | 100kBs | MBs | GBs | 100GBs | TBs |

# How do we make Address Translation Fast?

Cache results of recent translations !

Cache Page Table Entries using Virtual Page # as the key



V_Pg $M_1$ : <Phs_Frame #$_1$, V, ... >

V_Pg $M_2$ : <Phs_Frame #$_2$, V, ... >

V_Pg $M_k$ : <Phs_Frame #$_k$, V, ... >
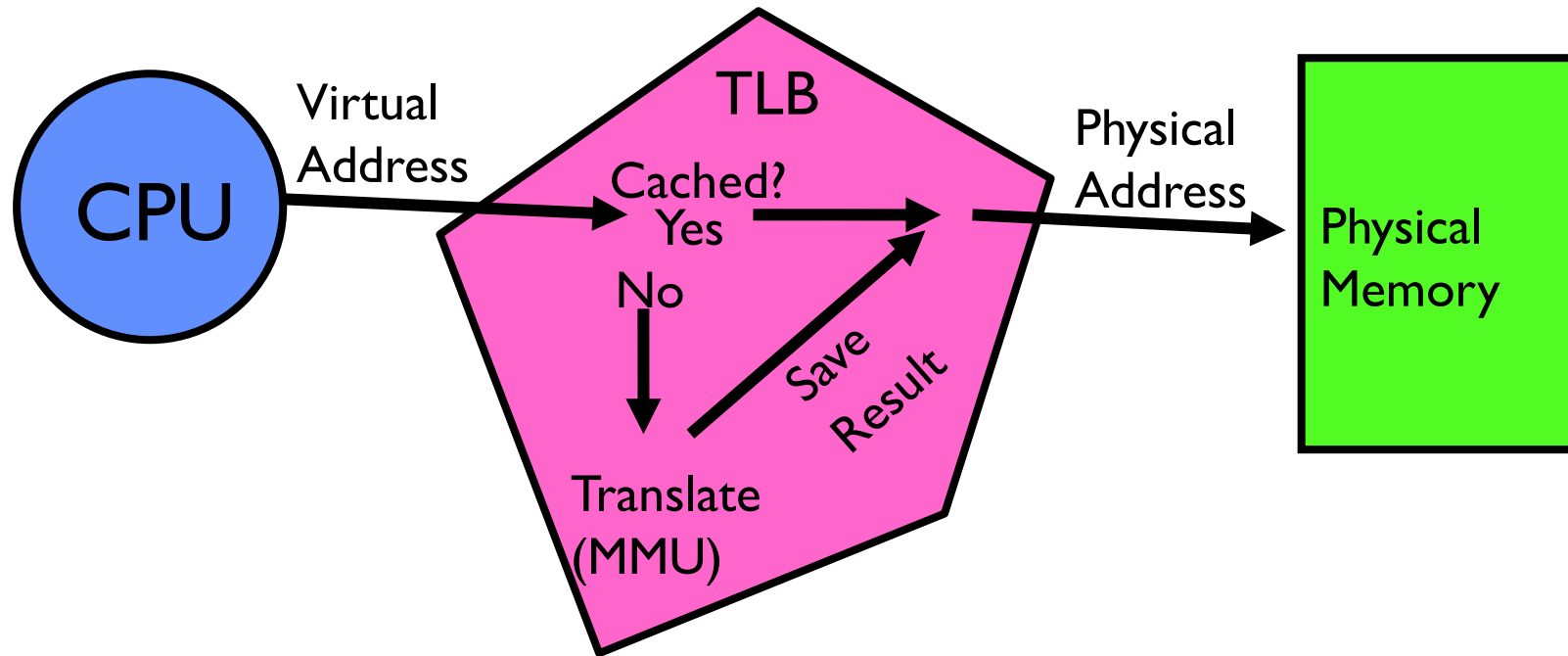
# Translation Look-Aside Buffer

Record recent Virtual Page # to Physical Frame # translation

If present, have the physical address without reading any of the page tables !!!

Caches the end-to-end result

# Caching Applied to Address Translation



Does page locality exist?

Instruction accesses spend a lot of time on the same page
(since accesses sequential)
Stack accesses have definite locality of reference