# Lecture 12:
# Memory Management

Address space, Address translation, Segments

# Abstraction is our Business

- What I have
  - A single (or a finite number) of CPUs
  - Many programs I would like to run
- What I want: a Thread
  - Each program has full control of one or more CPUs

# Abstraction is our Business

- What I have

    - A certain amount of physical memory

    - Multiple programs I would like to run

        - together, they may need more than the available physical memory

- What I want: an Address Space

    - Each program has as much memory as the machine's architecture will allow to name

    - All for itself

# Address Space

- Set of all names used to identify and manipulate unique instances of a given resource
  - memory locations (determined by the size of the machine's word)
    - for 32-bit-register machine, the address space goes from 0x00000000 to 0xFFFFFFFF
  - memory locations (determined by the number of memory banks mounted on the machine)
  - phone numbers (XXX) (YYY-YYYY)
  - colors: R (8 bits) + G (8 bits) + B (8 bits)

4

# Virtual Address Space: An Abstraction for Memory

- Virtual addresses start at 0

- Heap and stack can be placed far away from each other, so they can nicely grow

- Addresses are all contiguous

- Size is independent of physical memory on the machine

0xFFFFFFFF

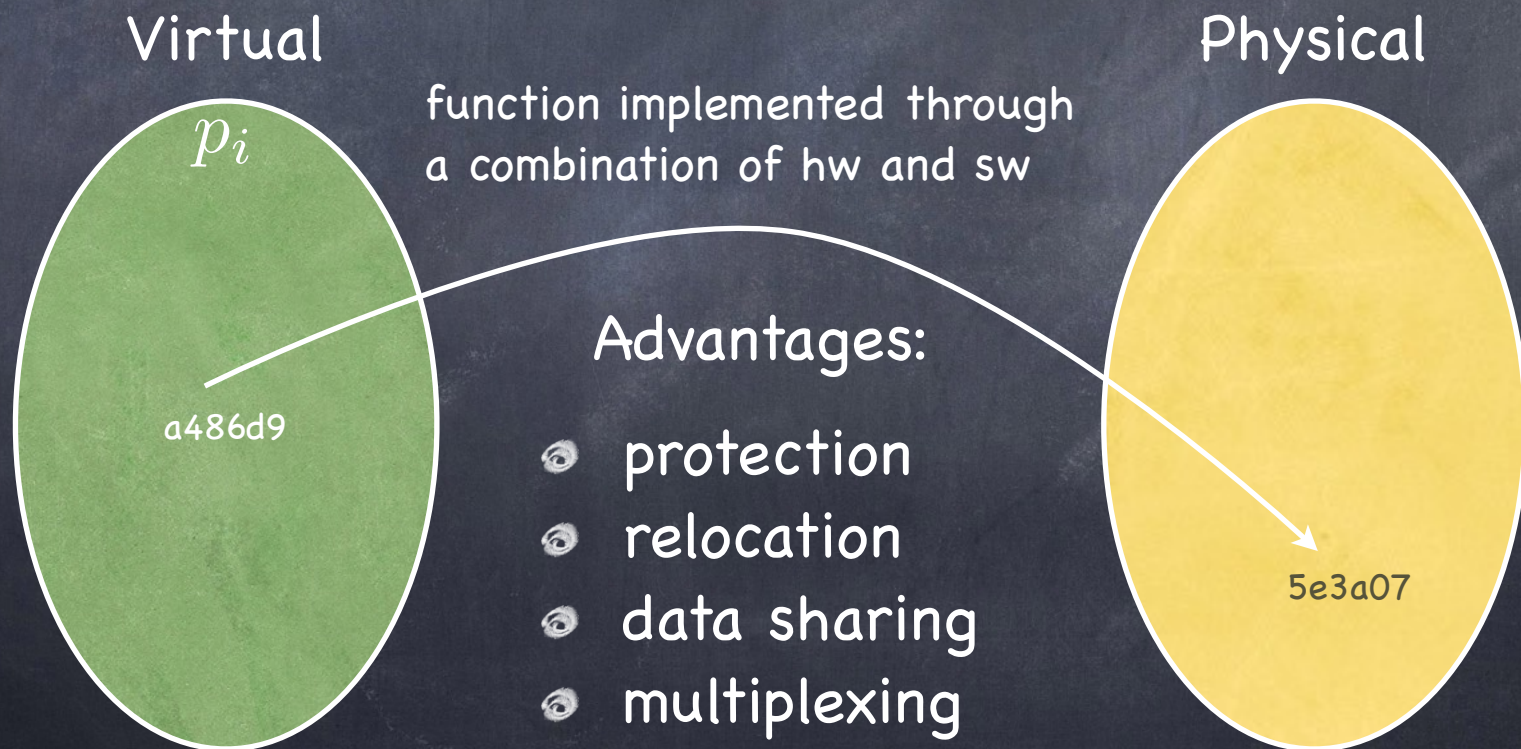| Stack |
| free |
| Heap |
| Data |
| Text |

0x00000000

# Physical Address Space: How memory actually looks

- Processes loaded in memory at some memory location

    - virtual address 0 is not loaded at physical address 0

- Multiple processes may be loaded in memory at the same time, and yet...

- ...physical memory may be too small to hold even a single virtual address space in its entirety
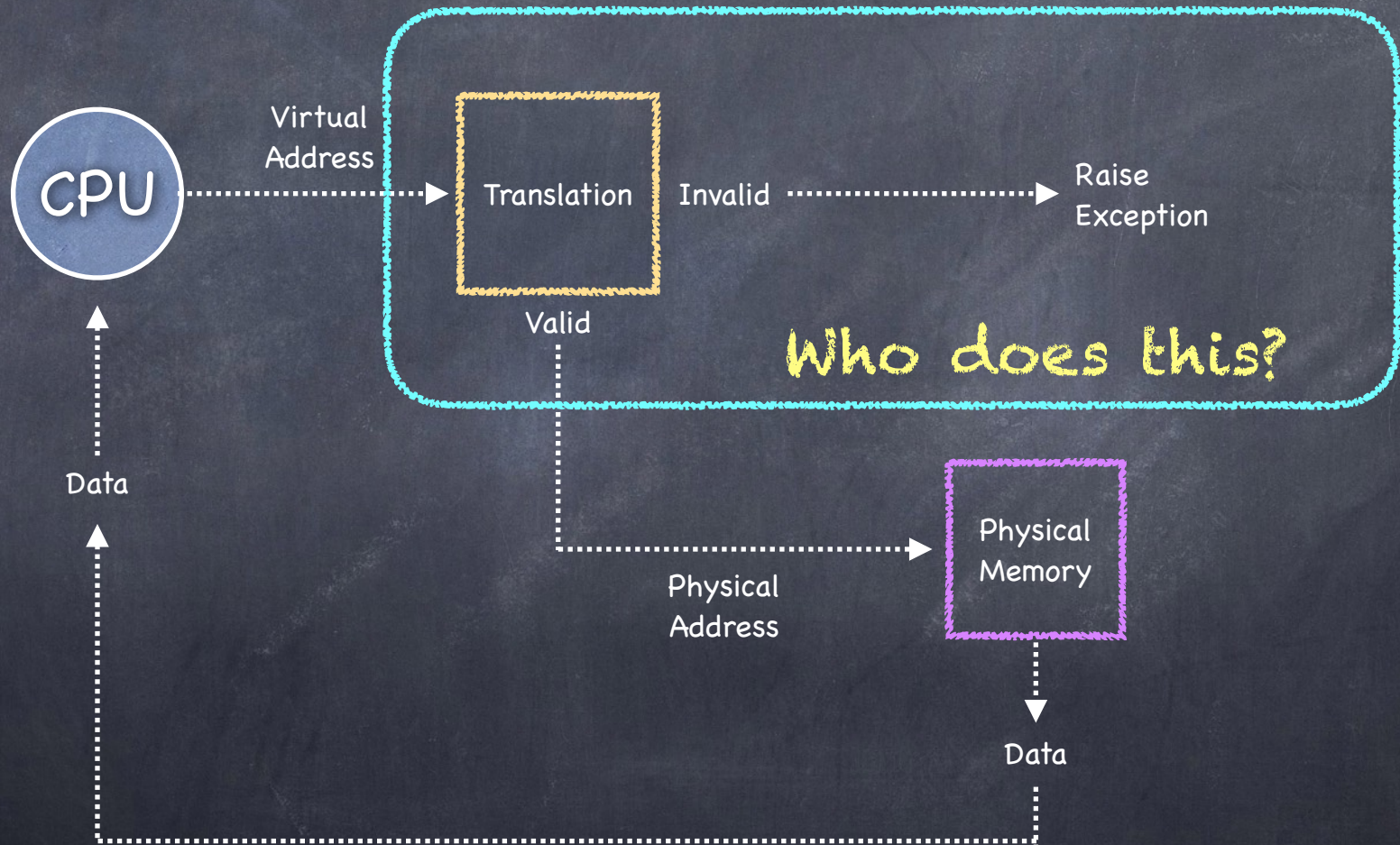
    - 64-bit, anyone?

| 0 | |
|---|---|
| OS | |
| **free** | |
| Process 2 | code, data, etc |
| **free** | |
| **free** | |
| Process 1 | code, data, etc |
| Process 3 | code, data, etc |
| **free** | |

512K

# Address Translation

- A function that maps $\langle pid, virtual\ address \rangle$ into a corresponding *physical address*

Virtual

Physical

$p_i$

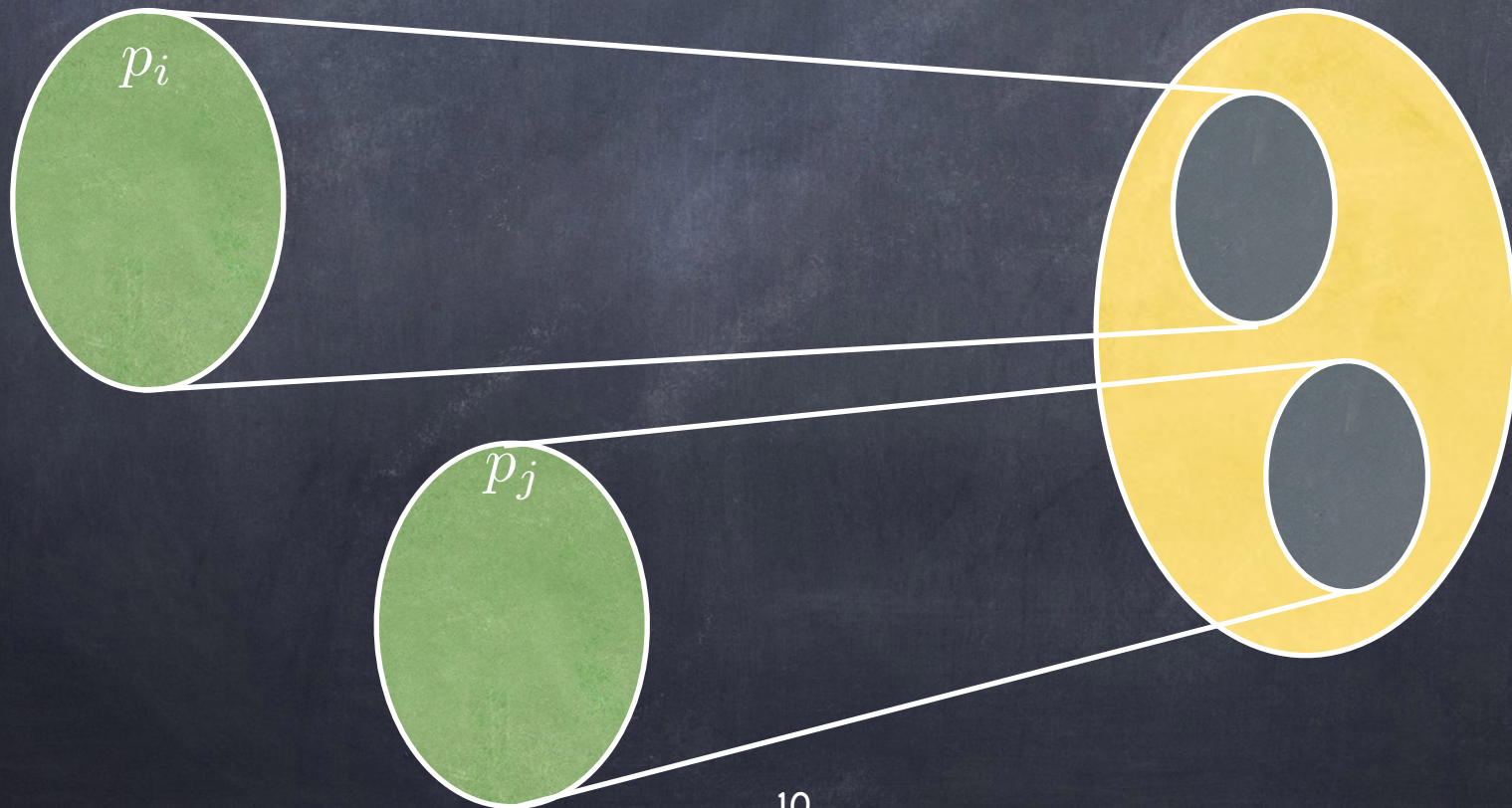function implemented through a combination of hw and sw

a486d9

Advantages:

- protection
- relocation
- data sharing
- multiplexing

5e3a07

7

# Address Translation, Conceptually

CPU

Virtual Address

Translation

Invalid

Raise Exception

Valid

Who does this?

Physical Address

Physical Memory

Data

Data

# Memory Management Unit (MMU)

- Hardware device

  - Maps virtual addresses to physical addresses

- User process

  - deals with virtual addresses

  - never sees the physical address

- Physical memory

  - deals with physical addresses

  - never sees the virtual address


Motorola 68000

# Protection

- The functions used by different processes map their virtual addresses to disjoint ranges of physical addresses
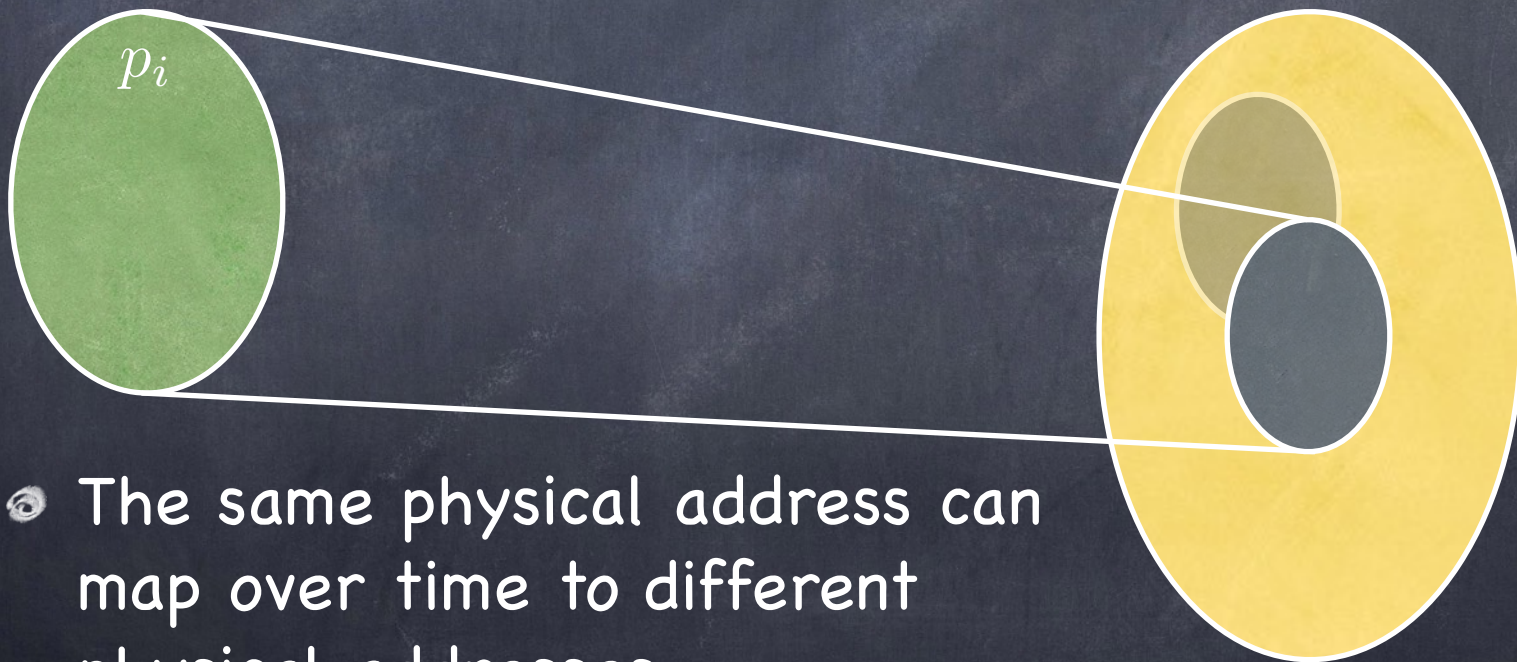


$p_i$

$p_j$

# Relocation

- The range of the function used by a process can change over time
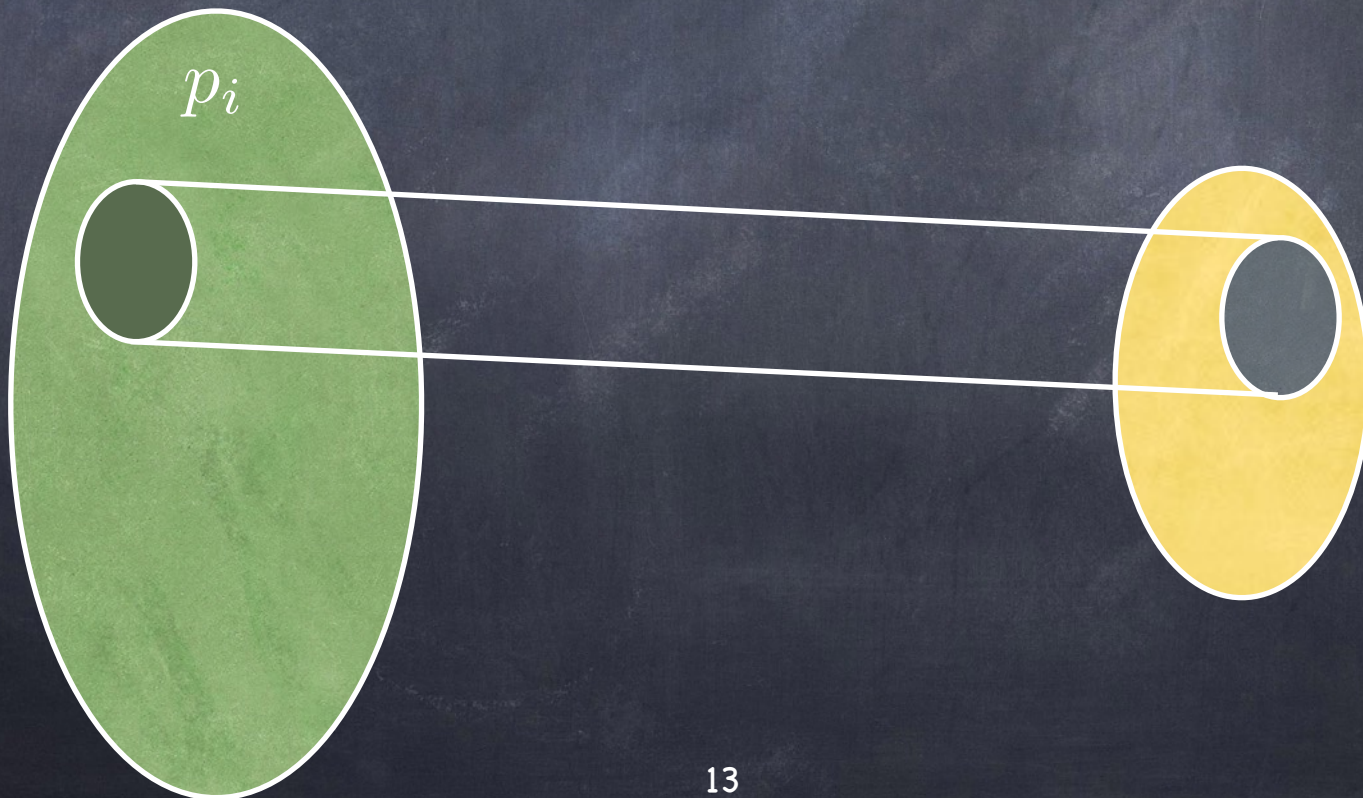


$p_i$

# Relocation

- The range of the function used by a process can change over time

$p_i$

- The same physical address can map over time to different physical addresses

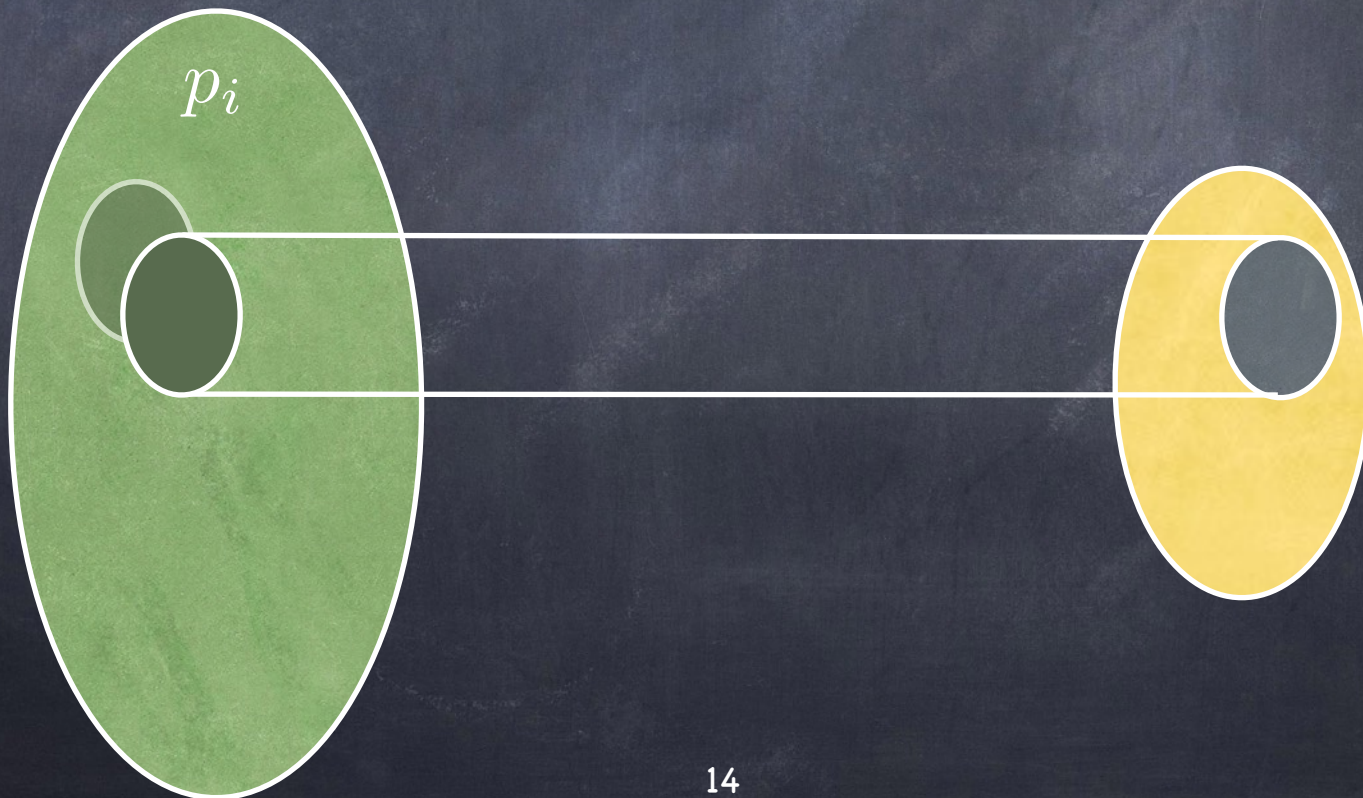  - or the mapping can be (temporarily) undefined

# Multiplexing

- The set of virtual addresses that map to a given range of physical addresses can change over time
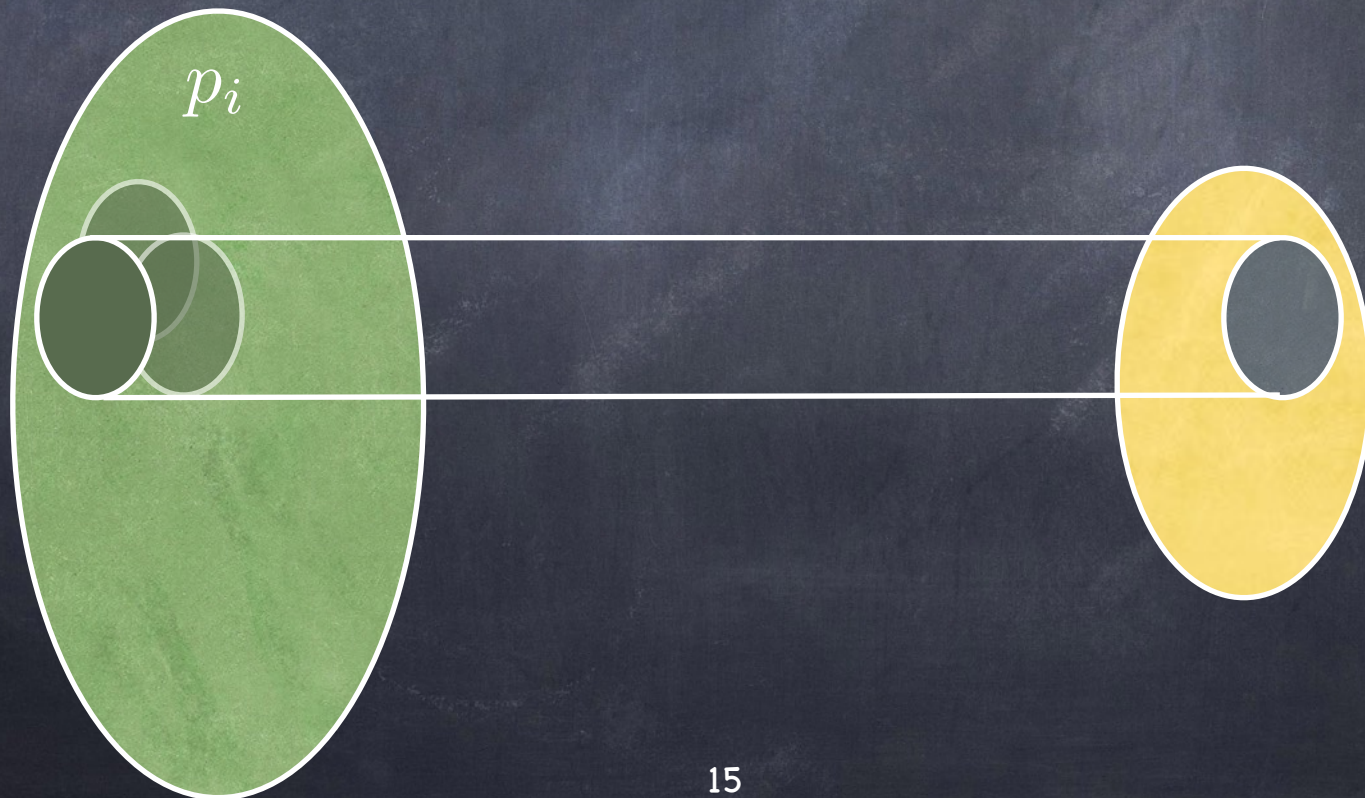
$p_i$

# Multiplexing

- The set of virtual addresses that map to a given range of physical addresses can change over time
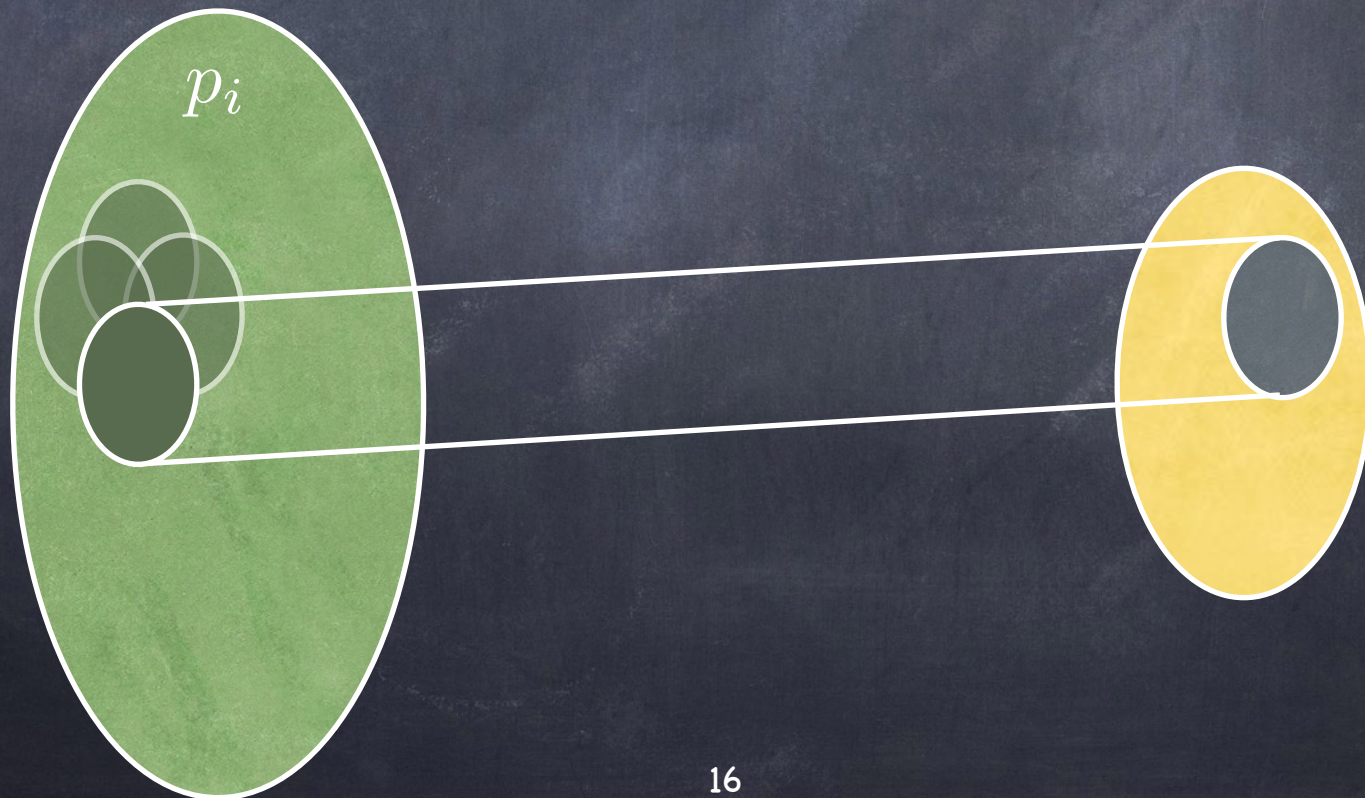
$p_i$

# Multiplexing

- The set of virtual addresses that map to a given range of physical addresses can change over time
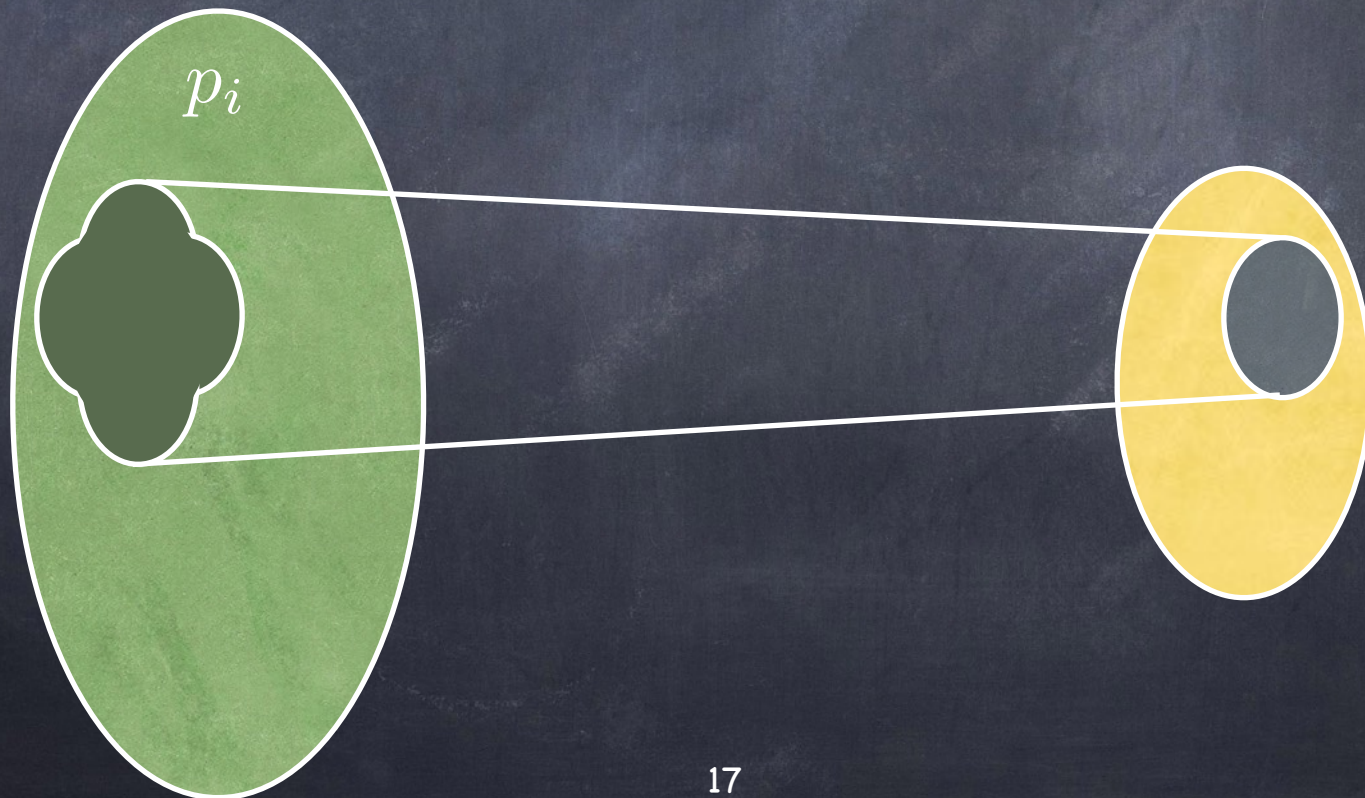
$p_i$

# Multiplexing

- The set of virtual addresses that map to a given range of physical addresses can change over time
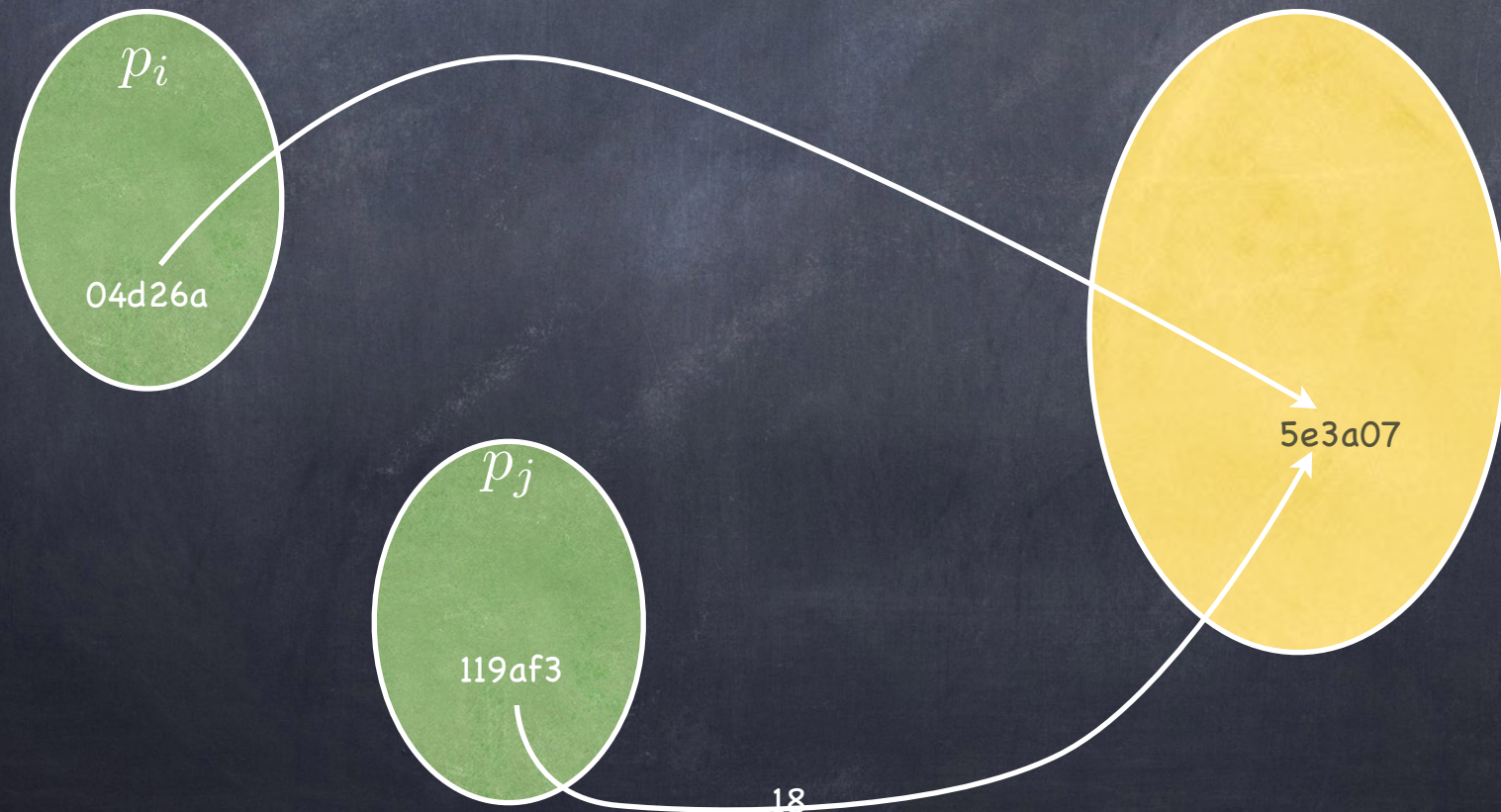
$p_i$

# Multiplexing

- The set of virtual addresses that map to a given range of physical addresses can change over time

$p_i$

# Data Sharing

- Map different virtual addresses of different processes to the same physical address
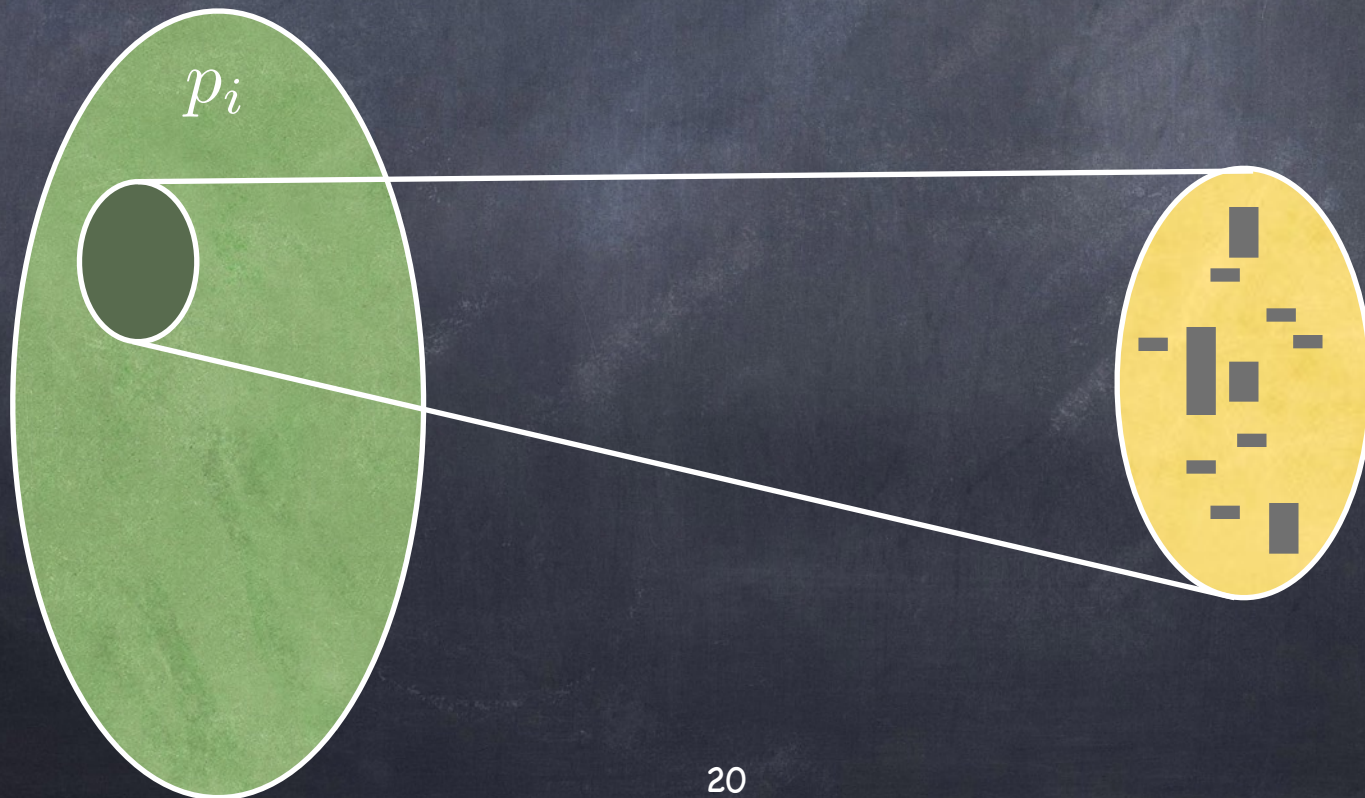
$p_i$

04d26a

$p_j$

119af3

5e3a07

# Contiguity

- Contiguous virtual addresses need not map to contiguous physical addresses

$p_i$

# Contiguity

- Contiguous virtual addresses need not map to contiguous physical addresses

$p_i$

# The Identity Mapping

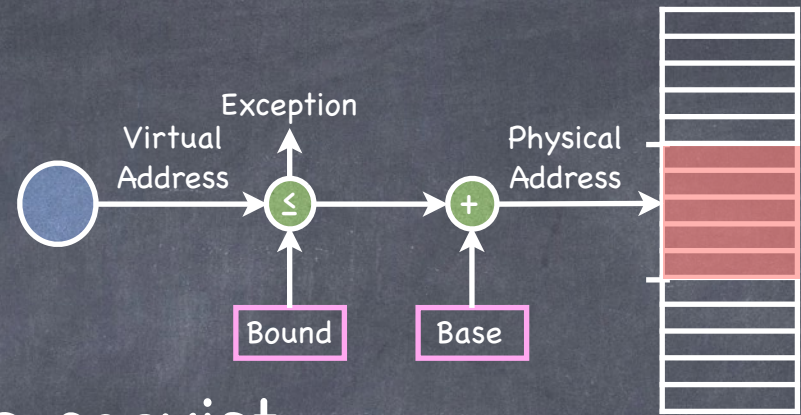- Map each virtual address onto the identical physical address

  - Virtual and physical address spaces have the same size

  - Run a single program at a time

    - OS can be a simple library

    - very early computers

- Friendly amendment: leave some of the physical address space for the OS

  - Use loader to relocate process

    - early PCs

Max

| Stack |
|-------|
| free |
| Heap |
| Text, Data, etc |
| OS |

0

# More sophisticated address translation

🌀 How to perform the mapping <u>efficiently</u>?

- ◻ So that it can be represented concisely?

- ◻ So that it can be computed quickly?

- ◻ So that it makes efficient use of the limited physical memory?

- ◻ So that multiple processes coexist in physical memory while guaranteeing isolation?

- ◻ So that it decouples the size of the virtual and physical addresses?

🌀 Ask hardware for help!

# Base & Bound

Exception

Virtual
Address

Physical
Address

Bound    Base

- Goal: let multiple processes coexist in memory while guaranteeing isolation

- Needed hardware
  - two registers: Base and Bound (a.k.a. Limit)
  - Stored in the PCB

- Mapping
  - pa = va + Base
    - as long as 0 ≤ va ≤ Bound
  - On context switch, change B&B (privileged instruction)

# Base & Bound

- $P_1$ : Base = 1000; Bound = 300
- $P_2$ : Base = 500; Bound = 400

Memory
Exception

no

CPU

Virtual
address

$\leq$

yes

$+$

Physical
address

Bound
Register

Base
Register

MAX$_{sys}$

1300

1000

0

# Base & Bound

- $P_1$ : Base = 1000; Bound = 300
- $P_2$ : Base = 500; Bound = 400

Memory
Exception

no

150

CPU

Virtual
address

$P_1$

< 

yes

+

Physical
address

300

Bound
Register

1000

Base
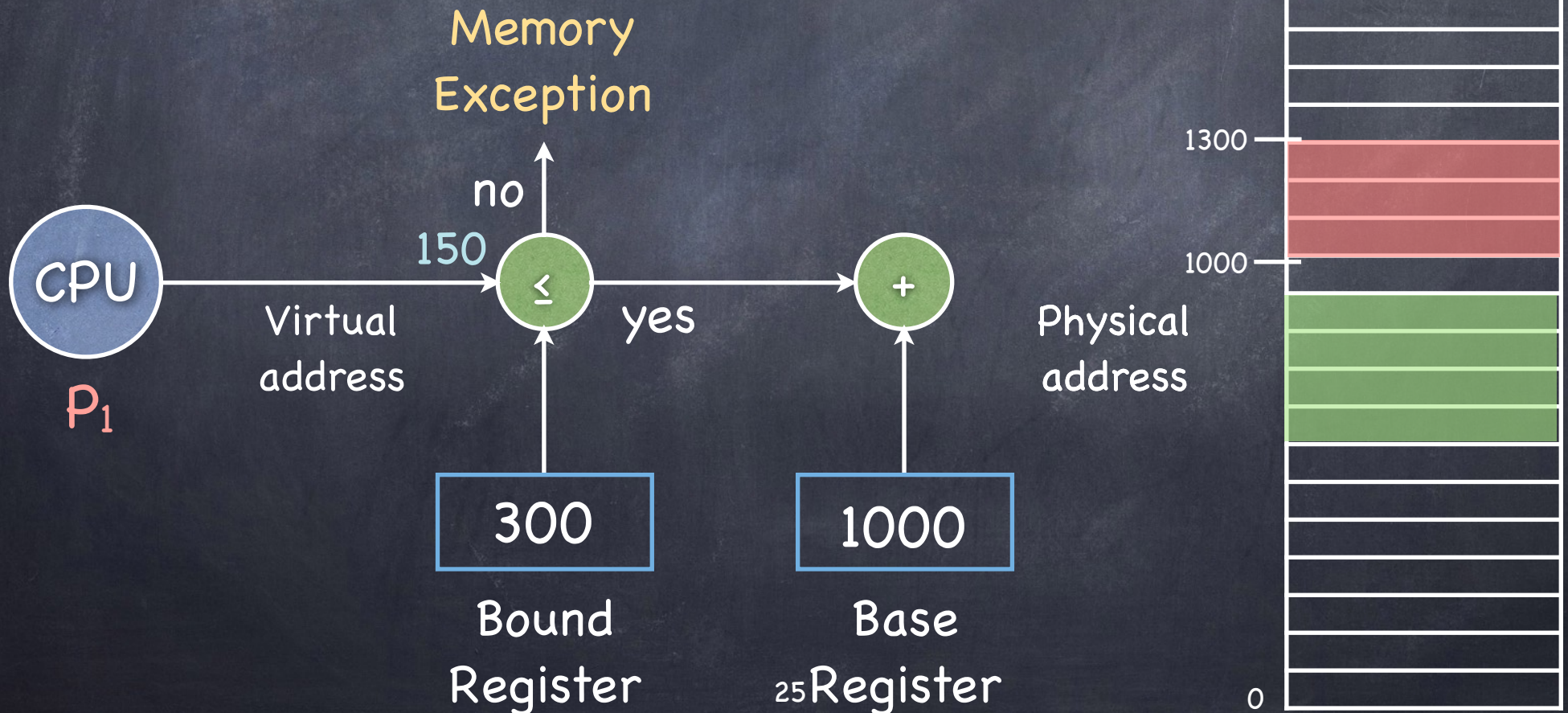25 Register

MAX$_{sys}$

1300

1000

0

# Base & Bound

- $P_1$ : Base = 1000; Bound = 300
- $P_2$ : Base = 500; Bound = 400

# Base & Bound

- $P_1$ : Base = 1000; Bound = 300
- $P_2$ : Base = 500; Bound = 400

Memory Exception

CPU

$P_1$

Virtual address

no

≤

yes

+

1150

Physical address

Context Switch

Base & Bound saved in $P_1$'s PCB

300

Bound Register

1000

Base

27 Register

MAX_sys

1300

1000

0

# Base & Bound

- $P_1$ : Base = 1000; Bound = 300
- $P_2$ : Base = 500; Bound = 400

# On Base & Bound

- Contiguous Allocation

  - <u>contiguous</u> virtual addresses are mapped to <u>contiguous</u> physical addresses

- But mapping entire address space to physical memory

  - is wasteful

    - lots of free space between heap and stack...

    - makes sharing hard

  - does not work if the address space is larger than physical memory

    - think 64-bit registers...

# E Pluribus Unum

- An address space comprises multiple segments

  - contiguous sets of virtual addresses, logically connected

    - heap, code, stack, (and also globals, libraries...)

  - each segment can be of a different size

| Stack |
|:-:|
| free |
| Heap |
| Data |
| Text |

# Segmentation: Generalizing Base & Bound

- Base & Bound registers to each segment

  - each segment is independently mapped to a set of contiguous addresses in physical memory

    - no need to map unused virtual addresses

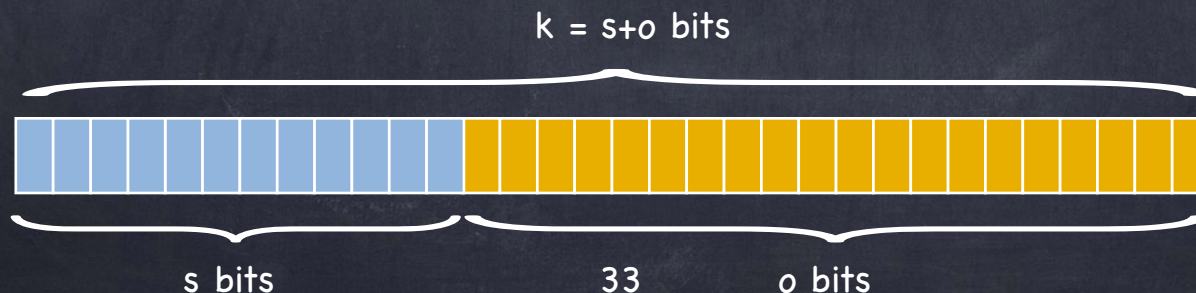| Segment | Base | Bound |
|---------|------|-------|
| Code | 10K | 2K |
| Stack | 28 | 2K |
| Heap | 35K | 3K |



(not to scale)

31

# Segmentation

- Goal: Supporting large address spaces (while allowing multiple processes to coexist in memory)

- Needed hardware

  - two registers (Base and Bound) per segment

    - values stored in the PCB

  - if many segments, a segment table, stored in memory, at an address pointed to by a Segment Table Register (STBR)

    - process' STBR value stored in the PCB

# Segmentation: Mapping

- How do we map a virtual address to the appropriate segment?
  - Read VA as having two components
    - s most significant bits identify the segment
      - at most $2^s$ segments
    - o remaining bits identify offset within segment
      - each segment's size can be at most $2^o$ bytes

k = s+o bits

s bits              33        o bits

# Segment Table

- Use s bits to index to the appropriate row of the segment table

|  | Base | Bound (Max 4k) | Access |
|---|---|---|---|
| Code$_{00}$ | 32K | 2K | Read/Execute |
| Heap$_{01}$ | 34K | 3K | Read/Write |
| Stack$_{10}$ | 28K | 3K | Read/Write |

- Segments can be shared by different processes
  - use protection bits to determine if shared Read only (maintaining isolation) or Read/Write (if shared, no isolation)
    - processes can share code segment while keeping data private

34

# Implementing Segmentation

## Segment table
generalizes Base & Bound

CPU

$s$  $o$

STBR  Segment Table Base Register

$s$

| Base | Bound | Access |
|------|-------|--------|
|      |       |        |
|      |       |        |
|      |       |        |
| 40K  | 512   | R/X    |
|      |       |        |
|      |       |        |

Logical addresses

Memory exception

no

$<$

yes

Bound
512

Base
40K

$+$

Physical addresses

MAX$_{sys}$

40K

0