# Operating Systems: Perspectives and Summary

**Computer Systems**

Nov. 25, 2017

Troels Henriksen

**Based on slides by:**
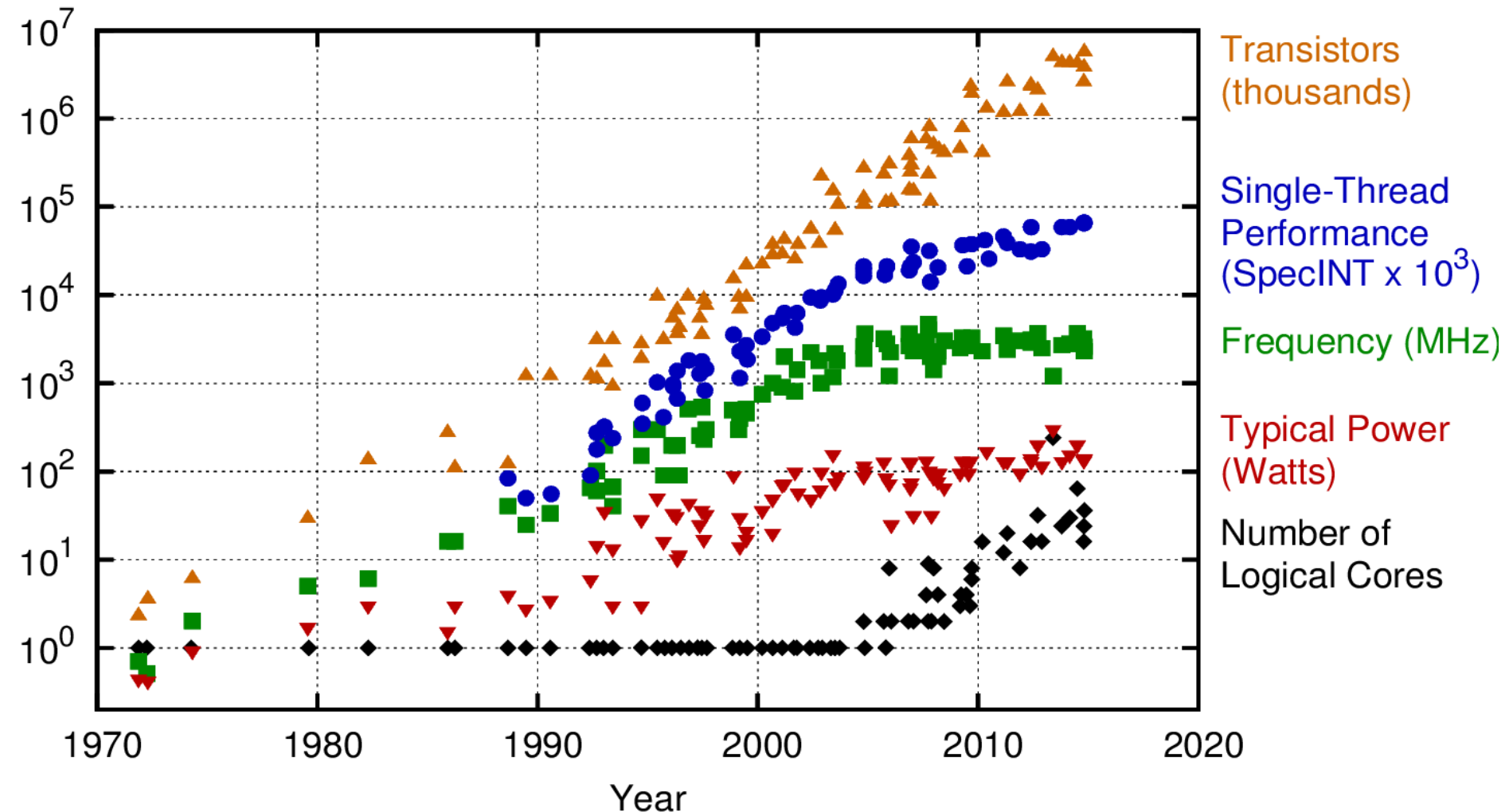
Randal E. Bryant and David R. O'Hallaron

# Motivation for Performance

- **The *only purpose* of a computing machine is to be faster than a human.**

- **All novel programs are the result of a good idea combined with a *performance surplus.***
  - Surplus can be generated by new/more/better machines.
  - … or by clever programming.

- **We can no longer (only) depend on engineers solving our problems by building better machines.**
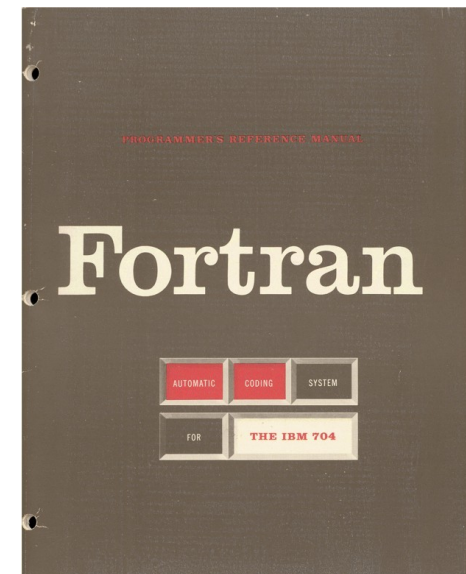
# Our Situation

## 40 Years of Microprocessor Trend Data



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

# Example: high-level languages

- **The performance surplus:** computers of the 50s got faster and faster (1000s of statements per second!).

- **The good idea:** a high-level language (FORTRAN) could improve productivity, in most cases offsetting the lower performance compared to hand-coding.

- **The edge:** an *optimizing* compiler (particularly CSE) was used to narrow the gap (see *THE FORTRAN AUTOMATIC CODING SYSTEM* from 1957).
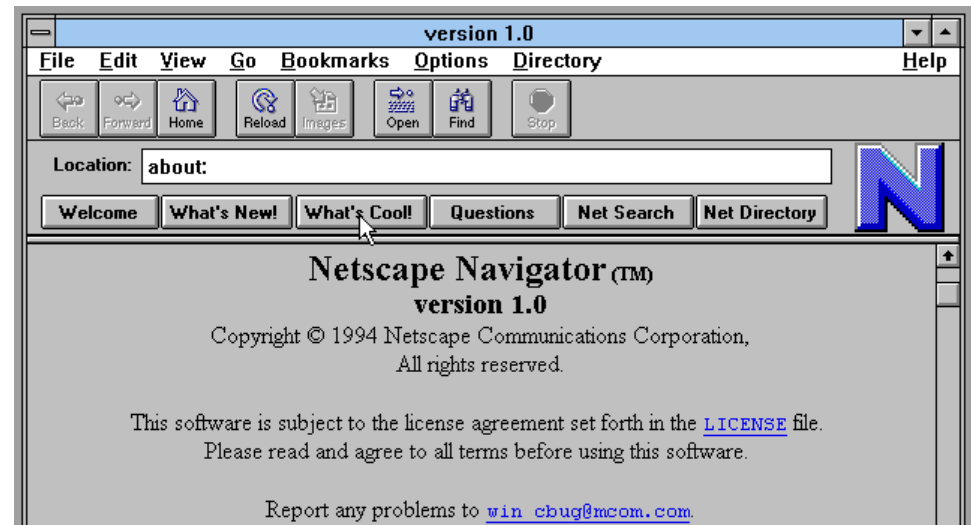
# Example: operating systems

- **The performance surplus:** increasing transistor budgets afforded non-computational circuits like MMUs.
- **The good idea:** impose a virtualisation layer that permitted running multiple applications simultaneously, safely.
- **The edge:** a clean distinction between API and implementation, making a single program runnable on machines of vastly different capabilities (IBM System/360).

# Example: Visual web browsers

- **The performance surplus:** personal computers of the early 90s got fast enough to run sophisticated GUIs.
- **The good idea:** accessible hypertext along with *really clever programming* in the browsers – particularly to handle concurrent network requests alongside rendering.
  - **Later:** JIT compilation of Javascript (particularly Chrome's V8 in 2008) sped up web applications to create a new *software-based* performance surplus.

# Example: Deep Learning

- **The good (and old) idea:** deep sequences of simple layers of *neurons* can be trained to perform input classification, if given sufficient (huge) numbers of examples.

- **The performance surplus:** cheap *massive parallelism* in the form of generally programmable graphics processors (GPUs), funded by millions of gamers in the 90s.

- **The edge:** programming tools and techniques that made GPUs accessible to more than just graphics.

hidden neurons

output neurons

input neurons

# Non-Unix operating systems

Not everything has to be as I have taught you.  It is healthy to look at very different ways of doing things.

- **Forth:** mostly on bare hardware
- **Lisp Machines:** hardware support for a high-level language
- **Smalltalk:** image-based development https://squeak.org/
- **VMS:** historical competitor to Unix
- **z/OS:** mainframe OS from IBM

# Summary: Control Flow

# Control Flow

- **Processors do only one thing:**
    - From startup to shutdown, a CPU simply reads and executes (interprets) a sequence of instructions, one at a time
    - This sequence is the CPU's *control flow* (or *flow of control*)

*Physical control flow*

**Time**

<startup>

$inst_1$

$inst_2$

$inst_3$

**...**

$inst_n$

# Exceptional Control Flow

- **Exists at all levels of a computer system**

- **Low level mechanisms**
  - 1. **Exceptions**
    - Change in control flow in response to a system event (i.e., change in system state)
    - Implemented using combination of hardware and OS software

- **Higher level mechanisms**
  - 2. **Process context switch**
    - Implemented by OS software and hardware timer
  - 3. **Signals**
    - Implemented by OS software
  - 4. **Nonlocal jumps**: `setjmp()` and `longjmp()`
    - Implemented by C runtime library

# Exceptions

- **An *exception* is a transfer of control to the OS *kernel* in response to some *event*  (i.e., change in processor state)**
    - Kernel is the memory-resident part of the OS
    - Examples of events: Divide by 0, arithmetic overflow, page fault, I/O request completes, typing Ctrl-C

*User code*                    *Kernel code*

*Event* ⟶ I_current          *Exception*

I_next                        *Exception processing by exception handler*

- Return to I_current
- Return to I_next
- Abort

# Asynchronous Exceptions (Interrupts)

- **Caused by events external to the processor**
  - Indicated by setting the processor's *interrupt pin*
  - Handler returns to "next" instruction

- **Examples:**
  - Timer interrupt
    - Every few ms, an external timer chip triggers an interrupt
    - Used by the kernel to take back control from user programs
  - I/O interrupt from external device
    - Hitting Ctrl-C at the keyboard
    - Arrival of a packet from a network
    - Arrival of data from a disk

# Synchronous Exceptions

- **Caused by events that occur as a result of executing an instruction:**
  - *Traps*
    - Intentional
    - Examples: **system calls**, breakpoint traps, special instructions
    - Returns control to "next" instruction
  - *Faults*
    - Unintentional but possibly recoverable
    - Examples: page faults (recoverable), protection faults (unrecoverable), floating point exceptions
    - Either re-executes faulting ("current") instruction or aborts
  - *Aborts*
    - Unintentional and unrecoverable
    - Examples: illegal instruction (actually SIGILL), parity error, machine check
    - Aborts current program

# System Calls

- **Each x86-64 system call has a unique ID number**
- **Examples:**

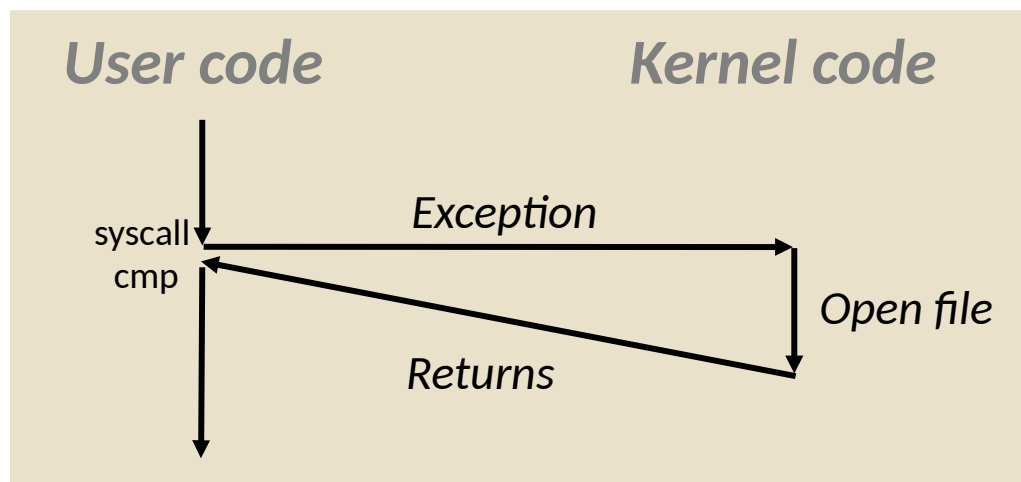| Number | Name | Description |
|--------|--------|----------------------|
| 0 | read | Read file |
| 1 | write | Write file |
| 2 | open | Open file |
| 3 | close | Close file |
| 4 | stat | Get info about file |
| 57 | fork | Create process |
| 59 | execve | Execute a program |
| 60 | _exit | Terminate process |
| 62 | kill | Send signal to process |

# System Call Example: Opening File

- User calls: `open(filename, options)`
- Calls `__open` function, which invokes system call instruction `syscall`

```
00000000000e5d70 <__open>:
...
e5d79:    b8 02 00 00 00        mov  $0x2,%eax  # open is syscall #2
e5d7e:    0f 05                 syscall         # Return value in %rax
e5d80:    48 3d 01 f0 ff ff     cmp  $0xfffffffffffff001,%rax
...
e5dfa:    c3                    retq
```

**User code**          **Kernel code**

syscall ↓
cmp
*Exception*

*Returns*          *Open file*

- `%rax` contains syscall number
- Other arguments in `%rdi`, `%rsi`, `%rdx`, `%r10`, `%r8`, `%r9`
- Return value in `%rax`
- Negative value is an error corresponding to negative `errno`

16
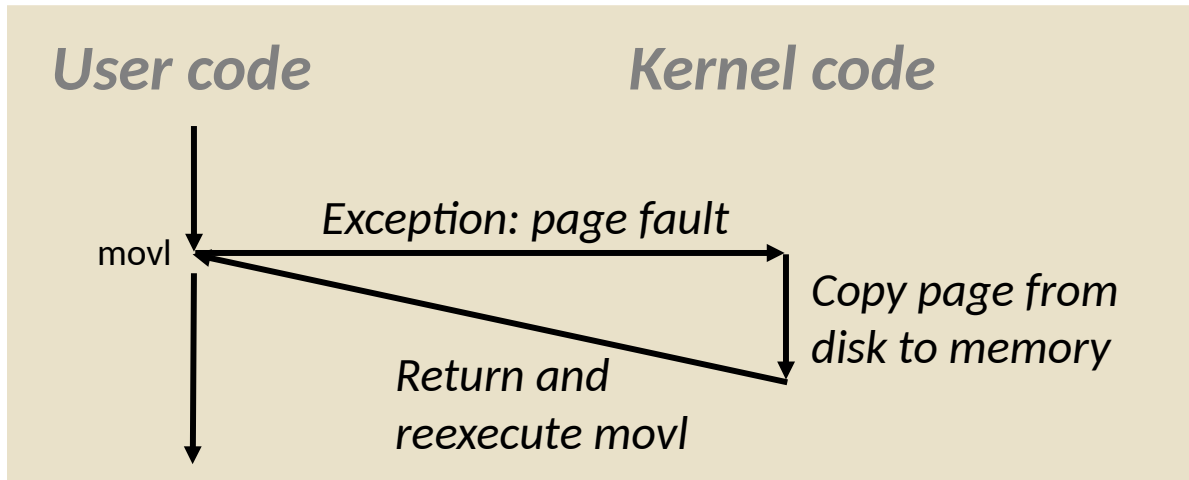
# Fault Example: Page Fault

- User writes to memory location
- That portion (page) of user's memory is currently on disk

```
int a[1000];
main ()
{
    a[500] = 13;
}
```

```
80483b7:   c7 05 10 9d 04 08 0d   movl   $0xd,0x8049d10
```

**User code**                **Kernel code**

movl

Exception: page fault

Copy page from disk to memory

Return and reexecute movl

17

# Summary: Processes

# Processes

- **Definition: A *process* is an instance of a running program.**
  - One of the most profound ideas in computer science
  - Not the same as "program" or "processor"

- **Process provides each program with two key abstractions:**
  - *Logical control flow*
    - Each program seems to have exclusive use of the CPU
    - Provided by kernel mechanism called *context switching*
  - *Private address space*
    - Each program seems to have exclusive use of main memory.
    - Provided by kernel mechanism called *virtual memory*

**Memory**

| Stack |
| Heap |
| Data |
| Code |

**CPU**

| Registers |

# Multiprocessing: The Illusion

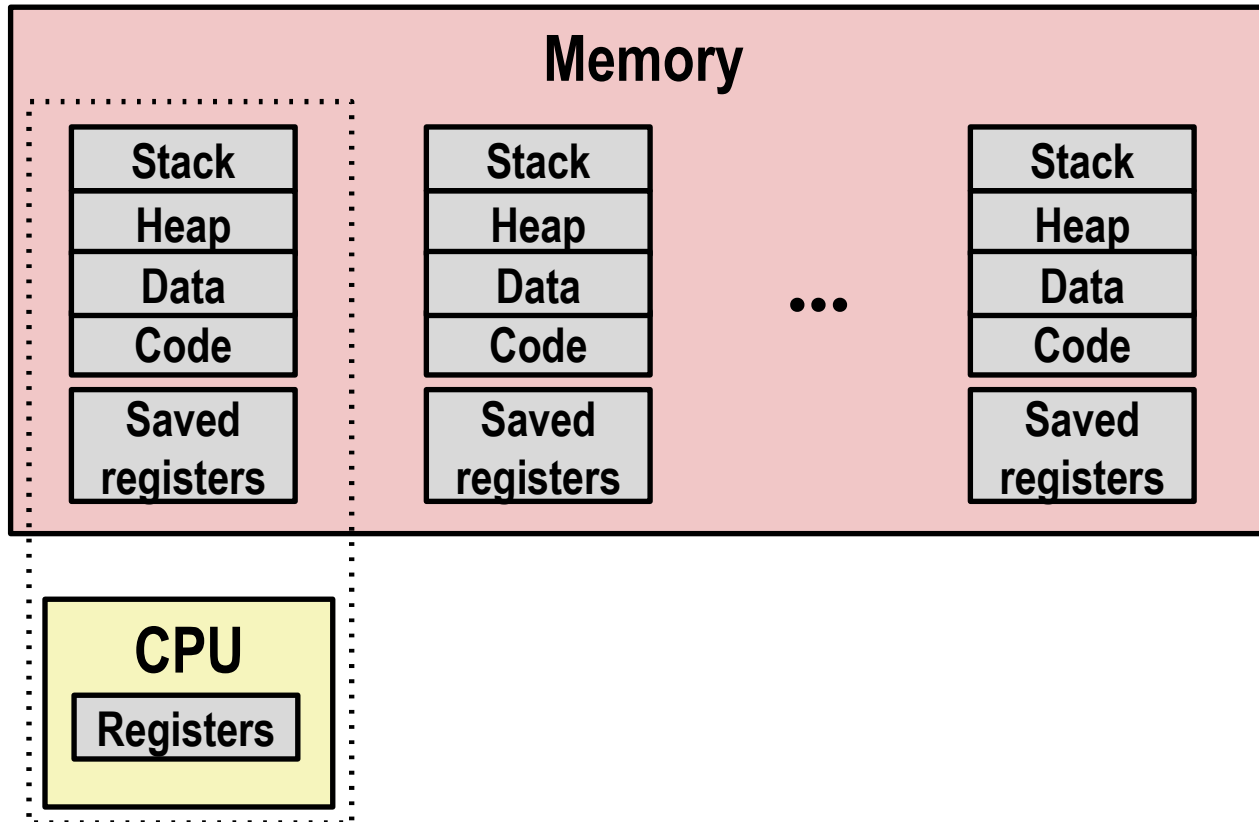| Memory | | Memory | | | Memory |
|---|---|---|---|---|---|
| Stack | | Stack | | | Stack |
| Heap | | Heap | ••• | | Heap |
| Data | | Data | | | Data |
| Code | | Code | | | Code |
| **CPU** | | **CPU** | | | **CPU** |
| Registers | | Registers | | | Registers |

- **Computer runs many processes simultaneously**
  - Applications for one or more users
    - Web browsers, email clients, editors, …
  - Background tasks
    - Monitoring network & I/O devices

# Multiprocessing: The (Traditional) Reality

**Memory**

| Stack | | Stack | | | Stack |
|-------|--|-------|--|--|-------|
| Heap | | Heap | | | Heap |
| Data | | Data | **. . .** | | Data |
| Code | | Code | | | Code |
| Saved registers | | Saved registers | | | Saved registers |

**CPU**

**Registers**

- **Single processor executes multiple processes concurrently**
  - Process executions interleaved (multitasking)
  - Address spaces managed by virtual memory system (later in course)
  - Register values for nonexecuting processes saved in memory

# Multiprocessing: The (Traditional) Reality

**Memory**

| | | | |
|---|---|---|---|
| **Stack** | **Stack** | | **Stack** |
| **Heap** | **Heap** | | **Heap** |
| **Data** | **Data** | **● ● ●** | **Data** |
| **Code** | **Code** | | **Code** |
| **Saved registers** | **Saved registers** | | **Saved registers** |

**CPU**

**Registers**

- **Save current registers in memory**

# Multiprocessing: The (Traditional) Reality

**Memory**

| Stack |
|-------|
| Heap |
| Data |
| Code |

| Saved registers |

| Stack |
|-------|
| Heap |
| Data |
| Code |

| Saved registers |

• • •

| Stack |
|-------|
| Heap |
| Data |
| Code |

| Saved registers |

**CPU**

| Registers |

- **Schedule next process for execution**

# Multiprocessing: The (Traditional) Reality

**Memory**

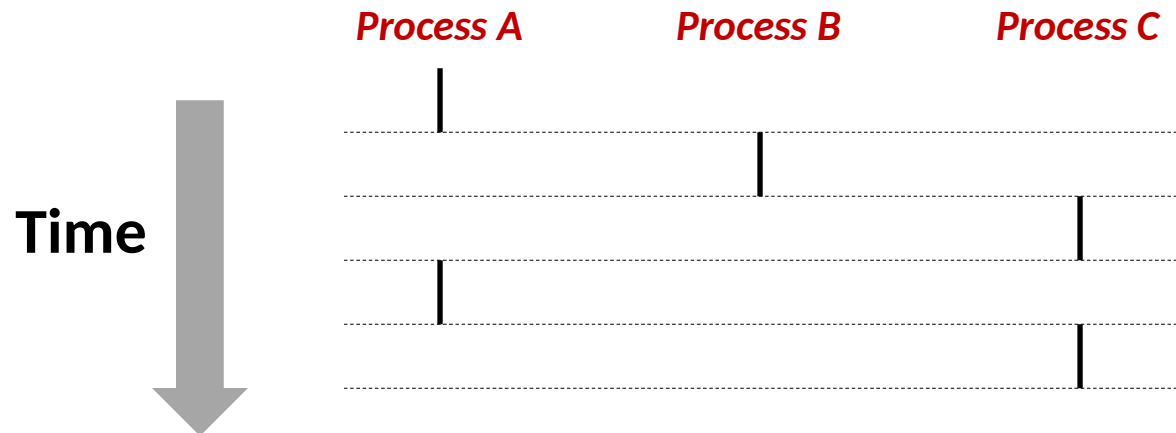| Stack | | Stack | | | Stack |
|-------|---|-------|---|---|-------|
| Heap | | Heap | | | Heap |
| Data | | Data | **• • •** | | Data |
| Code | | Code | | | Code |
| Saved registers | | Saved registers | | | Saved registers |

**CPU**

**Registers**

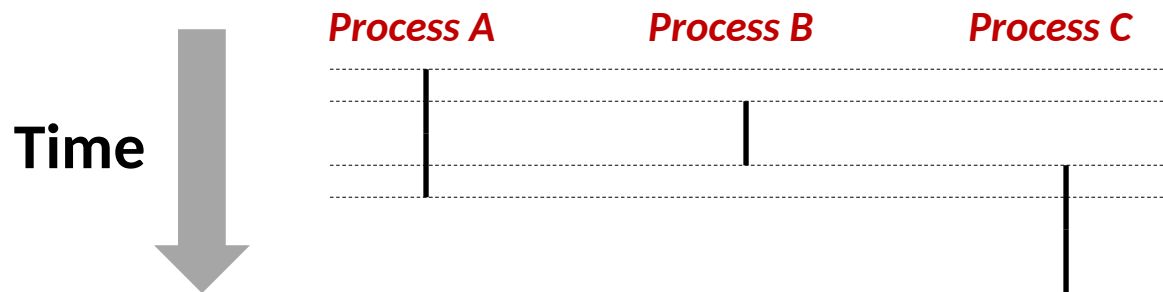■ **Load saved registers and switch address space (context switch)**

# Concurrent Processes

- **Each process is a logical control flow.**
- **Two processes *run* *concurrently* (*are concurrent)* if their flows overlap in time**
- **Otherwise, they are *sequential***
- **Examples (running on single core):**
  - Concurrent: A & B, A & C
  - Sequential: B & C

*Process A*    *Process B*    *Process C*

**Time**

# User View of Concurrent Processes

- **Control flows for concurrent processes are physically disjoint in time**

- **However, we can think of concurrent processes as running in parallel with each other**

*Process A*    *Process B*    *Process C*

**Time**

# Creating and Terminating Processes

**From a programmer's perspective, we can think of a process as being in one of three states**

- **Running**
  - Process is either executing, or waiting to be executed and will eventually be *scheduled* (i.e., chosen to execute) by the kernel

- **Stopped**
  - Process execution is *suspended* and will not be scheduled until further notice

- **Terminated**
  - Process is stopped permanently

# Terminating Processes

- **Process becomes terminated for one of three reasons:**
  - Receiving a signal whose default action is to terminate (next lecture)
  - Returning from the `main` routine
  - Calling the `exit` function
  - Last thread terminates with `pthread_exit`

- **`void exit(int status)`**
  - Terminates with an *exit status* of `status`
  - Convention: normal return status is 0, nonzero on error
  - Another way to explicitly set the exit status is to return an integer value from the main routine

- **`exit` is called <span style="color:red">once</span> but <span style="color:red">never</span> returns.**

# Creating Processes

- *Parent process* creates a new running *child process* by calling `fork`

- `int fork(void)`
  - Returns 0 to the child process, child's PID to parent process
  - Child is *almost* identical to parent:
    - Child get an identical (but separate) copy of the parent's virtual address space.
    - Child gets identical copies of the parent's open file descriptors
    - Child has a different PID than the parent

- `fork` is interesting (and often confusing) because it is called *once* but returns *twice*

# Reaping Child Processes

- **Idea**
  - When process terminates, it still consumes system resources
    - Examples: Exit status, various OS tables
  - Called a "zombie"
    - Living corpse, half alive and half dead
- **Reaping**
  - Performed by parent on terminated child (using `wait` or `waitpid`)
  - Parent is given exit status information
  - Kernel then deletes zombie child process
- **What if parent doesn't reap?**
  - If any parent terminates without reaping a child, then the orphaned child will be reaped by **init** process (pid == 1)
  - So, only need explicit reaping in long-running processes
    - e.g., shells and servers

# `wait`: Synchronizing with Children

- **Parent reaps a child by calling the `wait` function**


- **`int wait(int *child_status)`**
  - Suspends current process until one of its children terminates
  - Return value is the **`pid`** of the child process that terminated
  - If **`child_status != NULL`**, then the integer it points to will be set to a value that indicates reason the child terminated and the exit status:
    - Checked using macros defined in `wait.h`
      - `WIFEXITED, WEXITSTATUS, WIFSIGNALED, WTERMSIG, WIFSTOPPED, WSTOPSIG, WIFCONTINUED`
      - See textbook for details

# Summary: Signals

# Signals

- **A *signal* is a small message that notifies a process that an event of some type has occurred in the system**
  - Akin to exceptions and interrupts
  - Sent from the kernel (sometimes at the request of another process) to a process
  - Signal type is identified by small integer ID's (1-30)
  - Only information in a signal is its ID and the fact that it arrived

| ID | Name | Default Action | Corresponding Event |
|----|------|----------------|---------------------|
| 2 | SIGINT | Terminate | User typed ctrl-c |
| 9 | SIGKILL | Terminate | Kill program (cannot override or ignore) |
| 11 | SIGSEGV | Terminate | Segmentation violation |
| 14 | SIGALRM | Terminate | Timer signal |
| **17** | SIGCHLD | Ignore | **Child stopped or terminated** |

# Signal Concepts: Sending a Signal

- **Kernel *sends* (delivers) a signal to a *destination process* by updating some state in the context of the destination process**

- **Kernel sends a signal for one of the following reasons:**
  - Kernel has detected a system event such as divide-by-zero (SIGFPE) or the termination of a child process (SIGCHLD)
  - Another process has invoked the `kill` system call to explicitly request the kernel to send a signal to the destination process

# Signal Concepts: Receiving a Signal

- **A destination process *receives* a signal when it is forced by the kernel to react in some way to the delivery of the signal**

- **Some possible ways to react:**
  - ***Ignore*** the signal (do nothing)
  - ***Terminate*** the process (with optional core dump)
  - ***Catch*** the signal by executing a user-level function called ***signal handler***
    - Akin to a hardware exception handler being called in response to an asynchronous interrupt:

*(1) Signal received by process*

$I_{curr}$
$I_{next}$

*(2) Control passes to signal handler*

*(3) Signal handler runs*

*(4) Signal handler returns to next instruction*

35

# Signal Concepts: Pending and Blocked Signals

- **A signal is *pending* if sent but not yet received**
  - There can be at most one pending signal of any particular type
  - Important: Signals are not queued
    - If a process has a pending signal of type k, then subsequent signals of type k that are sent to that process are discarded

- **A process can *block* the receipt of certain signals**
  - Blocked signals can be delivered, but will not be received until the signal is unblocked

- **A pending signal is received at most once**
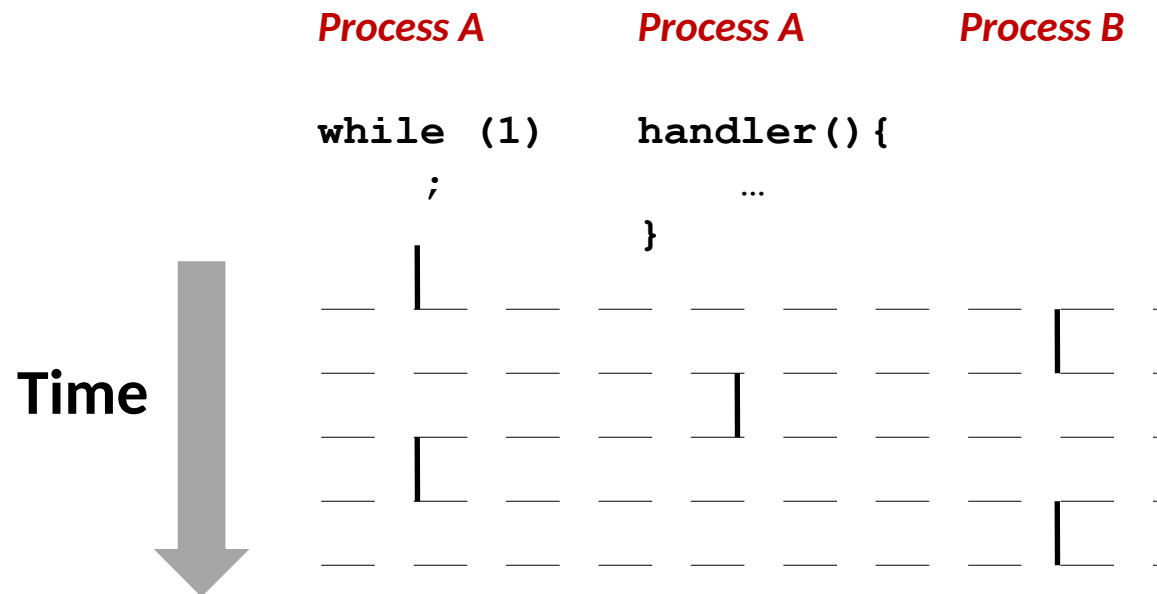
DIKU

# Signal Concepts: Pending/Blocked Bits

- **Kernel maintains `pending` and `blocked` bit vectors in the context of each process**
    - **`pending`: represents the set of pending signals**
        - Kernel sets bit k in **`pending`** when a signal of type k is delivered
        - Kernel clears bit k in **`pending`** when a signal of type k is received

    - **`blocked`: represents the set of blocked signals**
        - Can be set and cleared by using the **`sigprocmask`** function
        - Also referred to as the *signal mask*.

# Receiving Signals

- **Suppose kernel is returning from an exception handler and is ready to pass control to process *p***

- **Kernel computes `pnb = pending & ~blocked`**
  - The set of pending nonblocked signals for process *p*

- **If (`pnb == 0`)**
  - Pass control to next instruction in the logical flow for *p*

- **Else**
  - Choose least nonzero bit *k* in `pnb` and force process *p* to *receive* signal *k*
  - The receipt of the signal triggers some *action* by *p*
  - Repeat for all nonzero *k* in `pnb`
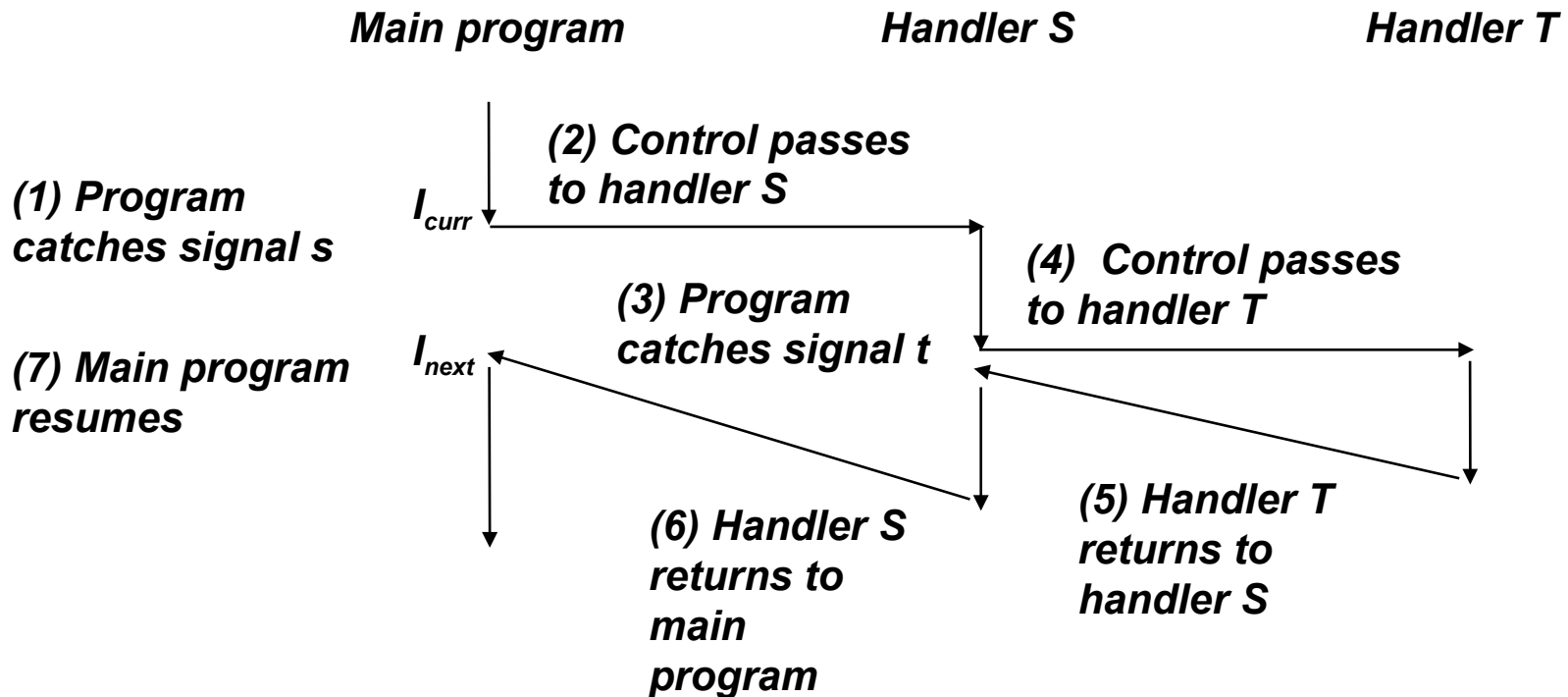  - Pass control to next instruction in logical flow for *p*

# Signals Handlers as Concurrent Flows

- **A signal handler is a separate logical flow (not process) that runs concurrently with the main program**

*Process A*          *Process A*          *Process B*

```
while (1)      handler(){
   ;                …
               }
```

**Time**

# Nested Signal Handlers

- **Handlers can be interrupted by other handlers**



**Main program**  **Handler S**  **Handler T**

*(1) Program catches signal s*

*I_curr*

*(2) Control passes to handler S*

*(4) Control passes to handler T*

*(7) Main program resumes*

*I_next*

*(3) Program catches signal t*

*(6) Handler S returns to main program*

*(5) Handler T returns to handler S*

# Summary: Virtual Memory

# A System Using Virtual Addressing

**Main memory**

*CPU Chip*

**CPU** → **Virtual address (VA)** `4100` → **MMU** → **Physical address (PA)** `4` → Main memory

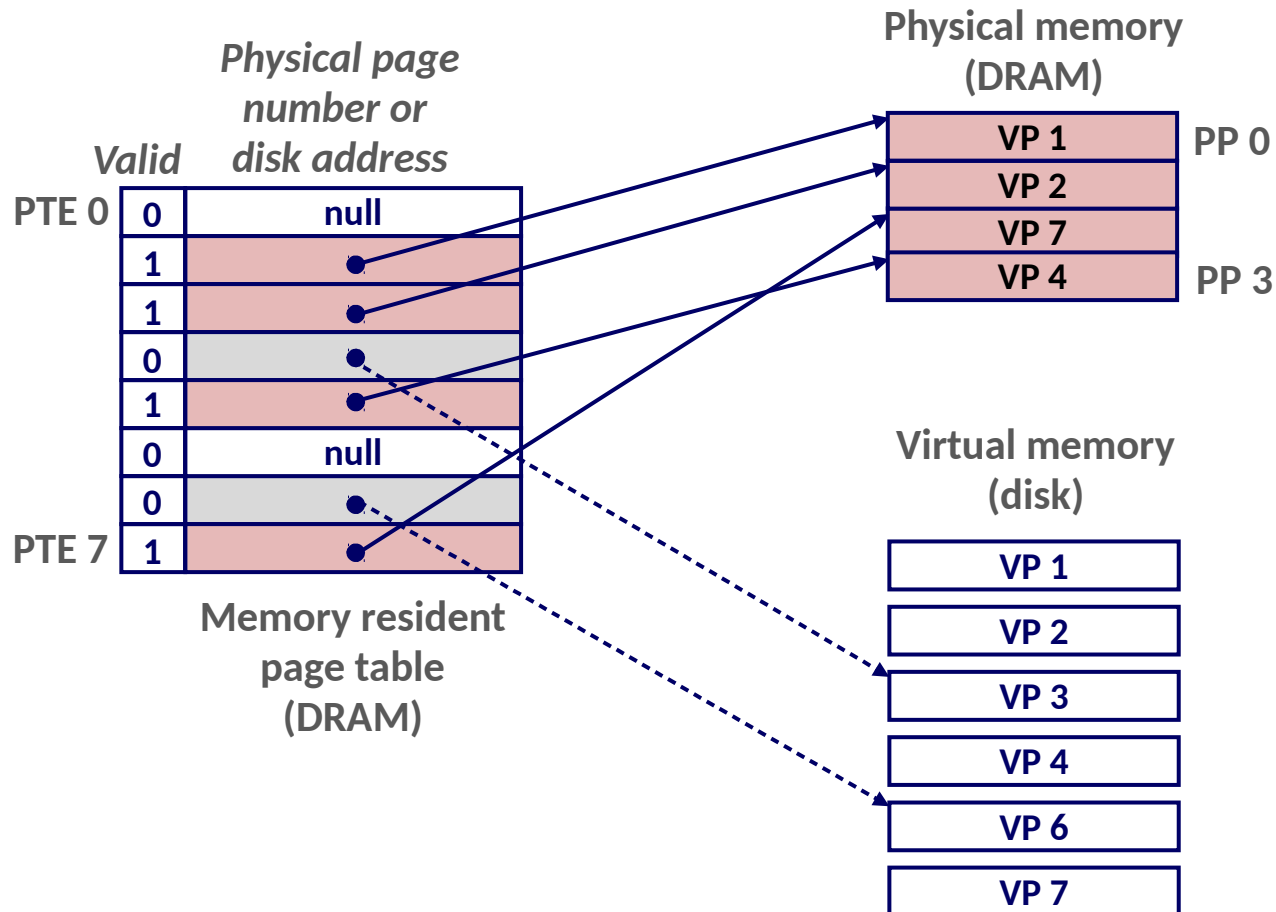0:
1:
2:
3:
4:
5:
6:
7:
8:

M-1:

**Data word**

- **Used in all modern servers, laptops, and smart phones**
- **One of the great ideas in computer science**

# Why Virtual Memory (VM)?

- **Uses main memory efficiently**
  - Use DRAM as a cache for parts of a virtual address space

- **Simplifies memory management**
  - Each process gets the same uniform linear address space

- **Isolates address spaces**
  - One process can't interfere with another's memory
  - User program cannot access privileged kernel information and code
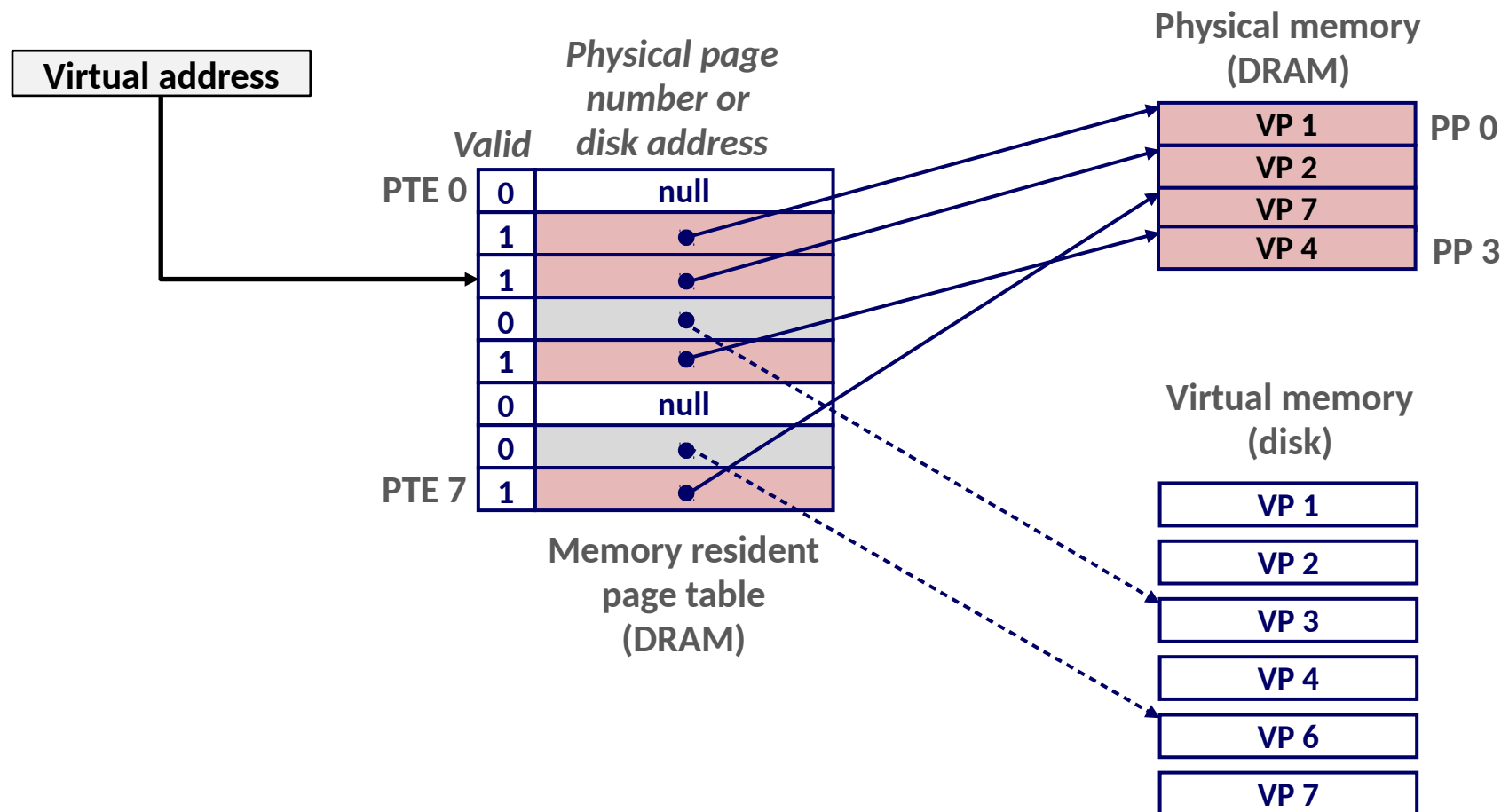
# Enabling Data Structure: Page Table

- A *page table* is an array of page table entries (PTEs) that maps virtual pages to physical pages.
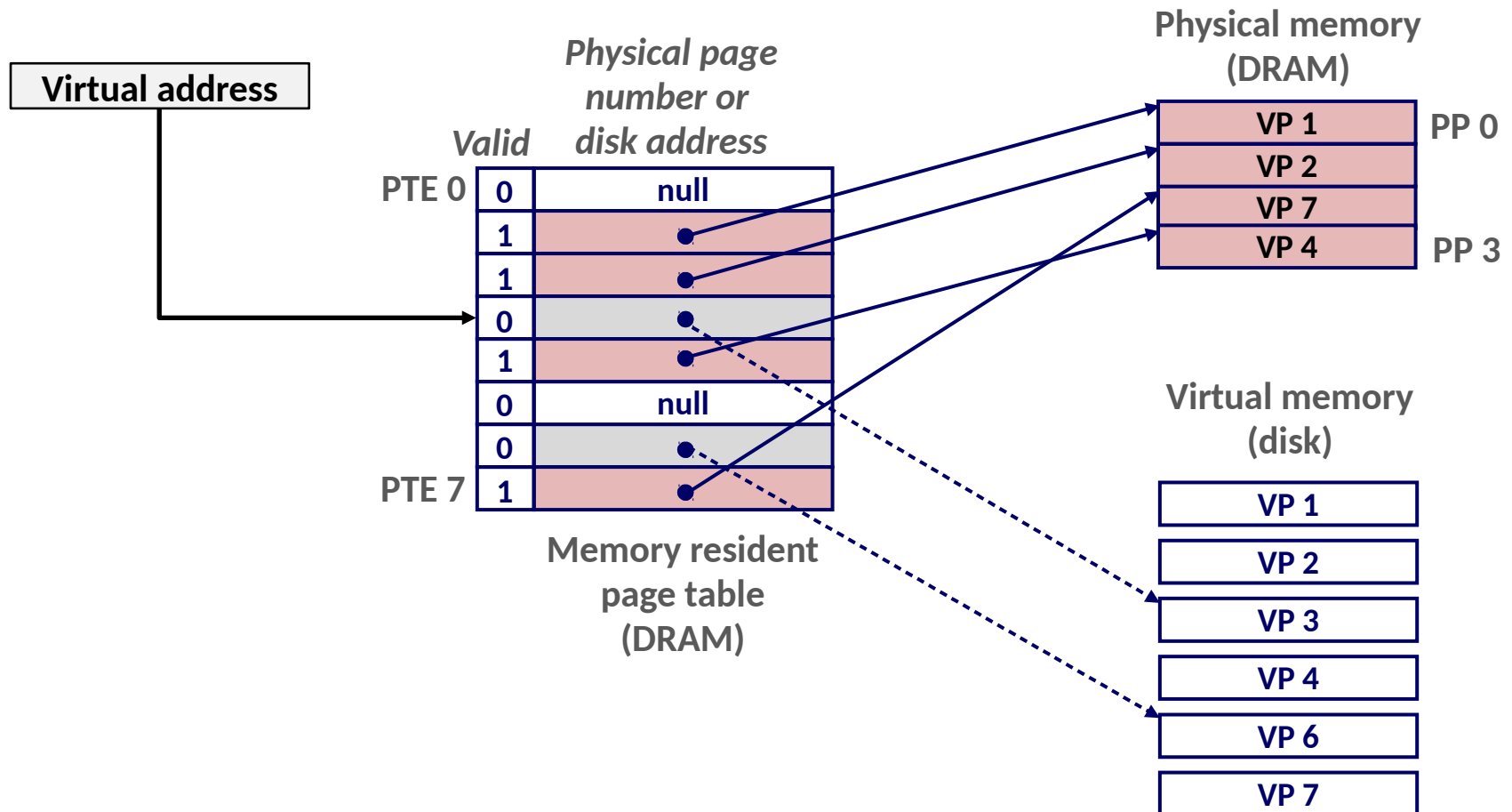  - Per-process kernel data structure in DRAM

# Page Hit

- *Page hit:* reference to VM word that is in physical memory (DRAM cache hit)

# Page Fault

- *Page fault:* reference to VM word that is not in physical memory (DRAM cache miss)

# Handling Page Fault

■ Page miss causes page fault (an exception)

**Virtual address**

**Physical page number or disk address**

**Physical memory (DRAM)**

**Valid**

| | |
|---|---|
| PTE 0 | 0 | null |
| | 1 | |
| | 1 | |
| | 0 | |
| | 1 | |
| | 0 | null |
| | 0 | |
| PTE 7 | 1 | |

**Memory resident page table (DRAM)**

| |
|---|
| VP 1 | PP 0 |
| VP 2 | |
| VP 7 | |
| VP 4 | PP 3 |

**Virtual memory (disk)**

| |
|---|
| VP 1 |
| VP 2 |
| VP 3 |
| VP 4 |
| VP 6 |
| VP 7 |

# Handling Page Fault

- Page miss causes page fault (an exception)
- Page fault handler selects a victim to be evicted (here VP 4)

# Handling Page Fault

- Page miss causes page fault (an exception)
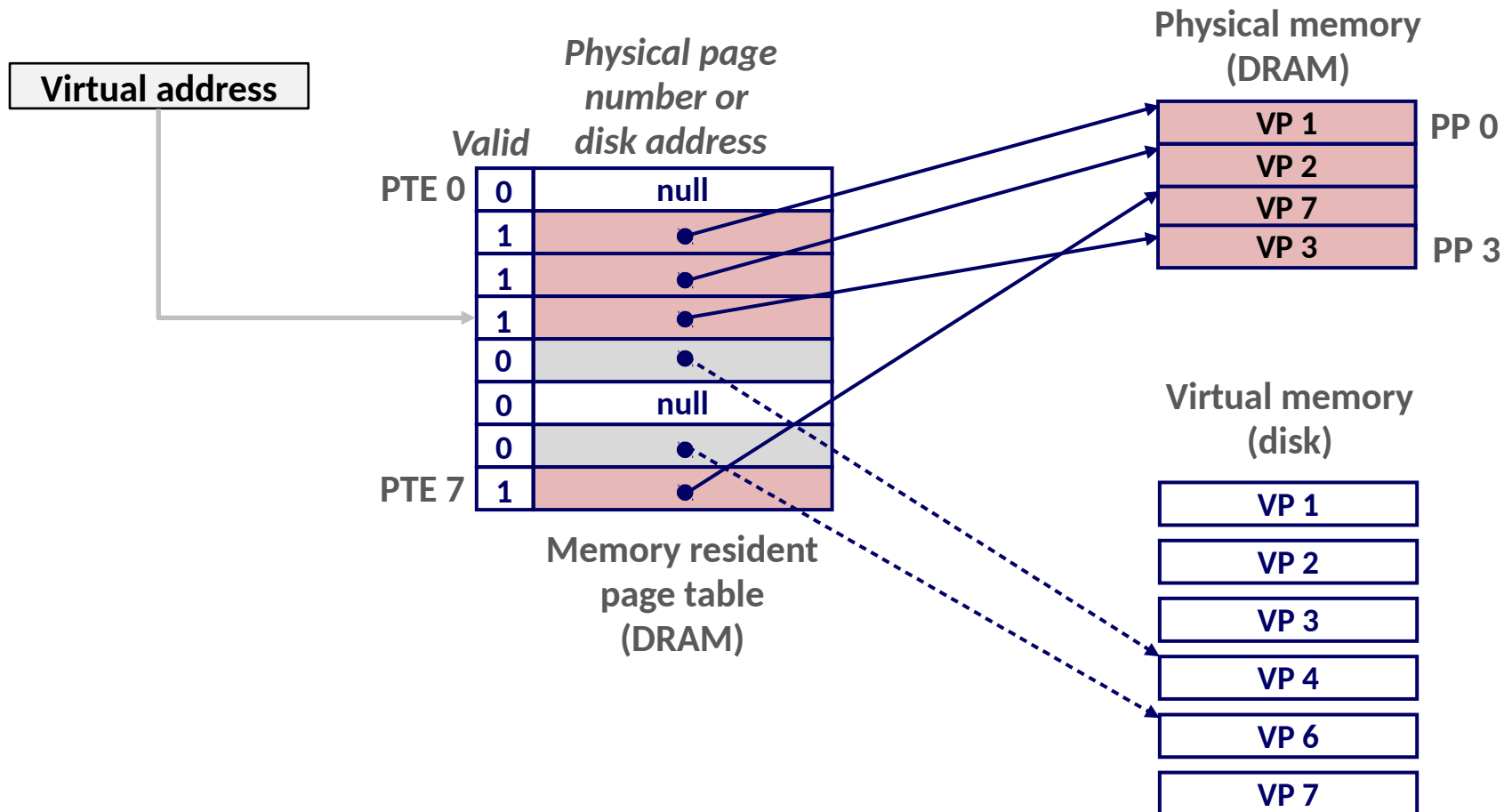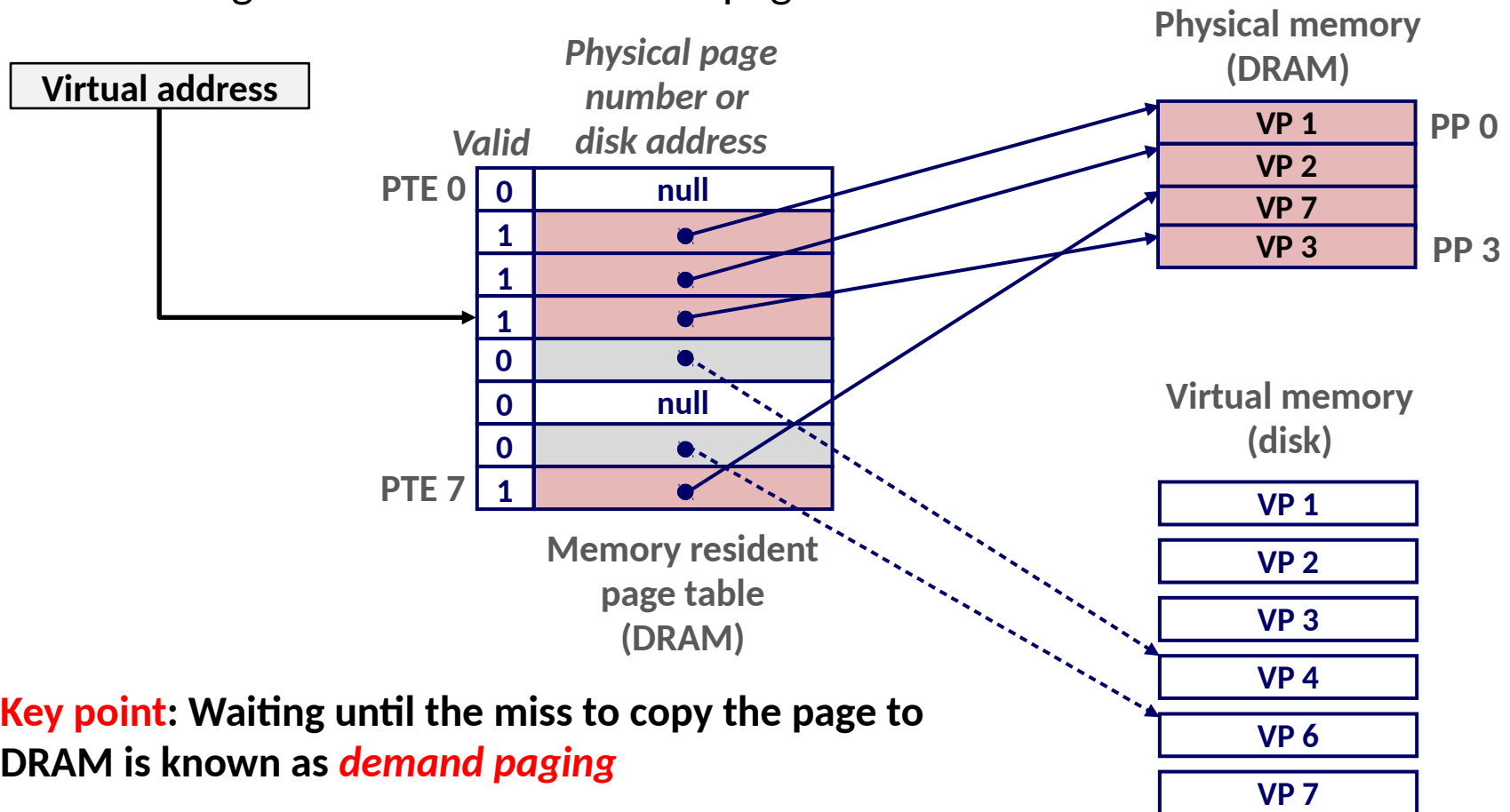- Page fault handler selects a victim to be evicted (here VP 4)

# Handling Page Fault

- Page miss causes page fault (an exception)
- Page fault handler selects a victim to be evicted (here VP 4)
- Offending instruction is restarted: page hit!



**Physical page number or disk address**

**Physical memory (DRAM)**

| | | |
|---|---|---|
| | VP 1 | PP 0 |
| | VP 2 | |
| | VP 7 | |
| | VP 3 | PP 3 |

**Virtual address**

*Valid*

| | | |
|---|---|---|
| PTE 0 | 0 | null |
| | 1 | |
| | 1 | |
| | 1 | |
| | 0 | |
| | 0 | null |
| | 0 | |
| PTE 7 | 1 | |

**Memory resident page table (DRAM)**

**Virtual memory (disk)**

| |
|---|
| VP 1 |
| VP 2 |
| VP 3 |
| VP 4 |
| VP 6 |
| VP 7 |

**Key point**: **Waiting until the miss to copy the page to DRAM is known as *demand paging***

50

# VM as a Tool for Memory Management

- **Key idea: each process has its own virtual address space**
  - It can view memory as a simple linear array
  - Mapping function scatters addresses through physical memory
    - Well-chosen mappings can improve locality

*Virtual Address Space for Process 1:*

*Address translation*

*Physical Address Space (DRAM)*

0

VP 1
VP 2
...
N-1

0

PP 2

PP 6

**(e.g., read-only library code)**

*Virtual Address Space for Process 2:*

0

VP 1
VP 2
...
N-1

PP 8

...

M-1

# VM as a Tool for Memory Protection

- **Extend PTEs with permission bits**
- **MMU checks these bits on each access**

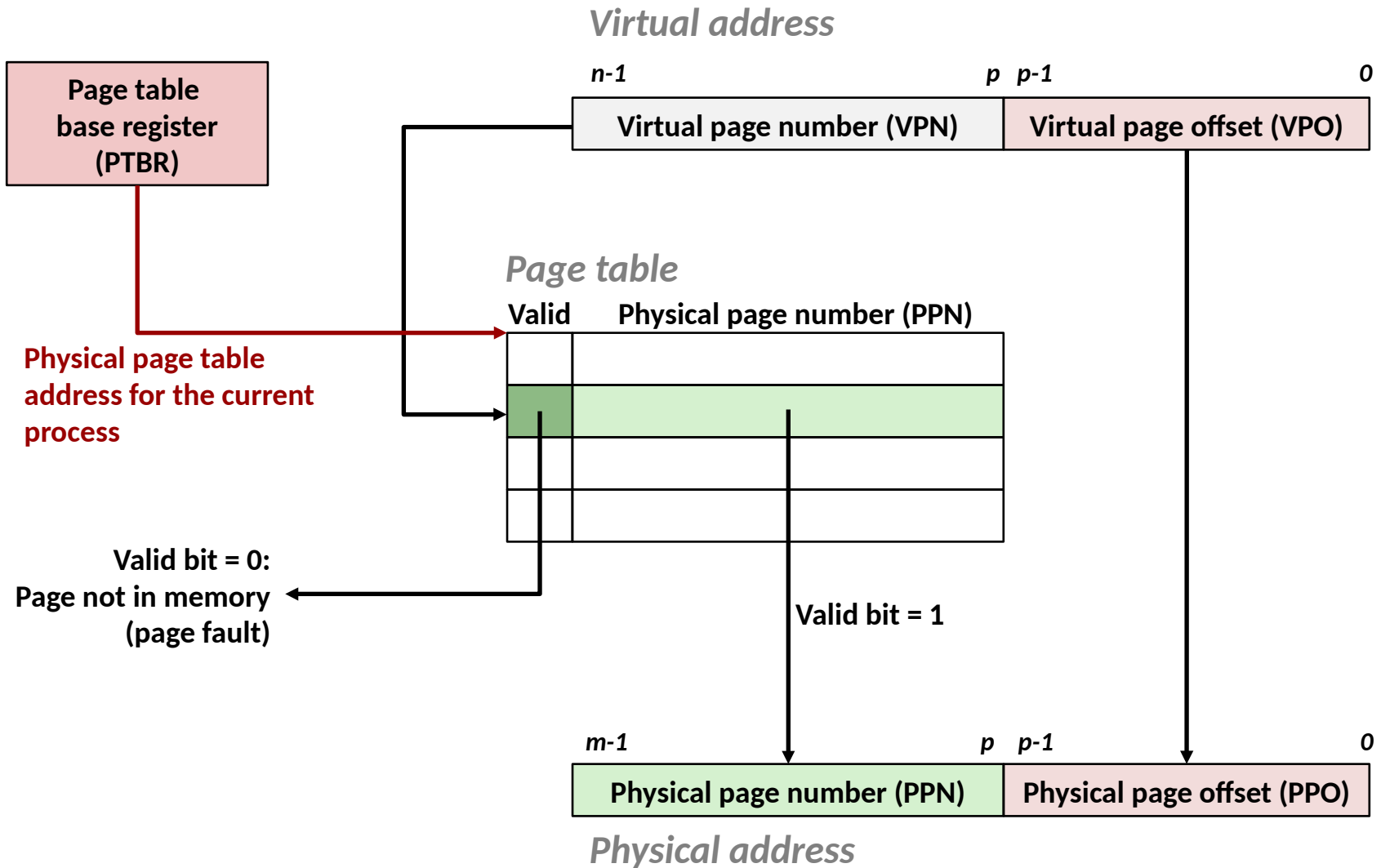**Physical Address Space**

**Process i:**

| | SUP | READ | WRITE | EXEC | Address |
|---|---|---|---|---|---|
| **VP 0:** | No | Yes | No | Yes | PP 6 |
| **VP 1:** | No | Yes | Yes | Yes | PP 4 |
| **VP 2:** | Yes | Yes | Yes | No | PP 2 |

**Process j:**

| | SUP | READ | WRITE | EXEC | Address |
|---|---|---|---|---|---|
| **VP 0:** | No | Yes | No | Yes | PP 9 |
| **VP 1:** | Yes | Yes | Yes | Yes | PP 6 |
| **VP 2:** | No | Yes | Yes | Yes | PP 11 |

| |
|---|
| |
| |
| PP 2 |
| |
| PP 4 |
| |
| PP 6 |
| |
| PP 8 |
| PP 9 |
| |
| PP 11 |

# Address Translation With a Page Table

*Virtual address*

| Page table base register (PTBR) |
|---|

*n-1*                          *p*   *p-1*                     *0*

| Virtual page number (VPN) | Virtual page offset (VPO) |
|---|---|

**Physical page table address for the current process**

*Page table*

**Valid**     **Physical page number (PPN)**

**Valid bit = 0:**
**Page not in memory**
**(page fault)**

**Valid bit = 1**

*m-1*                      *p*   *p-1*                      *0*

| Physical page number (PPN) | Physical page offset (PPO) |
|---|---|

*Physical address*
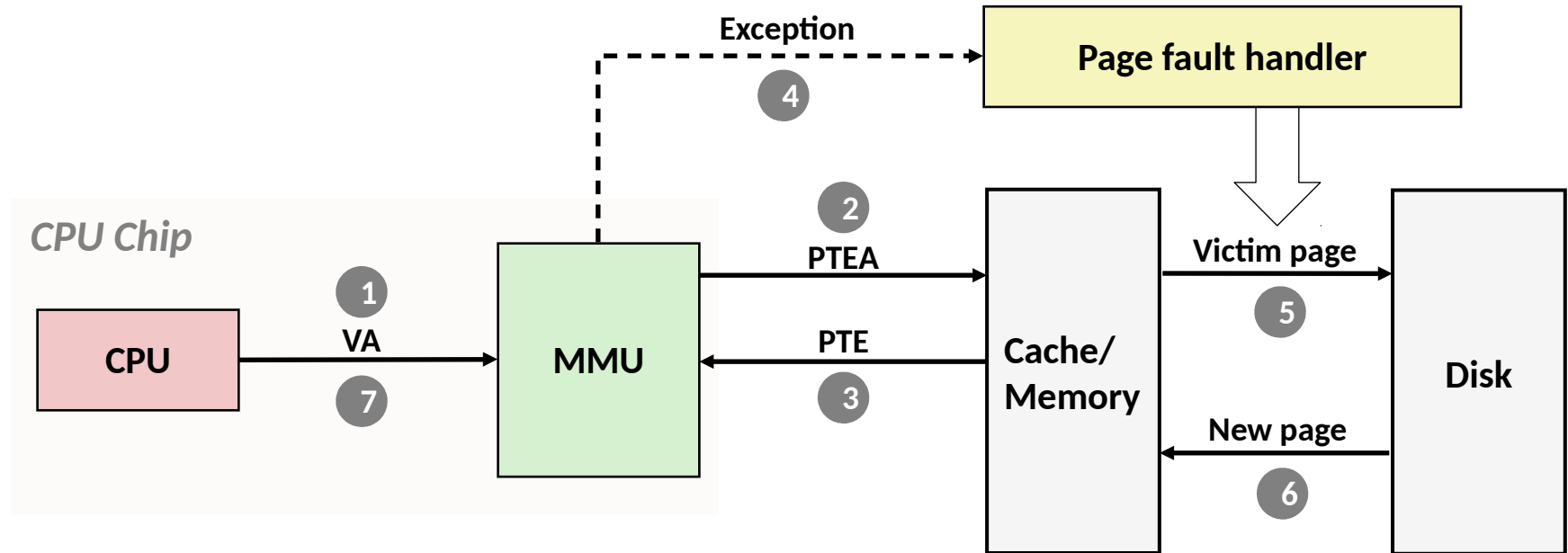
# Address Translation: Page Hit



1) Processor sends virtual address to MMU

2-3) MMU fetches PTE from page table in memory

4) MMU sends physical address to cache/memory

5) Cache/memory sends data word to processor

# Address Translation: Page Fault

**Exception**

**Page fault handler**

④

*CPU Chip*

②

**PTEA**

①

**VA**

**MMU**

**Victim page**

⑤

**PTE**

**Cache/Memory**

**Disk**

**CPU**

⑦

③

**New page**

⑥

1) Processor sends virtual address to MMU

2-3) MMU fetches PTE from page table in memory

4) Valid bit is zero, so MMU triggers page fault exception

5) Handler identifies victim (and, if dirty, pages it out to disk)

6) Handler pages in new page and updates PTE in memory

7) Handler returns to original process, restarting faulting instruction

# The `fork` Function Revisited

- **VM and memory mapping explain how `fork` provides private address space for each process.**

- **To create virtual address for new new process**
  - Create exact copies of current `mm_struct`, `vm_area_struct`, and page tables.
  - Flag each page in both processes as read-only
  - Flag each `vm_area_struct` in both processes as private COW

- **On return, each process has exact copy of virtual memory**

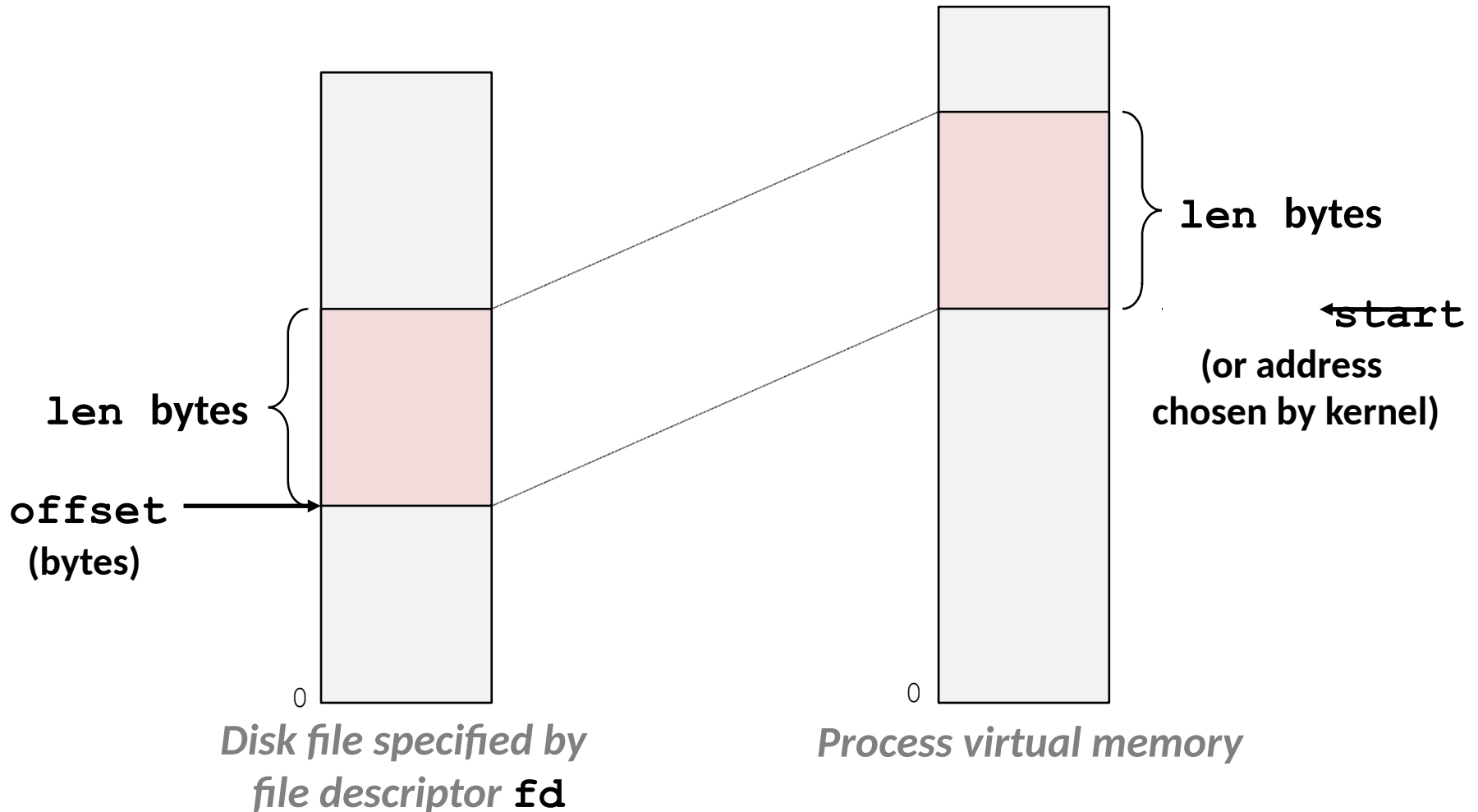- **Subsequent writes create new pages using COW mechanism.**

# User-Level Memory Mapping

```
void *mmap(void *start, int len,
           int prot, int flags, int fd, int offset)
```

- **Map `len` bytes starting at offset `offset` of the file specified by file description `fd`, preferably at address `start`**
  - **`start:`** may be 0 for "pick an address"
  - **`prot`**: PROT_READ, PROT_WRITE, …
  - **`flags`**: MAP_ANON, MAP_PRIVATE, MAP_SHARED, …

- **Return a pointer to start of mapped area (may not be `start`)**
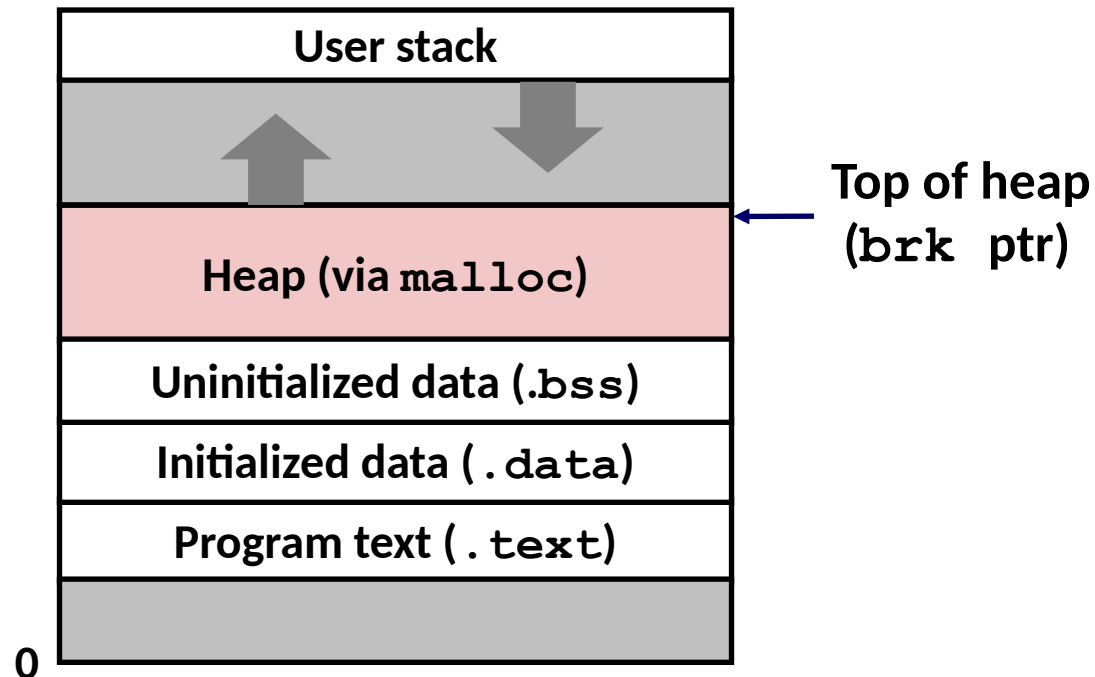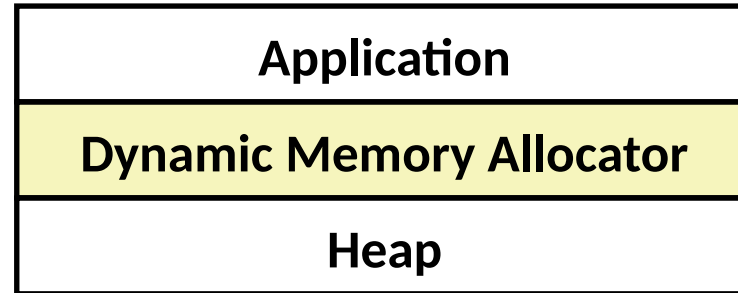
# User-Level Memory Mapping

```
void *mmap(void *start, int len,
           int prot, int flags, int fd, int offset)
```

**len bytes**

**len bytes**

**offset**
**(bytes)**

**start**

**(or address
chosen by kernel)**

0

0

*Disk file specified by
file descriptor* **fd**

*Process virtual memory*

# Summary: Dynamic Memory Allocation

# Dynamic Memory Allocation

- **Programmers use *dynamic memory allocators* (such as `malloc`) to acquire VM at run time.**

  - For data structures whose size is only known at runtime.

- **Dynamic memory allocators manage an area of process virtual memory known as the *heap*.**

| Application |
|:---:|
| **Dynamic Memory Allocator** |
| Heap |

| User stack |
|:---:|
| |
| Heap (via `malloc`) |
| Uninitialized data (`.bss`) |
| Initialized data (`.data`) |
| Program text (`.text`) |
| |

**Top of heap (`brk ptr`)**

0

# The `malloc` Package

`#include <stdlib.h>`

`void *malloc(size_t size)`
- Successful:
    - Returns a pointer to a memory block of at least `size` bytes aligned to an 8-byte (x86) or 16-byte (x86-64) boundary
    - If `size == 0`, returns NULL
- Unsuccessful: returns NULL (0) and sets `errno`

`void free(void *p)`
- Returns the block pointed at by `p` to pool of available memory
- `p` must come from a previous call to `malloc` or `realloc`

**Other functions**
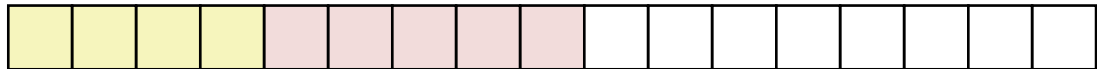- `calloc:` Version of `malloc` that initializes allocated block to zero.
- `realloc:` Changes the size of a previously allocated block.
- `sbrk:` Used internally by allocators to grow or shrink the heap
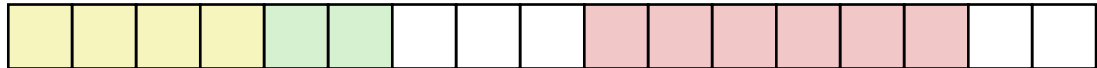
# Allocation Example

**p1 = malloc(4)**

**p2 = malloc(5)**

**p3 = malloc(6)**

**free(p2)**

**p4 = malloc(2)**

# Constraints

- **Applications**
  - Can issue arbitrary sequence of **malloc** and **free** requests
  - **free** request must be to a **malloc**'d block (or NULL)

- **Allocators**
  - Can't control number or size of allocated blocks
  - Must respond immediately to **malloc** requests
    - *i.e.*, can't reorder or buffer requests
  - Must allocate blocks from free memory
    - *i.e.*, can only place allocated blocks in free memory
  - Must align blocks so they satisfy all alignment requirements
    - 8-byte (x86) or 16-byte (x86-64) alignment on Linux boxes
  - Can manipulate and modify only free memory
  - Can't move the allocated blocks once they are **malloc**'d
    - *i.e.*, compaction is not allowed (why?)

# Performance Goal: Throughput

- **Given some sequence of `malloc` and `free` requests:**
  - $R_0, R_1, ..., R_k, ..., R_{n-1}$

- **Goals: maximize throughput and peak memory utilization**
  - These goals are often conflicting

- **Throughput:**
  - Number of completed requests per unit time
  - Example:
    - 5,000 `malloc` calls and 5,000 `free` calls in 10 seconds
    - Throughput is 1,000 operations/second
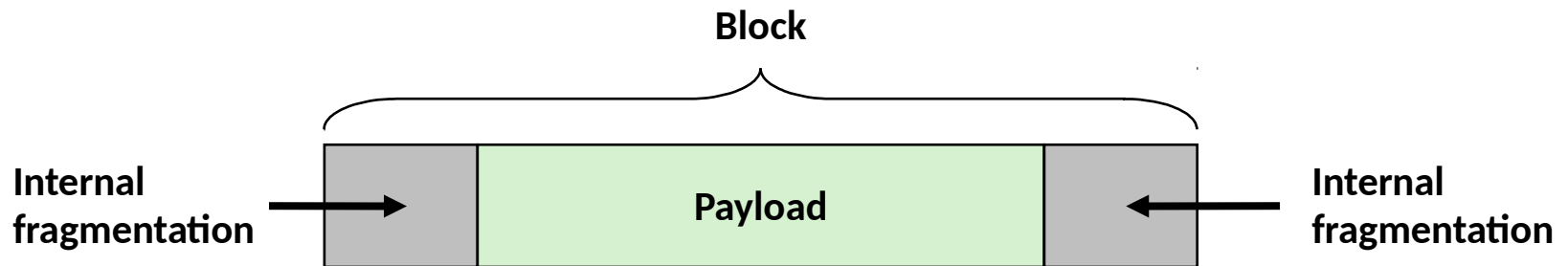
# Performance Goal: Peak Memory Utilization

- **Given some sequence of `malloc` and `free` requests:**
  - $R_0, R_1, ..., R_k, ... , R_{n-1}$

- *Def: Aggregate payload $P_k$*
  - `malloc(p)` results in a block with a ***payload*** of `p` bytes
  - After request $R_k$ has completed, the ***aggregate payload*** $P_k$ is the sum of currently allocated payloads

- *Def: Current heap size $H_k$*
  - Assume $H_k$ is monotonically nondecreasing
    - i.e., heap only grows when allocator uses `sbrk`

- *Def: Peak memory utilization after k+1 requests*
  - $U_k = ( max_{i<=k} P_i ) / H_k$

# Fragmentation

- **Poor memory utilization caused by *fragmentation***
  - *internal* fragmentation
  - *external* fragmentation

# Internal Fragmentation

- **For a given block, *internal fragmentation* occurs if payload is smaller than block size**

**Block**

**Internal fragmentation** →          **Payload**          ← **Internal fragmentation**

- **Caused by**
  - Overhead of maintaining heap data structures
  - Padding for alignment purposes
  - Explicit policy decisions
    (e.g., to return a big block to satisfy a small request)

- **Depends only on the pattern of *previous* requests**
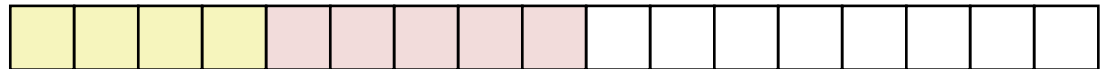  - Thus, easy to measure

# External Fragmentation

- **Occurs when there is enough aggregate heap memory, but no single free block is large enough**
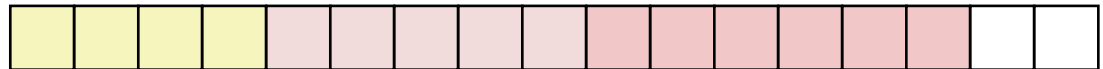
`p1 = malloc(4)`

`p2 = malloc(5)`

`p3 = malloc(6)`

`free(p2)`

`p4 = malloc(6)`   *Oops! (what would happen now?)*

- **Depends on the pattern of future requests**
  - Thus, difficult to measure

# Implementation Issues

- **How do we know how much memory to free given just a pointer?**

- **How do we keep track of the free blocks?**

- **What do we do with the extra space when allocating a structure that is smaller than the free block it is placed in?**

- **How do we pick a block to use for allocation -- many might fit?**

- **How do we reinsert freed block?**

# Knowing How Much to Free

- **Standard method**
    - Keep the length of a block in the word preceding the block.
        - This word is often called the ***header field*** or ***header***
    - Requires an extra word for every allocated block



p0

p0 = malloc(4)  | | | | | | | | | 5 | | | | | | |

block size        payload

free(p0)

# Keeping Track of Free Blocks

- **Method 1: *Implicit list* using length—links all blocks**



- **Method 2: *Explicit list* among the free blocks using pointers**



- **Method 3: *Segregated free list***
  - Different free lists for different size classes

- **Method 4: *Blocks sorted by size***
  - Can use a balanced tree (e.g. Red-Black tree) with pointers within each free block, and the length used as a key

# Summary: Concurrent Programming is Hard!

# Concurrent Programming is Hard!

- **The human mind tends to be sequential**

- **The notion of time is often misleading**

- **Thinking about all possible sequences of events in a computer system is at least error prone and frequently impossible**

# Concurrent Programming is Hard!

- **Classical problem classes of concurrent programs:**
  - *Races:* outcome depends on arbitrary scheduling decisions elsewhere in the system
    - Example: who gets the last seat on the airplane?
  - *Deadlock:* improper resource allocation prevents forward progress
    - Example: traffic gridlock
  - *Livelock / Starvation / Fairness*: external events and/or system scheduling decisions can prevent sub-task progress
    - Example: people always jump in front of you in line

# A Process With Multiple Threads

- **Multiple threads can be associated with a process**
  - Each thread has its own logical control flow
  - Each thread shares the same code, data, and kernel context
  - Each thread has its own stack for local variables
    - but not protected from other threads
  - Each thread has its own thread id (TID)

**Thread 1 (main thread)**

**Thread 2 (peer thread)**

**Shared code and data**

| stack 1 |
|---|

| Thread 1 context: |
|---|
| Data registers |
| Condition codes |
| SP1 |
| PC1 |

| stack 2 |
|---|

| Thread 2 context: |
|---|
| Data registers |
| Condition codes |
| SP2 |
| PC2 |

| shared libraries |
|---|
| |
| run-time heap |
| read/write data |
| read-only code/data |
| |

0

| Kernel context: |
|---|
| VM structures |
| Descriptor table |
| brk pointer |

# Logical View of Threads

- **Threads associated with process form a pool of peers**
    - Unlike processes which form a tree hierarchy

**Threads associated with process foo**

**Process hierarchy**

# Concurrent Thread Execution

- **Single Core Processor**
  - Simulate parallelism by time slicing

- **Multi-Core Processor**
  - Can have true parallelism

Thread A    Thread B    Thread C

Time

Thread A    Thread B    Thread C

**Run 3 threads on 2 cores**

# Threads vs. Processes

- **How threads and processes are similar**
  - Each has its own logical control flow
  - Each can run concurrently with others (possibly on different cores)
  - Each is context switched

- **How threads and processes are different**
  - Threads share all code and data (except local stacks)
    - Processes (typically) do not
  - Threads are somewhat less expensive than processes
    - Process control (creating and reaping) twice as expensive as thread control
    - Linux numbers:
      - ~20K cycles to create and reap a process
      - ~10K cycles (or less) to create and reap a thread
      - *Much* larger difference on non-Unices.

# Shared Variables in Threaded C Programs

- **Question: Which variables  in a threaded C program are shared among threads?**
  - The answer is not as simple as "*global variables are shared*" and "*stack variables are private*"

- ***Def*: A variable `x` is *shared* if and only if multiple threads reference some instance of `x`.**

- **Requires answers to the following questions:**
  - What is the memory model for threads?
  - How are instances of variables mapped to memory?
  - How many threads might reference each of these instances?

# Threads Memory Model

- **Conceptual model:**
  - Multiple threads run within the context of a single process
  - Each thread has its own separate thread context
    - Thread ID, stack, stack pointer, PC, condition codes, and GP registers
  - All threads share the remaining process context
    - Code, data, heap, and shared library segments of the process virtual address space
    - Open files and installed handlers
- **Operationally, this model is not strictly enforced:**
  - Register values are truly separate and protected, but...
  - Any thread can read and write the stack of any other thread

*The mismatch between the conceptual and operation model is a source of confusion and errors*

# Example Program to Illustrate Sharing

```c
char **ptr;  /* global var */

int main()
{
    long i;
    pthread_t tid;
    char *msgs[2] = {
        "Hello from foo",
        "Hello from bar"
    };

    ptr = msgs;
    for (i = 0; i < 2; i++)
        Pthread_create(&tid,
            NULL,
            thread,
            (void *)i);
    Pthread_exit(NULL);
}
```
sharing.c

```c
void *thread(void *vargp)
{
    long myid = (long)vargp;
    static int cnt = 0;

    printf("[%ld]:  %s (cnt=%d)\n",
            myid, ptr[myid], ++cnt);
    return NULL;
}
```

*Peer threads reference main thread's stack indirectly through global ptr variable*

# Mapping Variable Instances to Memory

- **Global variables**
  - *Def:* Variable declared outside of a function
  - **Virtual memory contains exactly one instance of any global variable**

- **Local variables (including *thread-local variables*)**
  - *Def:* Variable declared inside function without `static` attribute.
    - Or global variable with `__thread` in GCC.
  - **Each thread stack contains one instance of each local variable**

- **Local static variables**
  - *Def:* Variable declared inside function with the `static` attribute
  - **Virtual memory contains exactly one instance of any local static variable.**

# Synchronizing Threads

- **Shared variables are handy…**

- **…but introduce the possibility of nasty *synchronization* errors.**

# `badcnt.c`: Improper Synchronization

```c
/* Global shared variable */
volatile long cnt = 0; /* Counter */

int main(int argc, char **argv)
{
    long niters;
    pthread_t tid1, tid2;

    niters = atoi(argv[1]);
    Pthread_create(&tid1, NULL,
        thread, &niters);
    Pthread_create(&tid2, NULL,
        thread, &niters);
    Pthread_join(tid1, NULL);
    Pthread_join(tid2, NULL);

    /* Check result */
    if (cnt != (2 * niters))
        printf("BOOM! cnt=%ld\n", cnt);
    else
        printf("OK cnt=%ld\n", cnt);
    exit(0);
}
```
badcnt.c

```c
/* Thread routine */
void *thread(void *vargp)
{
    long i, niters =
            *((long *)vargp);

    for (i = 0; i < niters; i++)
        cnt++;

    return NULL;
}
```

```
linux> ./badcnt 10000
OK cnt=20000
linux> ./badcnt 10000
BOOM! cnt=13051
linux>
```

**`cnt` should equal 20,000.**

**What went wrong?**

# Assembly Code for Counter Loop

**C code for counter loop in thread i**

```
for (i = 0; i < niters; i++)
    cnt++;
```

*Asm code for thread i*

```
        movq   (%rdi), %rcx
        testq %rcx,%rcx
        jle    .L2
        movl   $0, %eax
.L3:
        movq   cnt(%rip),%rdx
        addq   $1, %rdx
        movq   %rdx, cnt(%rip)
        addq   $1, %rax
        cmpq   %rcx, %rax
        jne    .L3
.L2:
```

$H_i$ : Head

$L_i$ : Load `cnt`
$U_i$ : Update `cnt`
$S_i$ : Store `cnt`

$T_i$ : Tail

# Semaphores

- *Semaphore:* **non-negative global integer synchronization variable. Manipulated by *P (passering)* and *V (vrijgave)* erations.**

- **P(s):**
  - If *s* is nonzero, then decrement *s* by 1 and return immediately.
    - Test and decrement operations occur atomically (indivisibly)
  - If *s* is zero, then suspend thread until *s* becomes nonzero and the thread is restarted by a V operation.
  - After restarting, the P operation decrements *s* and returns control to the caller.

- *V(s):*
  - Increment *s* by 1.
    - Increment operation occurs atomically
  - If there are any threads blocked in a P operation waiting for *s* to become non-zero, then restart exactly one of those threads, which then completes its P operation by decrementing *s*.

- **Semaphore invariant:** *(s >= 0)*

# C Semaphore Operations

**Pthreads functions:**

```
#include <semaphore.h>

int sem_init(sem_t *s, 0, unsigned int val);} /* s = val */

int sem_wait(sem_t *s);   /* P(s) */
int sem_post(sem_t *s);   /* V(s) */
```

**CS:APP wrapper functions:**

```
#include "csapp.h"

void P(sem_t *s); /* Wrapper function for sem_wait */
void V(sem_t *s); /* Wrapper function for sem_post */
```

# Using Semaphores for Mutual Exclusion

- **Basic idea:**
  - Associate a unique semaphore *mutex*, initially 1, with each shared variable (or related set of shared variables).
  - Surround corresponding critical sections with *P(mutex)* and *V(mutex)* operations.

- **Terminology:**
  - *Binary semaphore*: semaphore whose value is always 0 or 1
  - *Mutex:* binary semaphore used for mutual exclusion
    - P operation: "locking" the mutex
    - V operation: "unlocking" or "releasing" the mutex
    - *"Holding"* a mutex: locked and not yet unlocked.
  - *Counting semaphore*: used as a counter for set of available resources.

# `goodcnt.c`: Proper Synchronization

■ **Define and initialize a mutex for the shared variable `cnt`:**

```
volatile long cnt = 0;    /* Counter */
sem_t mutex;              /* Semaphore that protects cnt */

Sem_init(&mutex, 0, 1); /* mutex = 1 */
```

■ **Surround critical section with *P* and *V*:**

```
for (i = 0; i < niters; i++) {
    P(&mutex);
    cnt++;
    V(&mutex);
}
                              goodcnt.c
```

```
linux> ./goodcnt 10000
OK cnt=20000
linux> ./goodcnt 10000
OK cnt=20000
linux>
```

**Warning: It's orders of magnitude slower than `badcnt.c`.**

# Summary: Unix I/O

# Unix I/O Overview

- **A Linux *file* is a sequence of *m* bytes:**
  - $B_0$ , $B_1$ , .... , $B_k$ , .... , $B_{m-1}$

- **Cool fact: All I/O devices are represented as files:**
  - **/dev/sda2**   (**/usr** disk partition)
  - **/dev/tty2**   (terminal)

- **Even the kernel is represented as a file:**
  - **/boot/vmlinuz-3.13.0-55-generic**  (kernel image)
  - **/proc**                                    (kernel data structures)
  - **/sys**                              (other kernel data structures)

# Unix I/O Overview

- **Elegant mapping of files to devices allows kernel to export simple interface called *Unix I/O:***
  - Opening and closing files
    - **`open()`** and **`close()`**
  - Reading and writing a file
    - **`read()`** and **`write()`**
  - Changing the *current file position* (seek)
    - indicates next offset into file to read or write
    - **`lseek()`**

| $B_0$ | $B_1$ | ● ● ● | $B_{k-1}$ | $B_k$ | $B_{k+1}$ | ● ● ● |

**Current file position = k**

# File Types

- **Each file has a *type* indicating its role in the system**
    - *Regular file:* Contains arbitrary data
    - *Directory:*  Index for a related group of files
    - *Socket:* For communicating with a process on another machine

- **Other file types beyond our scope**
    - *Named pipes (FIFOs)*
    - *Symbolic links*
    - *Character and block devices*

# Regular Files

- **A regular file contains arbitrary data**
- **Applications often distinguish between *text files* and *binary files***
  - Text files are regular files with only ASCII or Unicode characters
  - Binary files are everything else
    - e.g., object files, JPEG images
  - Kernel doesn't know the difference!
- **Text file is sequence of *text lines***
  - Text line is sequence of chars terminated by *newline char* ('\n')
    - Newline is 0xa, same as ASCII line feed character (LF)
- **End of line (EOL) indicators in other systems**
  - Linux and Mac OS: '\n' (0xa)
    - line feed (LF)
  - Windows and Internet protocols: '\r\n' (0xd 0xa)
    - Carriage return (CR) followed by line feed (LF)



carriage return

line feed

# Directories

- **Directory consists of an array of *links***
  - Each link maps a *filename* to a file
- **Each directory contains at least two entries**
  - `.` (dot) is  a link to itself
  - `..` (dot dot) is a link to *the parent directory* in the *directory hierarchy* (next slide)
- **Commands for manipulating directories**
  - `mkdir`: create empty directory
  - `ls`: view directory contents
  - `rmdir`: delete empty directory

# Opening Files

- **Opening a file informs the kernel that you are getting ready to access that file**

```
int fd;    /* file descriptor */

if ((fd = open("/etc/hosts", O_RDONLY)) < 0) {
    perror("open");
    exit(1);
}
```

- **Returns a small identifying integer *file descriptor***
  - **`fd == -1`** indicates that an error occurred

- **Each process created by a Linux shell begins life with three open files associated with a terminal:**
  - 0: standard input (stdin)
  - 1: standard output (stdout)
  - 2: standard error (stderr)

# Closing Files

- **Closing a file informs the kernel that you are finished accessing that file**

```
int fd;      /* file descriptor */
int retval; /* return value */

if ((retval = close(fd)) < 0) {
   perror("close");
   exit(1);
}
```

- **Closing an already closed file is a recipe for disaster in threaded programs**
- **Moral: Always check return codes, even for seemingly benign functions such as `close()`**

# Reading Files

- **Reading a file copies bytes from the current file position to memory, and then updates file position**

```
char buf[512];
int fd;          /* file descriptor */
int nbytes;      /* number of bytes read */

/* Open file fd ...  */
/* Then read up to 512 bytes from file fd */
if ((nbytes = read(fd, buf, sizeof(buf))) < 0) {
    perror("read");
    exit(1);
}
```

- **Returns number of bytes read from file `fd` into `buf`**
  - Return type `ssize_t` is signed integer
  - `nbytes < 0` indicates that an error occurred
  - *Short counts* (`nbytes < sizeof(buf)` ) are possible and are not errors!

# On Short Counts

- **Short counts often occurs in these situations:**
  - Encountering (end-of-file) EOF on reads
  - Reading text lines from a terminal
  - Reading and writing network sockets

- **Short counts rarely occurs in these situations:**
  - Reading from disk files (except for EOF)
    - …but may happen for huge reads.
  - Writing to disk files
    - …similarly.

- **Best practice is to always allow for short counts.**

# Implementation of `rio_readn`

```
/*
 * rio_readn - Robustly read n bytes (unbuffered)
 */
ssize_t rio_readn(int fd, void *usrbuf, size_t n)
{
    size_t nleft = n;
    ssize_t nread;
    char *bufp = usrbuf;

    while (nleft > 0) {
        if ((nread = read(fd, bufp, nleft)) < 0) {
            if (errno == EINTR) /* Interrupted by sig handler return */
                nread = 0;      /* and call read() again */
            else
                return -1;      /* errno set by read() */
        }
        else if (nread == 0)
            break;              /* EOF */
        nleft -= nread;
        bufp += nread;
    }
    return (n - nleft);      /* Return >= 0 */
}
```

csapp.c

# Buffered I/O: Motivation

- **Applications often read/write one character at a time**
  - `getc, putc, ungetc`
  - `gets, fgets`
    - Read line of text one character at a time, stopping at newline
- **Implementing as Unix I/O calls expensive**
  - `read` and `write` require Unix kernel calls
    - > 10,000 clock cycles
- **Solution: Buffered read**
  - Use Unix `read` to grab block of bytes
  - User input functions take one byte at a time from buffer
    - Refill buffer when empty

| *Buffer* | **already read** | **unread** | |
|---|---|---|---|

# Buffering in Standard I/O

- **Standard I/O functions use buffered I/O**

```
                printf("h");
                  printf("e");
                    printf("l");
                      printf("l");
                        printf("o");
                          printf("\n");
 buf
```

| h | e | l | l | o | \n | . | . |

```
                  fflush(stdout);


              write(1, buf, 6);
```

- **Buffer flushed to output fd on "\n", call to `fflush` or `exit`, or return from `main`.**

# Buffered I/O: Declaration

■ **All information contained in `struct`**



```
typedef struct {
    int rio_fd;                 /* descriptor for this internal buf */
    int rio_cnt;                /* unread bytes in internal buf */
    char *rio_bufptr;           /* next unread byte in internal buf */
    char rio_buf[RIO_BUFSIZE];  /* internal buffer */
} rio_t;
```

# Buffered I/O: Read some bytes

```c
static ssize_t rio_read(rio_t *rp, char *usrbuf, size_t n)
{
    int cnt;

    while (rp->rio_cnt <= 0) {  /* Refill if buf is empty */
        rp->rio_cnt = read(rp->rio_fd, rp->rio_buf,
                            sizeof(rp->rio_buf));
        if (rp->rio_cnt < 0) {
            if (errno != EINTR) /* Interrupted by sig handler return */
                return -1;
        }
        else if (rp->rio_cnt == 0)  /* EOF */
            return 0;
        else
            rp->rio_bufptr = rp->rio_buf; /* Reset buffer ptr */
    }

    /* Copy min(n, rp->rio_cnt) bytes from internal buf to user buf */
    cnt = n;
    if (rp->rio_cnt < n)
        cnt = rp->rio_cnt;
    memcpy(usrbuf, rp->rio_bufptr, cnt);
    rp->rio_bufptr += cnt;
    rp->rio_cnt -= cnt;
    return cnt;
}
```

`csapp.c`

# Buffered I/O: Read *n* bytes robustly

```c
ssize_t rio_readnb(rio_t *rp, void *usrbuf, size_t n)
{
    size_t nleft = n;
    ssize_t nread;
    char *bufp = usrbuf;

    while (nleft > 0) {
        if ((nread = rio_read(rp, bufp, nleft)) < 0)
            return -1;          /* errno set by read() */
        else if (nread == 0)
            break;              /* EOF */
        nleft -= nread;
        bufp += nread;
    }
    return (n - nleft);         /* return >= 0 */
}
                                              csapp.c
```
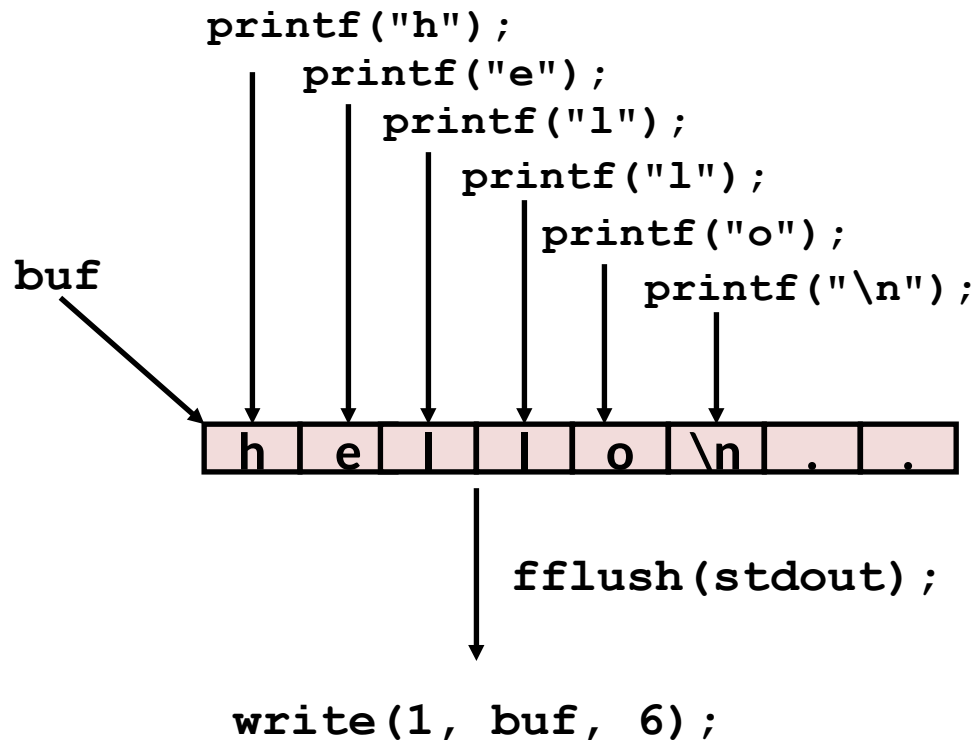
# File Metadata

- *Metadata* is data about data, in this case file data
- **Per-file metadata maintained by kernel**
  - accessed by users with the `stat` and `fstat` functions

```
/* Metadata returned by the stat and fstat functions */
struct stat {
    dev_t          st_dev;        /* Device */
    ino_t          st_ino;        /* inode */
    mode_t         st_mode;       /* Protection and file type */
    nlink_t        st_nlink;      /* Number of hard links */
    uid_t          st_uid;        /* User ID of owner */
    gid_t          st_gid;        /* Group ID of owner */
    dev_t          st_rdev;       /* Device type (if inode device) */
    off_t          st_size;       /* Total size, in bytes */
    unsigned long st_blksize;     /* Blocksize for filesystem I/O */
    unsigned long st_blocks;      /* Number of blocks allocated */
    time_t         st_atime;      /* Time of last access */
    time_t         st_mtime;      /* Time of last modification */
    time_t         st_ctime;      /* Time of last change */
};
```

# How the Unix Kernel Represents Open Files

- **Two descriptors referencing two distinct open files. Descriptor 1 (stdout) points to terminal, and descriptor 4 points to open disk file**



Descriptor table [one table per process]

Open file table [shared by all processes]

v-node table [shared by all processes]

stdin  fd 0
stdout fd 1
stderr fd 2
fd 3
fd 4

File A (terminal)

File pos
refcnt=1

File access
File size
File type

Info in stat struct

File B (disk)

File pos
refcnt=1

File access
File size
File type

# File Sharing

- **Two distinct descriptors sharing the same disk file through two distinct open file table entries**
  - E.g., Calling `open` twice with the same `filename` argument

**Descriptor table**
**[one table per process]**

**Open file table**
**[shared by all processes]**

**v-node table**
**[shared by all processes]**

**File A (disk)**

stdin    fd 0
stdout   fd 1
stderr   fd 2
         fd 3
         fd 4

**File pos**
**refcnt=1**

**File access**
**File size**
**File type**

**File B (disk)**

**File pos**
**refcnt=1**

# How Processes Share Files: `fork`

- **A child process inherits its parent's open files**
  - Note: situation unchanged by **exec** functions (use **fcntl** to change)
- **_Before_ fork call:**

**Descriptor table**
**[one table per process]**

**Open file table**
**[shared by all processes]**

**v-node table**
**[shared by all processes]**

| | | |
|---|---|---|
| stdin | fd 0 | |
| stdout | fd 1 | |
| stderr | fd 2 | |
| | fd 3 | |
| | fd 4 | |

**File A (terminal)**

File pos

refcnt=1

**File B (disk)**

File pos

refcnt=1

**File access**

**File size**

**File type**

**File access**

**File size**

**File type**

# How Processes Share Files: `fork`

- **A child process inherits its parent's open files**

- *After* `fork`:
  - Child's table same as parent's, and +1 to each refcnt

**Descriptor table**
**[one table per process]**

**Open file table**
**[shared by all processes]**

**v-node table**
**[shared by all processes]**

**Parent**

| fd 0 |
| fd 1 |
| fd 2 |
| fd 3 |
| fd 4 |

**Child**

| fd 0 |
| fd 1 |
| fd 2 |
| fd 3 |
| fd 4 |

**File A (terminal)**

| |
| **File pos** |
| `refcnt=2` |
| |

**File B (disk)**

| |
| **File pos** |
| `refcnt=2` |
| |

| **File access** |
| **File size** |
| **File type** |
| |

| **File access** |
| **File size** |
| **File type** |
| |

# I/O Redirection

■ **Question: How does a shell implement I/O redirection?**
```
linux> ls > foo.txt
```

■ **Answer: By calling the `dup2(oldfd, newfd)` function**
  ▪ Copies (per-process) descriptor table entry `oldfd` to entry `newfd`

**Descriptor table**
*before* `dup2(4,1)`

| fd 0 |     |
|------|-----|
| fd 1 | a   |
| fd 2 |     |
| fd 3 |     |
| fd 4 | b   |

**Descriptor table**
*after* `dup2(4,1)`

| fd 0 |     |
|------|-----|
| fd 1 | b   |
| fd 2 |     |
| fd 3 |     |
| fd 4 | b   |

# I/O Redirection Example

- **Step #1: open file to which stdout should be redirected**
  - Happens in child executing shell code, before `exec`

**Descriptor table**
**[one table per process]**

**Open file table**
**[shared by all processes]**

**v-node table**
**[shared by all processes]**



File A

| stdin | fd 0 |
| stdout | fd 1 |
| stderr | fd 2 |
| | fd 3 |
| | fd 4 |

File pos

refcnt=1

File access

File size

File type

File B

File pos

refcnt=1

File access

File size

File type

# I/O Redirection Example (cont.)

- **Step #2: call `dup2(4,1)`**
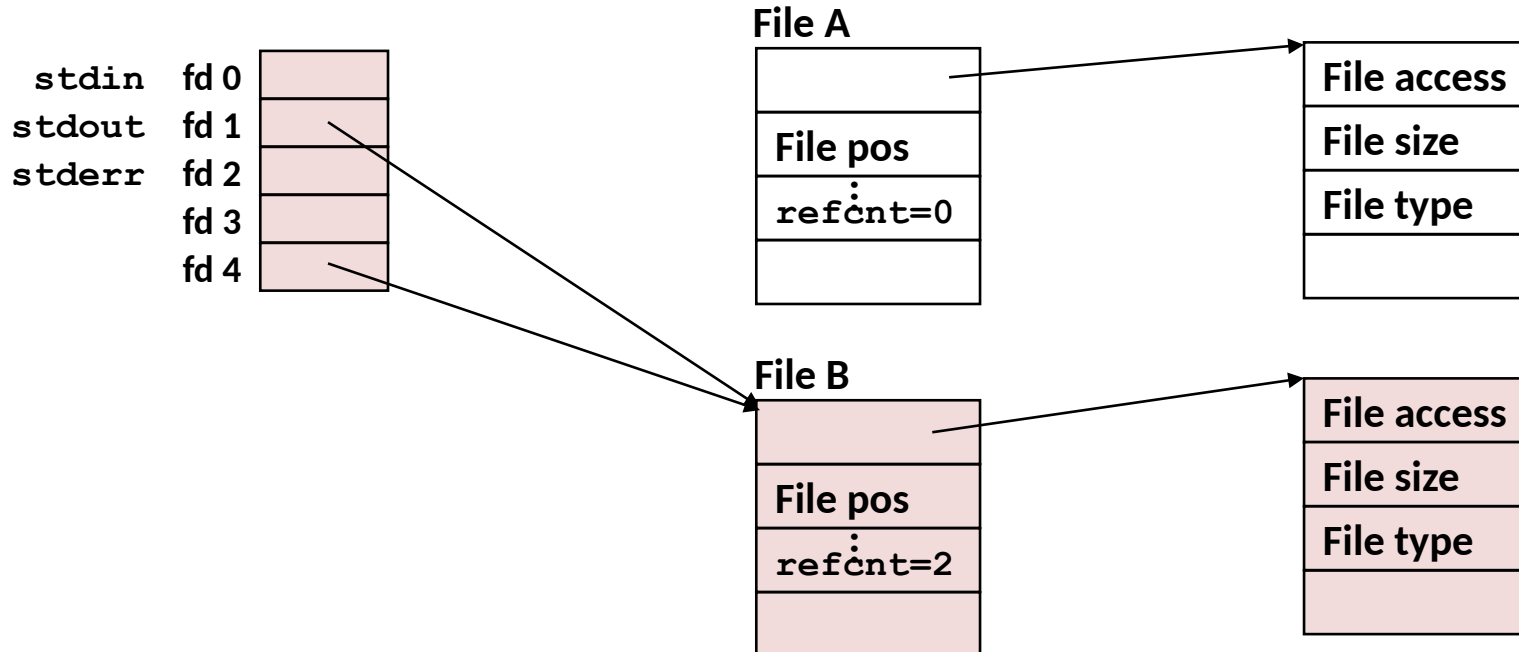  - cause fd=1 (stdout) to refer to disk file pointed at by fd=4

**Descriptor table**
**[one table per process]**

**Open file table**
**[shared by all processes]**

**v-node table**
**[shared by all processes]**

File A

| | |
|---|---|
| `stdin` fd 0 | |
| `stdout` fd 1 | |
| `stderr` fd 2 | |
| fd 3 | |
| fd 4 | |

File A:
| |
|---|
| |
| **File pos** |
| **refcnt=0** |
| |

| |
|---|
| **File access** |
| **File size** |
| **File type** |
| |

File B:
| |
|---|
| |
| **File pos** |
| **refcnt=2** |
| |

| |
|---|
| **File access** |
| **File size** |
| **File type** |
| |

# Unix I/O vs. Standard I/O vs. RIO

- **Standard I/O and RIO are implemented using low-level Unix I/O**

```
fopen   fdopen
fread   fwrite
fscanf  fprintf
 sscanf
sprintf fgets
fputs   fflush
fseek
fclose
```

```
open    read
write   lseek
stat    close
```

**C application program**

Standard I/O
functions

RIO
functions

Unix I/O functions
(accessed via system calls)

```
rio_readn
rio_writen
rio_readinitb
rio_readlineb
rio_readnb
```

- **Which ones should you use in your programs?**

# Choosing I/O Functions

- **General rule: use the highest-level I/O functions you can**
  - Many C programmers are able to do all of their work using the standard I/O functions
  - But, be sure to understand the functions you use!

- **When to use standard I/O**
  - When working with disk or terminal files

- **When to use raw Unix I/O**
  - Inside signal handlers, because Unix I/O is async-signal-safe
  - In rare cases when you need absolute highest performance

- **When to use RIO**
  - When you are reading and writing network sockets
  - Avoid using standard I/O on sockets

# Goodbye for now

- **This was my last lecture this year. Final advice:**
  - Check your return codes.
  - Use `assert()` liberally.
  - Use a better language than C if at all possible.
  - If such a language does not exist, *invent it*.
  - No, C++ is not that language.

- **If you like this kind of stuff, there are two courses you should take:**
  - *Programming Massively Parallel Hardware* (PMPH), in block 1.
  - *Parallel Functional Programming* in block 2 (may change its name next year).
  - (Nominally master's courses, but don't let that stop you; it is a moral imperative to disobey bad rules/guidelines.)
  - Also, check out my research: https://futhark-lang.org