# Operating Systems Overview

# operating system

- no clear traditional definition
  - each definition cover a distinct aspect

- an interface between applications and hardware
  - true, this was the first reason for having an OS
- a set of programs that provides basic functionalities for managing system resources
  - true, but the OS is not only a "library" of functionalities
- a program that "controls" the execution of application programs
  - e.g. it decides which program is running and when

# my definition

- a software that does not **depend** on any other software in the computer
  - do you know what "depend" means?
  - in a certain sense it only depends on hardware

in current PC and servers it also is...

- a software that...
  - considers "activities" to be carried on and...
  - assigns "resources" to them
    - "activities" usually means "processes"
- the latter is the most interesting definition

# Operating System as a Resource Manager

- **resource**: anything needed for program execution, e.g.
  - cpu time
  - I/O devices
  - memory
  - executable code
  - etc.
- an os manages...
  - resources
    - comprising "internal" data structures
  - processes

# my definition

if the hardware is so small that cannot carry on many tasks (e.g. a mobile phone)...

- – less and less frequent

- ...it makes simple for applicative software to interface with hardware

- hw interface on bigger systems is not the primary concern

- – handled with "drivers"

# kernel

- its the OS as defined above
- it is always in main memory
- contains most frequently used functionality
- also called "nucleus"

# kernel

- a very interesting piece of software!
- can be studied from several point of views
  - structural
    - which concepts realize?
    - which data structures adopt?
    - what modules it contains?
    - what are the interfaces among modules?
  - behavioral
    - which strategies adopt to optimize usage resources?
    - which algorithms it implements to maintain its data structures?
  - synchronization
    - oops! many things happens in the kernel at the same time. Consistency is hard to maintain. Large number of techniques uses. Completeley indpendent by OS study, but largely developed by the same community.

# kernel

- didactic approach vs. reality
  - structural
    - what is the easiest way to realize a "decent" kernel?
    - what is there in a real kernel?
  - behavioral
    - what are the easiest strategies?
    - what are the ones that are adopted in practice?
  - synchronization
    - too complicated... there is a specific course for this ;)
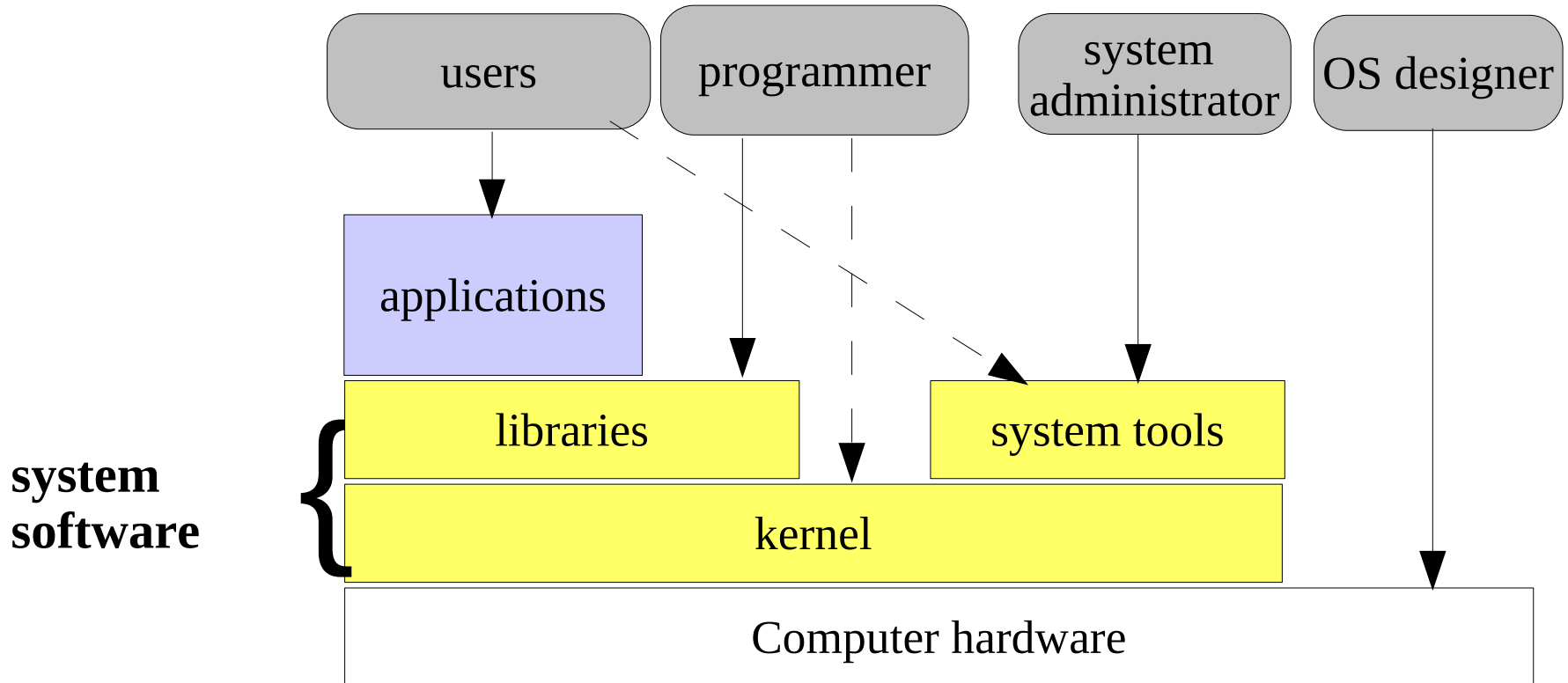
# system software

- it is the software that is usually installed when the OS is installed on a computer system

  – e.g.
    basic C runtime libraries,
    sw to access the hard disk directory,
    graphical user interface,
    etc.

9

# overloaded language

- actually "operating system" is an overloaded word

- it means
  - the **kernel** or
  - the kernel **plus** the other system software

- meaning should be clear from the context

# computer system layers

not on
the book

users    programmer    system administrator    OS designer

applications

**system software**

libraries    system tools

kernel

Computer hardware

# OS objectives

- convenience
  - makes the computer more convenient to use (for programmers and users)

- efficiency
  - allows computer system resources to be used in an efficient manner

- ability to evolve
  - permit effective development, testing, and introduction of new software (debug, isolation)

# services provided by the os

- services for users
  - **program execution, usually many at the same time**
  - error detection and response
  - support for program development
  - security (user login, user confinement, etc.)
  - accounting
- services for programs (or programmers)
  - resource management
    - es. memory and cpu time
  - access to I/O devices
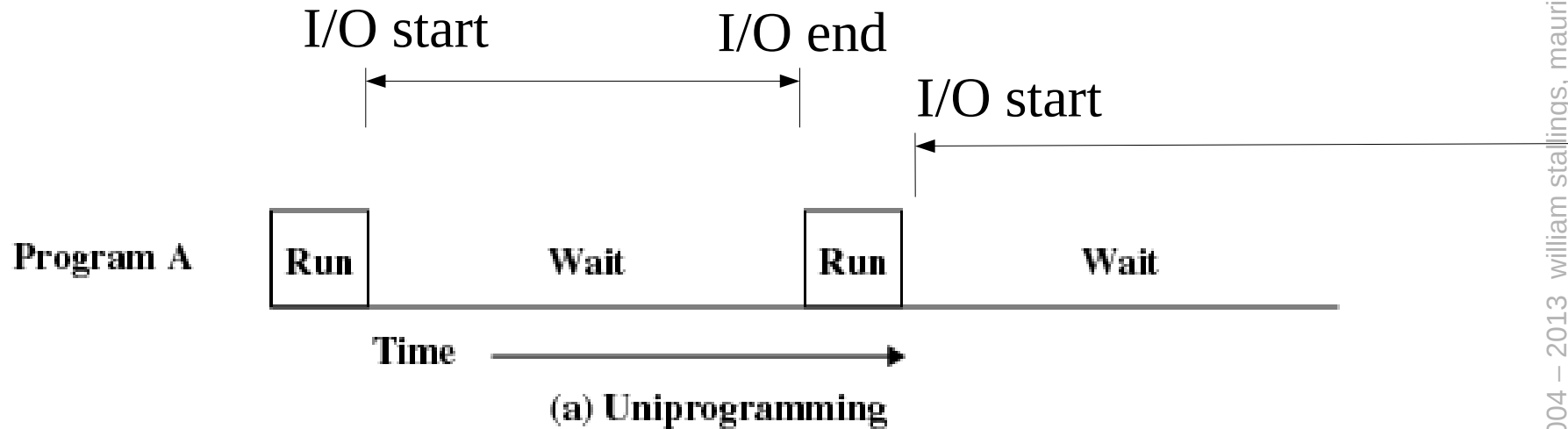    - es. files

# efficiency: i/o vs. cpu

Read one record from file 15 µs
Execute 100 instructions 1 µs
Write one record to file 15 µs
TOTAL 31 µs

Percent CPU Utilization $= \dfrac{1}{31} = 0.032 = 3.2\%$
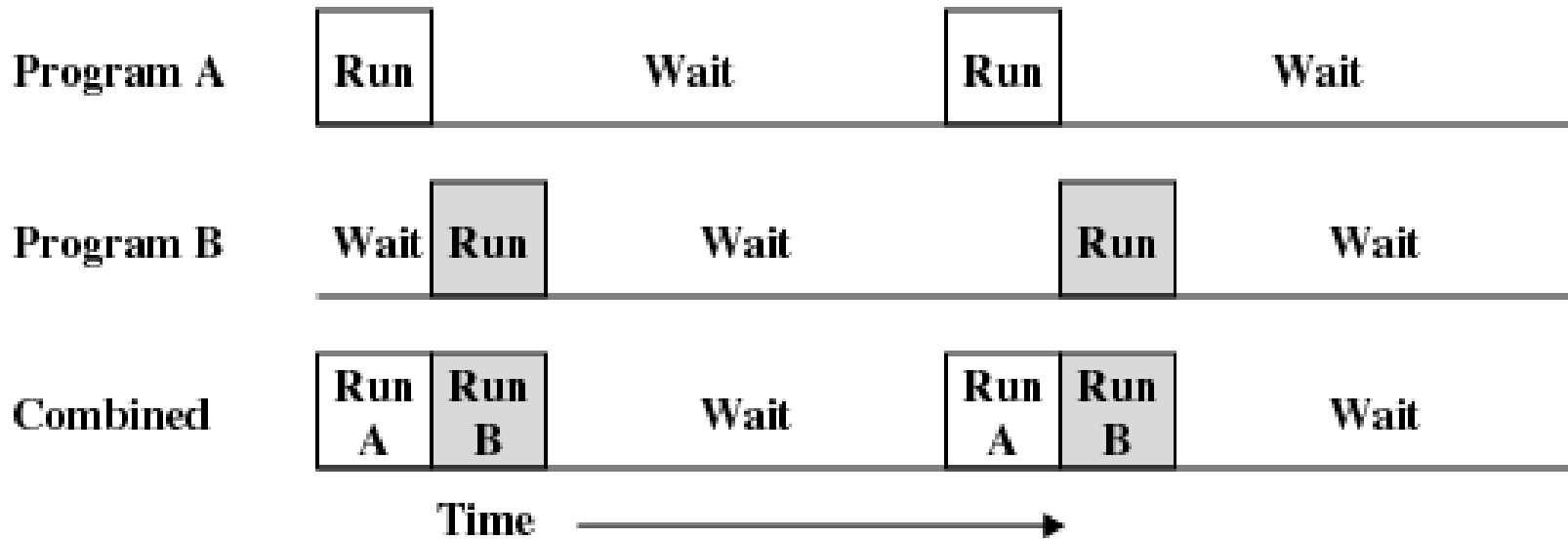
**Figure 2.4 System Utilization Example**

# uniprogramming

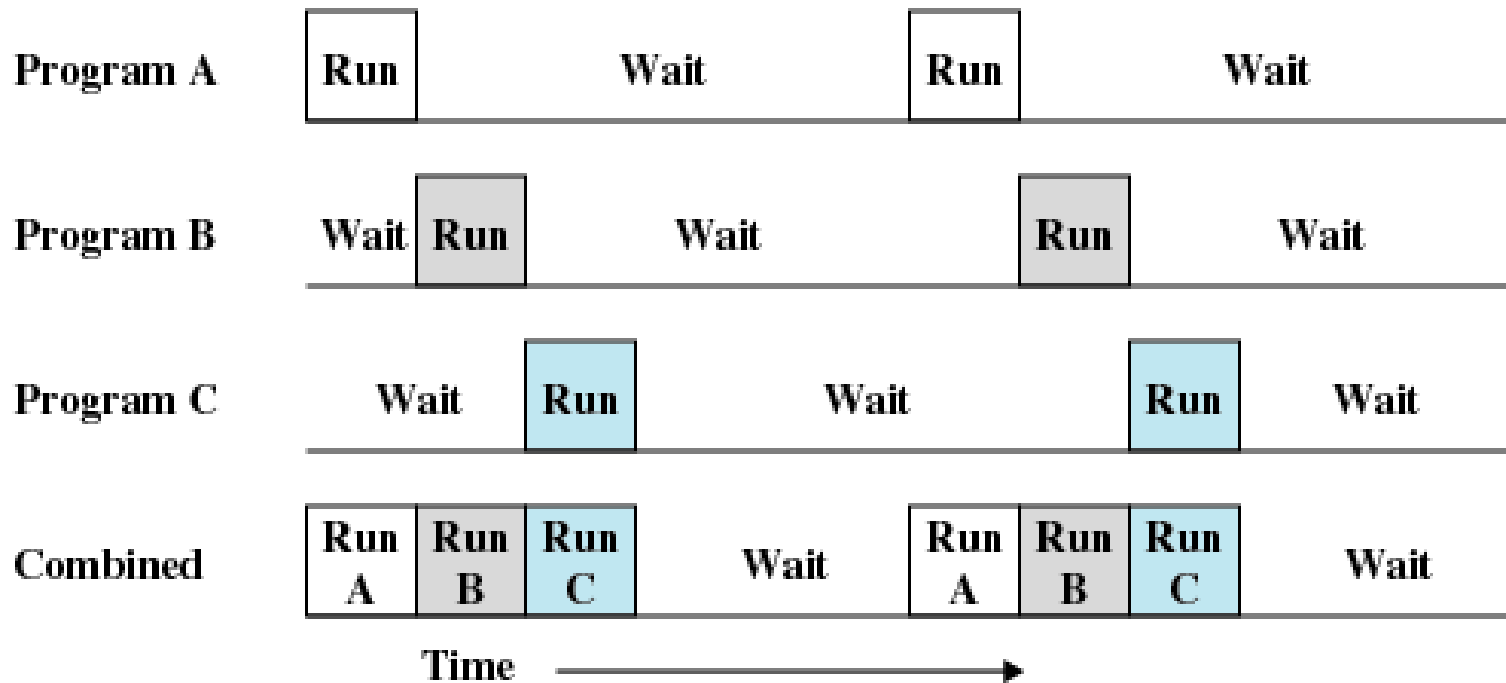- simple approach: execution must wait for I/O to complete before preceding

I/O start        I/O end

I/O start

Program A    | Run |    Wait    | Run |    Wait

Time ──────────►

(a) Uniprogramming

# multiprogramming

- when one job needs to wait for I/O, the processor can switch to the other job

| Program A | Run | Wait | Run | Wait |

| Program B | Wait | Run | Wait | Run | Wait |

| Combined | Run A | Run B | Wait | Run A | Run B | Wait |

Time →

(b) Multiprogramming with two programs

# multiprogramming



(c) Multiprogramming with three programs

17

# the user's point of view

- from the point of view of the user is a way to keep many applications active at the same time
  - e.g. I got frustrated by "uniprogrammed" smartphones need to close an app before open another one loosing the work
- PC started to be multiprogrammed in 90s (about)
  - windows 3.1, Mac classic
  - GUIs greatly rise the demand of the user
- kernels that support multiprogramming are much more complex
  - but, it is a "must have"
  - … and it is one of the main topics of this course

# I/O bound vs. CPU bound processes

- I/O bound
  - mostly waiting for some data to arrive
  - use cpu for requesting more data from devices
  - e.g. DBMS, interactive applications
- cpu bound
  - mostly perform computation
  - rarely performs i/o to get data to compute
  - e.g. multimedia coding/encoding, complex graphic rendering, etc.

# I/O bound vs. CPU bound processes

- a more complex example
  - sharable resources: cpu time and memory
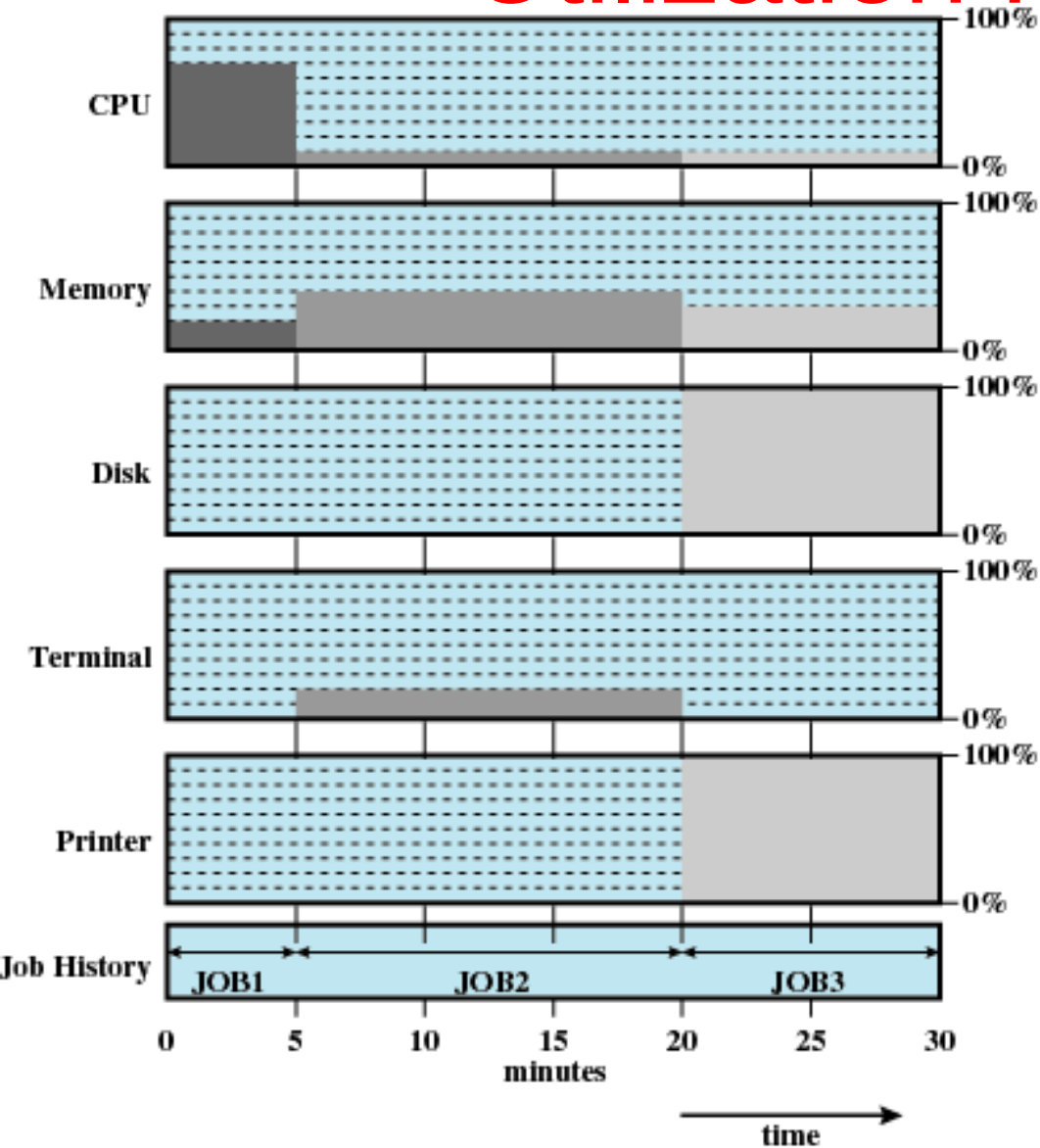  - non-sharable resources: disk, terminal, printer

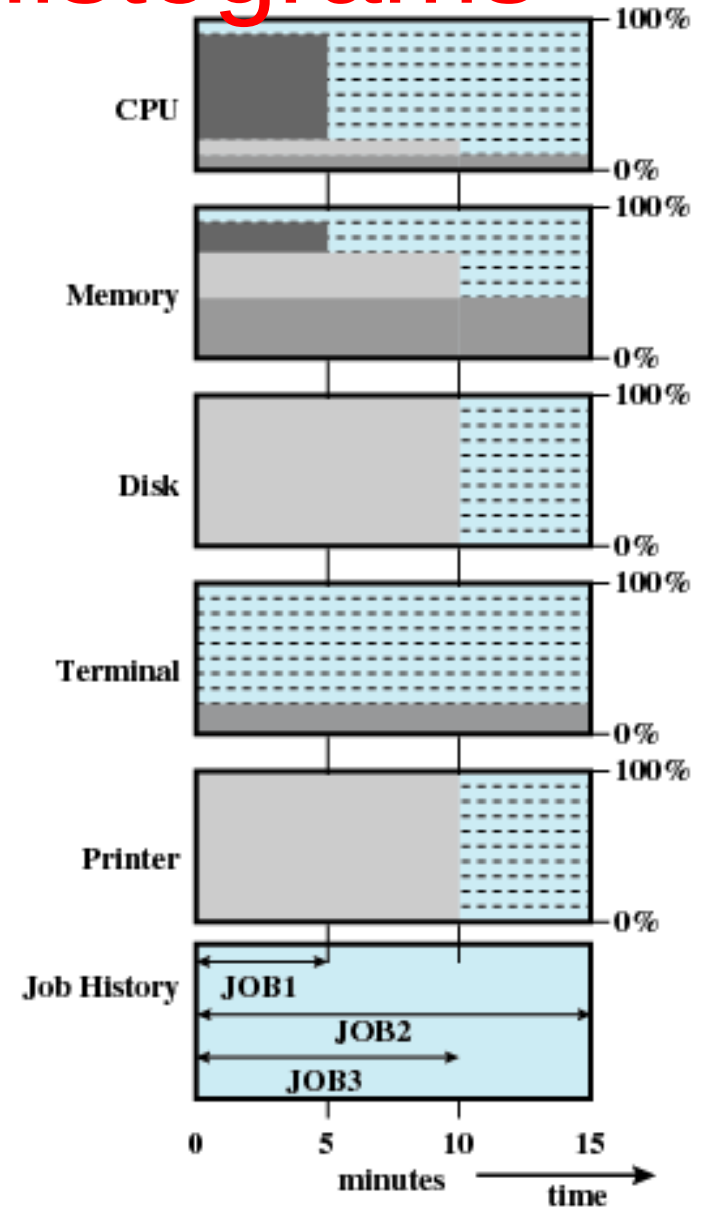| | JOB1 | JOB2 | JOB3 |
|---|---|---|---|
| **Type of job** | Heavy compute | Heavy I/O | Heavy I/O |
| **Duration** | 5 min | 15 min | 10 min |
| **Memory required** | 50 M | 100 M | 75 M |
| **Need disk?** | No | No | Yes |
| **Need terminal?** | No | Yes | No |
| **Need printer?** | No | No | Yes |
| | CPU bound | I/O bound | I/O bound |

# Utilization Histograms



(a) Uniprogramming

(b) Multiprogramming

# Major Achievements of Modern OSes

- Processes
- Memory Management
- Information protection and security
- Scheduling and resource management
- System structure

# processes

- a program in execution (running) on a computer
- a unit of activity characterized by
  - an associated set of system resources
    - memory regions
    - open files
    - etc.
  - at least one execution **thread** with its current state of CPU

# threads

- the entity that can be assigned to, and executed on, a processor
  - it is meaningful only within a process
  - described by
    - the value of the program counter
    - the value of the CPU regisers
- in modern operating systems a process may contains one or more thread
- we always assume it contains one thread

# security: user/system mode

- processes execute in **user mode**
  - certain *privileged* machine instructions may not be executed
  - only a restricted part of main memory can be accessed (**user space**)
- kernel executes in **system mode**
  - a.k.a. **kernel mode** or **supervisor mode**
  - privileged instructions can be executed
  - protected areas of memory may be accessed (**kernel space**)

# many processes, one CPU

- only one process can be executed by one CPU
- the other processes are somehow "forzen"
- what is needed to resume it is saved somewere
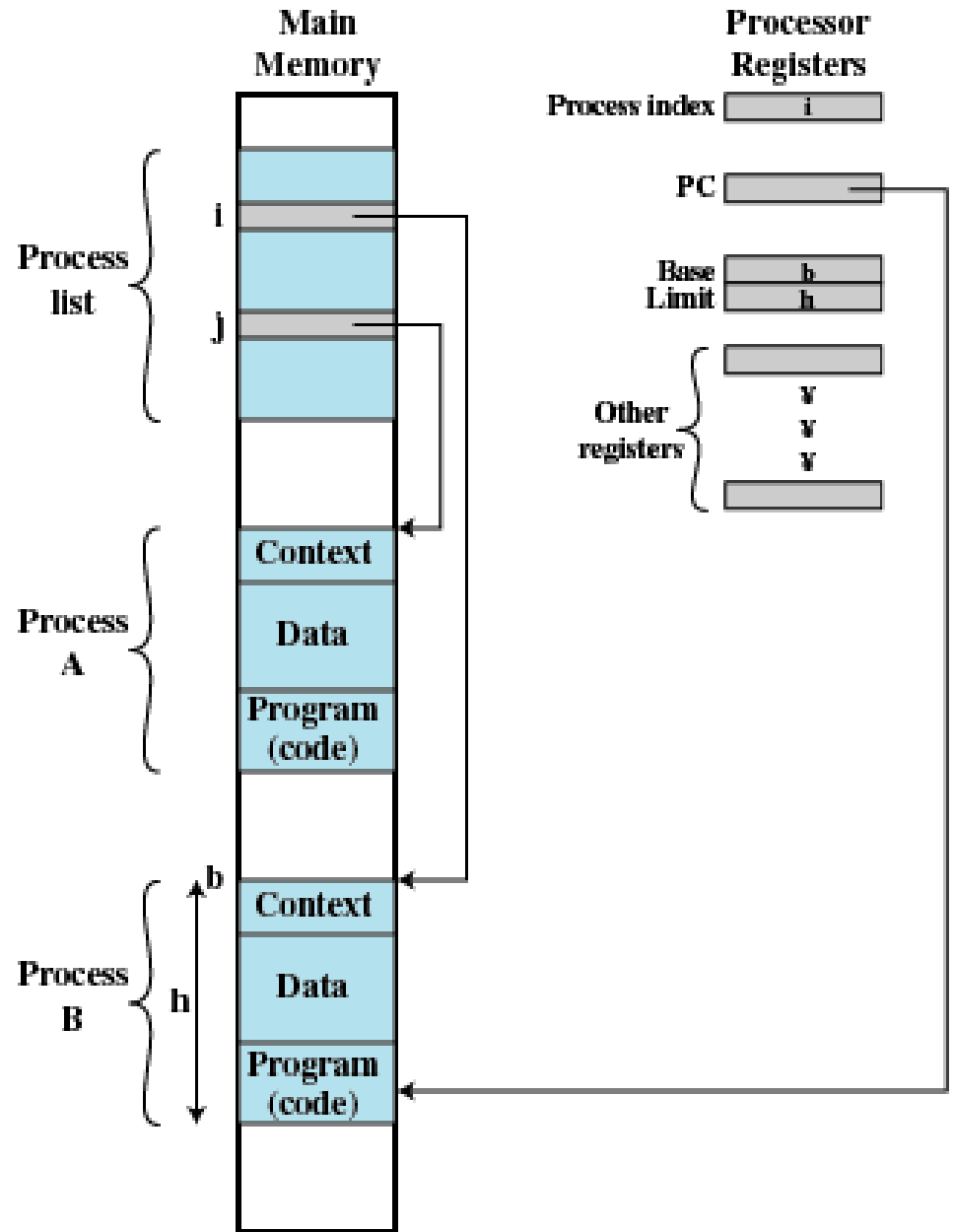  - execution context

# time sharing

- CPU time is shared among multiple users or processes
- illusion of more CPUs

# execution context

- cpu registers
- priority of the process
- is the process waiting for I/O? on which device?
- etc.
- etc.
- etc.
- etc.
- ...

# process representation

# Memory Management

- Process isolation
- Automatic allocation and management
- Support of modular programming
- Protection and access control
- Long-term storage

# Virtual Memory

- Allows programmers to address memory from a logical point of view

- Virtual memory can be much larger than Real Memory
  - processes see a large virtual address space

- Real Memory is used only for (part of the) processes that are really need it
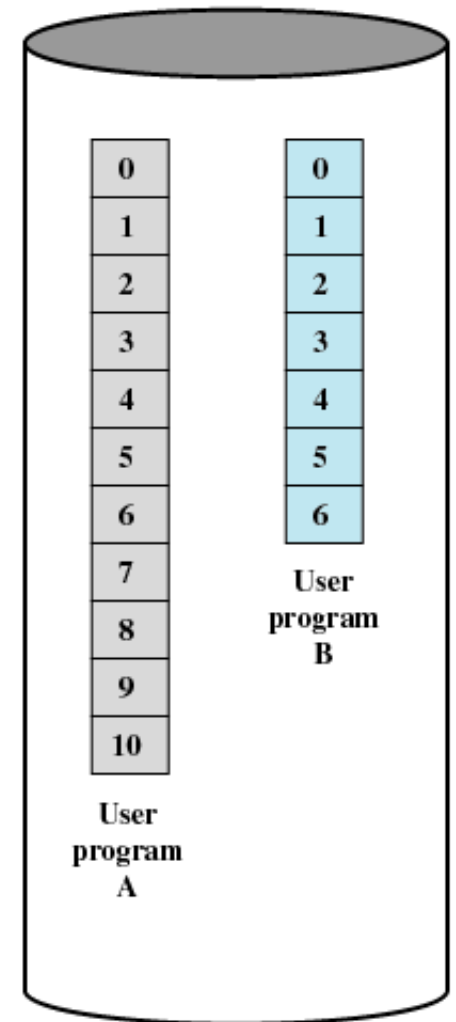
# Paging

- Allows process to be comprised of a number of fixed-size blocks, called pages
- Virtual address is a page number and an offset within the page
- Each page may be located any where in main memory
- Real address or physical address in main memory are managed only by the kernel

# Virtual Memory and Main Memory

| | | | |
|---|---|---|---|
| A.1 | | | |
| | A.0 | A.2 | |
| | A.5 | | |
| | | | |
| B.0 | B.1 | B.2 | B.3 |
| | | | |
| | | | |
| | | | |
| | | A.7 | |
| | A.9 | | |
| | | | |
| | | A.8 | |
| | | | |
| | | | |
| | | | |
| | B.5 | B.6 | |
| | | | |

**Main Memory**

Main memory consists of a number of fixed-length frames, each equal to the size of a page. For a program to execute, some or all of its pages must be in main memory.
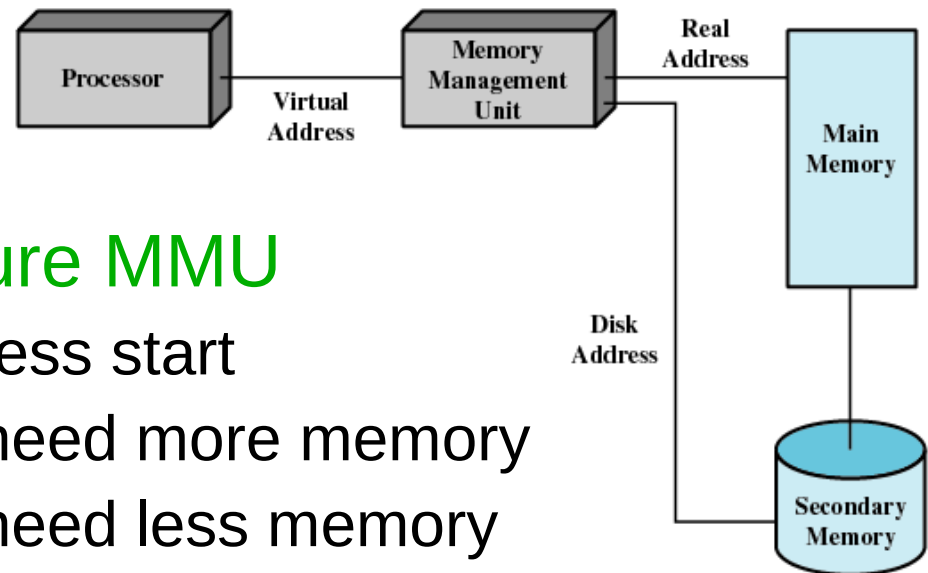
| 0 | 0 |
|---|---|
| 1 | 1 |
| 2 | 2 |
| 3 | 3 |
| 4 | 4 |
| 5 | 5 |
| 6 | 6 |
| 7 | User program B |
| 8 | |
| 9 | |
| 10 | |

User program A

**Disk**

Secondary memory (disk) can hold many fixed-length pages. A user program consists of some number of pages. Pages for all programs plus the operating system are on disk, as are files.

# Virtual Memory Addressing

- virtual addresses management
  - the Memory Management Unit of the CPU translates addresses from virtual to real
    - each time an machine instruction refers to memory
    - this is very very frequent!



  - the kernel configure MMU
    - when a new process start
    - when a process need more memory
    - when a process need less memory

# Scheduling and Resource Management

- Fairness
  - Give equal and fair access to resources
- Differential responsiveness
  - Discriminate among different classes of jobs
- Efficiency
  - Maximize throughput, minimize response time, and accommodate as many uses as possible
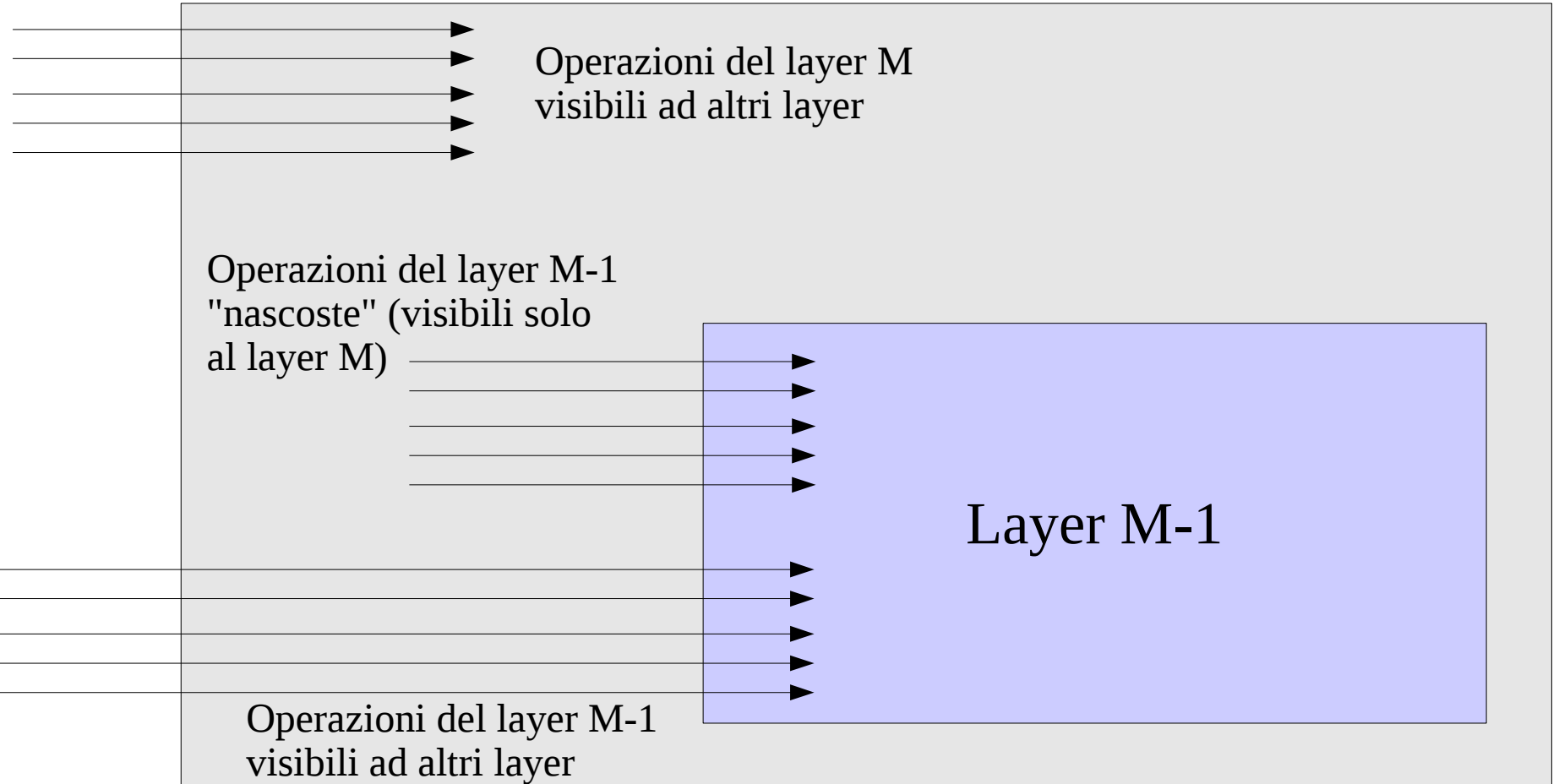
# Scheduling Elements

- queues
  - at least one for each resource
- CPU
  - short term
    - contains processes in main memory and ready to run
    - short term scheduler / dispatcher
      - simple approach: round robin (circular queue)
  - long term
    - new jobs waiting for the processor
    - long term scheduler
- I/O
  - at least queue for each device
  - interrupts

# System Structure

- View the system as a series of levels
- Each level performs a related subset of functions
- Each level relies on the next lower level to perform more primitive functions
- This decomposes a problem into a number of more manageable subproblems

# Layered Systems

Operazioni del layer M
visibili ad altri layer

Operazioni del layer M-1
"nascoste" (visibili solo
al layer M)

Layer M-1

Operazioni del layer M-1
visibili ad altri layer

© 2002, 2003 Renzo Davoli e Alberto Montresor
GNU FDL

# Hypothetical System Structure
## *Hardware Levels*

- Level 1
  - Electronic circuits
  - Objects are registers, memory cells, and logic gates
  - Operations are clearing a register or reading a memory location
- Level 2
  - Processor's instruction set
  - Operations such as add, subtract, load, and store

# Hypothetical System Structure
## *Hardware Levels*

- Level 3
  - Adds the concept of a procedure or subroutine, plus call/return operations
    - at this level everything of what you know in language like C or Java can be done, but I/O
- Level 4
  - **Interrupts    <<<< very important**

# Hypothetical System Structure
## *Basic Multiprogramming and Memory Management*

- Level 5
  - Process management from the point of view of CPU (no I/O support)
  - Suspend and resume processes
    - execution context
- Level 6
  - Secondary storage devices
  - Transfer of blocks of data
- Level 7
  - Creates logical address space for processes
  - Organizes virtual address space into blocks

# Hypothetical System Structure
## *Process comunication, I/O, and Inter Process Comunication*

- Level 8
  - Communication of information and messages between processes
    - allow client/server within a single machine
- Level 9
  - Supports long-term storage of named files
- Level 10
  - Provides access to external devices using standardized interfaces

# Hypothetical System Structure
## *Process comunication, I/O, and Inter Process Comunication*

- Level 11
  - Responsible for maintaining the association between the external and internal identifiers

- Level 12
  - Provides full-featured facility for the support of processes

# Hypothetical System Structure
## *User Interface*

- Level 13
  - Provides an interface to the operating system for the user (shell)
  - windows and pointer GUI

# Modern Operating Systems

- Microkernel architecture
  - Assigns only a few essential functions to the kernel
    - Address spaces
    - Interprocess communication (IPC)
    - Basic scheduling
  - everything else is implemented as a process

# Modern Operating Systems

- Multithreading
  - Process is divided into threads that can run concurrently
    - Thread
      - Dispatchable unit of work
      - executes sequentially and is interruptable
    - Process is a collection of one or more threads

# Modern Operating Systems

- Symmetric multiprocessing (SMP)
  - There are multiple processors
  - These processors share same main memory and I/O facilities
  - All processors can perform the same functions