# Components of a Linux System
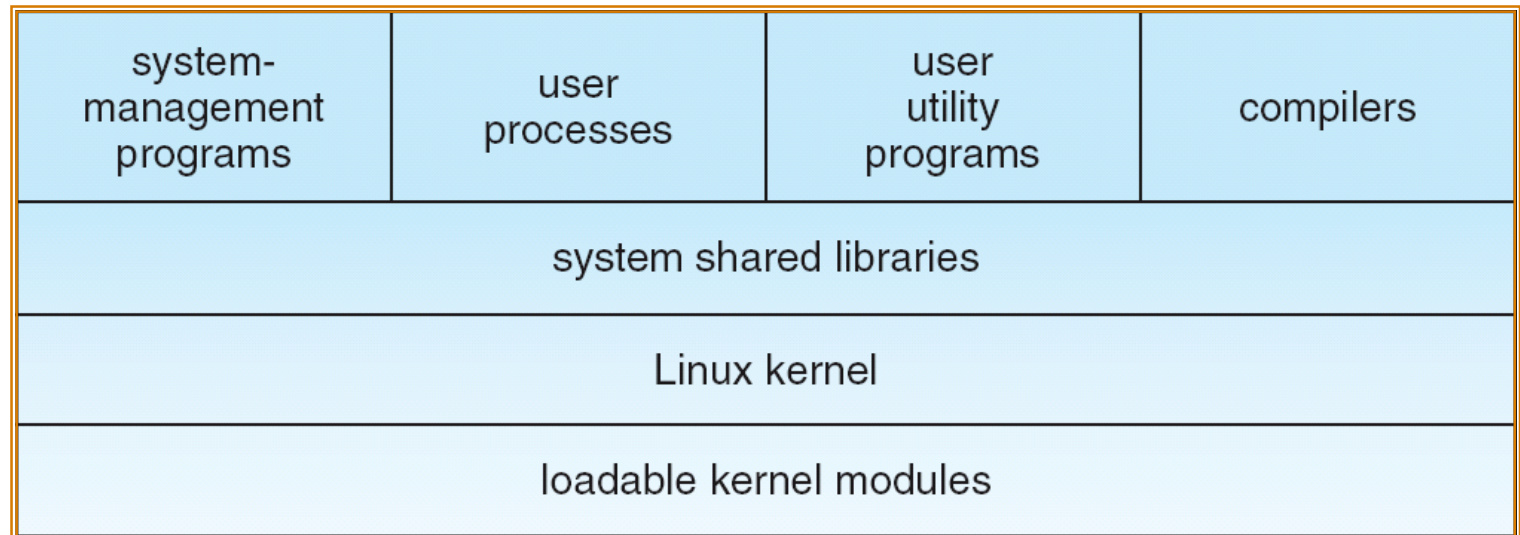
▶ Linux uses many tools developed as part of Berkeley's BSD operating system, MIT's X Window System, GNOME/KDE Window manager and the Free Software Foundation's GNU project

| system-management programs | user processes | user utility programs | compilers |
|---|---|---|---|
| system shared libraries | | | |
| Linux kernel | | | |
| loadable kernel modules | | | |

# Processes and Threads

▶ Linux uses the same internal representation for processes and threads; a thread is simply a new process that happens to share the same address space as its parent

▶ A distinction is only made when a new thread is created by the **clone** system call

- **fork** creates a new process with its own entirely new process context

- **clone** creates a new process with its own identity, but that is allowed to share the data structures of its parent

▶ Using **clone** gives an application fine-grained control over exactly what is shared between two threads

# Process Scheduling

▸ Linux uses two process-scheduling algorithms:

- A time-sharing algorithm for fair preemptive scheduling between multiple processes
- A real-time algorithm for tasks where absolute priorities are more important than fairness

▸ A process's scheduling class defines which algorithm to apply

▸ For time-sharing processes, Linux uses a prioritized, credit based algorithm

- The crediting rule

$$credits := \frac{credits}{2} + priority$$

factors in both the process's history and its priority

- This crediting system automatically prioritizes interactive or I/O-bound processes

# Process Scheduling (Cont.)

▸ Linux implements the FIFO and round-robin real-time scheduling classes; in both cases, each process has a priority in addition to its scheduling class

  ■ The scheduler runs the process with the highest priority; for equal-priority processes, it runs the process waiting the longest

  ■ FIFO processes continue to run until they either exit or block

  ■ A round-robin process will be preempted after a while and moved to the end of the scheduling queue, so that round-robing processes of equal priority automatically time-share between themselves

# Symmetric Multiprocessing

▶ Linux 2.0 was the first Linux kernel to support SMP hardware; separate processes or threads can execute in parallel on separate processors

▶ To preserve the kernel's nonpreemptible synchronization requirements, SMP imposes the restriction, via a single kernel spinlock, that only one processor at a time may execute kernel-mode code

# Memory Management

▸ Linux's physical memory-management system deals with allocating and freeing pages, groups of pages, and small blocks of memory

▸ It has additional mechanisms for handling virtual memory, memory mapped into the address space of running processes

▸ Splits memory into 3 different **zones** due to hardware characteristics

# Virtual Memory

▸ The VM system maintains the address space visible to each process:  It creates pages of virtual memory on demand, and manages the loading of those pages from disk or their swapping back out to disk as required

▸ The VM manager maintains two separate views of a process's address space:

  ■ A logical view describing instructions concerning the layout of the address space

    ● The address space consists of a set of nonoverlapping regions, each representing a continuous, page-aligned subset of the address space

  ■ A physical view of each address space which is stored in the hardware page tables for the process

# Virtual Memory (Cont.)

▶ Virtual memory regions are characterized by:

- The backing store, which describes from where the pages for a region come; regions are usually backed by a file or by nothing (*demand-zero* memory)

- The region's reaction to writes (page sharing or copy-on-write)

▶ The kernel creates a new virtual address space

1. When a process runs a new program with the **exec** system call

2. Upon creation of a new process by the **fork** system call

# Virtual Memory (Cont.)

▸ On executing a new program, the process is given a new, completely empty virtual-address space; the program-loading routines populate the address space with virtual-memory regions

▸ Creating a new process with **fork** involves creating a complete copy of the existing process's virtual address space

  ■ The kernel copies the parent process's VMA descriptors, then creates a new set of page tables for the child

  ■ The parent's page tables are copied directly into the child's, with the reference count of each page covered being incremented

  ■ After the fork, the parent and child share the same physical pages of memory in their address spaces

# Virtual Memory (Cont.)

▸ The VM paging system relocates pages of memory from physical memory out to disk when the memory is needed for something else

▸ The VM paging system can be divided into two sections:

  ▪ The pageout-policy algorithm decides which pages to write out to disk, and when

  ▪ The paging mechanism actually carries out the transfer, and pages data back into physical memory as needed

# Virtual Memory (Cont.)

▸ The Linux kernel reserves a constant, architecture-dependent region of the virtual address space of every process for its own internal use

▸ This kernel virtual-memory area contains two regions:

  ■ A static area that contains page table references to every available physical page of memory in the system, so that there is a simple translation from physical to virtual addresses when running kernel code

  ■ The reminder of the reserved section is not reserved for any specific purpose; its page-table entries can be modified to point to any other areas of memory
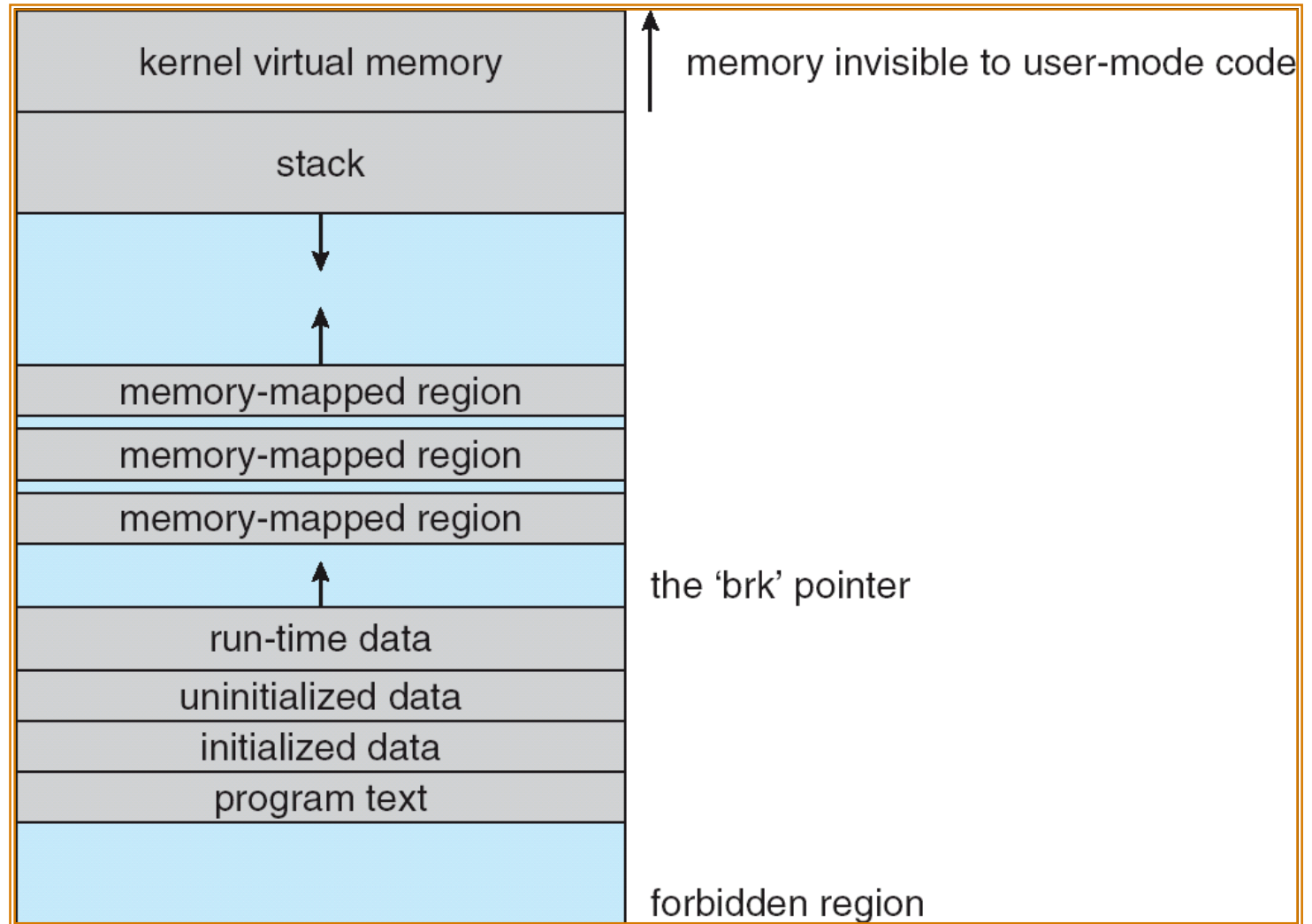
# Executing and Loading User Programs

▶ Linux maintains a table of functions for loading programs; it gives each function the opportunity to try loading the given file when an exec system call is made

▶ The registration of multiple loader routines allows Linux to support both the ELF and **a.out** binary formats

▶ Initially, binary-file pages are mapped into virtual memory

- Only when a program tries to access a given page will a page fault result in that page being loaded into physical memory

▶ An ELF-format binary file consists of a header followed by several page-aligned sections

- The ELF loader works by reading the header and mapping the sections of the file into separate regions of virtual memory

# Memory Layout for **ELF** Programs

| | |
|---|---|
| kernel virtual memory | memory invisible to user-mode code |
| stack | |
| ↓ ↑ | |
| memory-mapped region | |
| memory-mapped region | |
| memory-mapped region | |
| ↑ | the 'brk' pointer |
| run-time data | |
| uninitialized data | |
| initialized data | |
| program text | |
| | forbidden region |

# Static and Dynamic Linking

▸ A program whose necessary library functions are embedded directly in the program's executable binary file is *statically* linked to its libraries

▸ The main disadvantage of static linkage is that every program generated must contain copies of exactly the same common system library functions

▸ *Dynamic* linking is more efficient in terms of both physical memory and disk-space usage because it loads the system libraries into memory only once

# File Systems

▶ To the user, Linux's file system appears as a hierarchical directory tree obeying UNIX semantics

▶ Internally, the kernel hides implementation details and manages the multiple different file systems via an abstraction layer, that is, the *virtual file system (VFS)*

▶ The Linux VFS is designed around object-oriented principles and is composed of two components:

- A set of definitions that define what a file object is allowed to look like
  - The *inode-object* and the *file-object* structures represent individual files
  - the *file system object* represents an entire file system
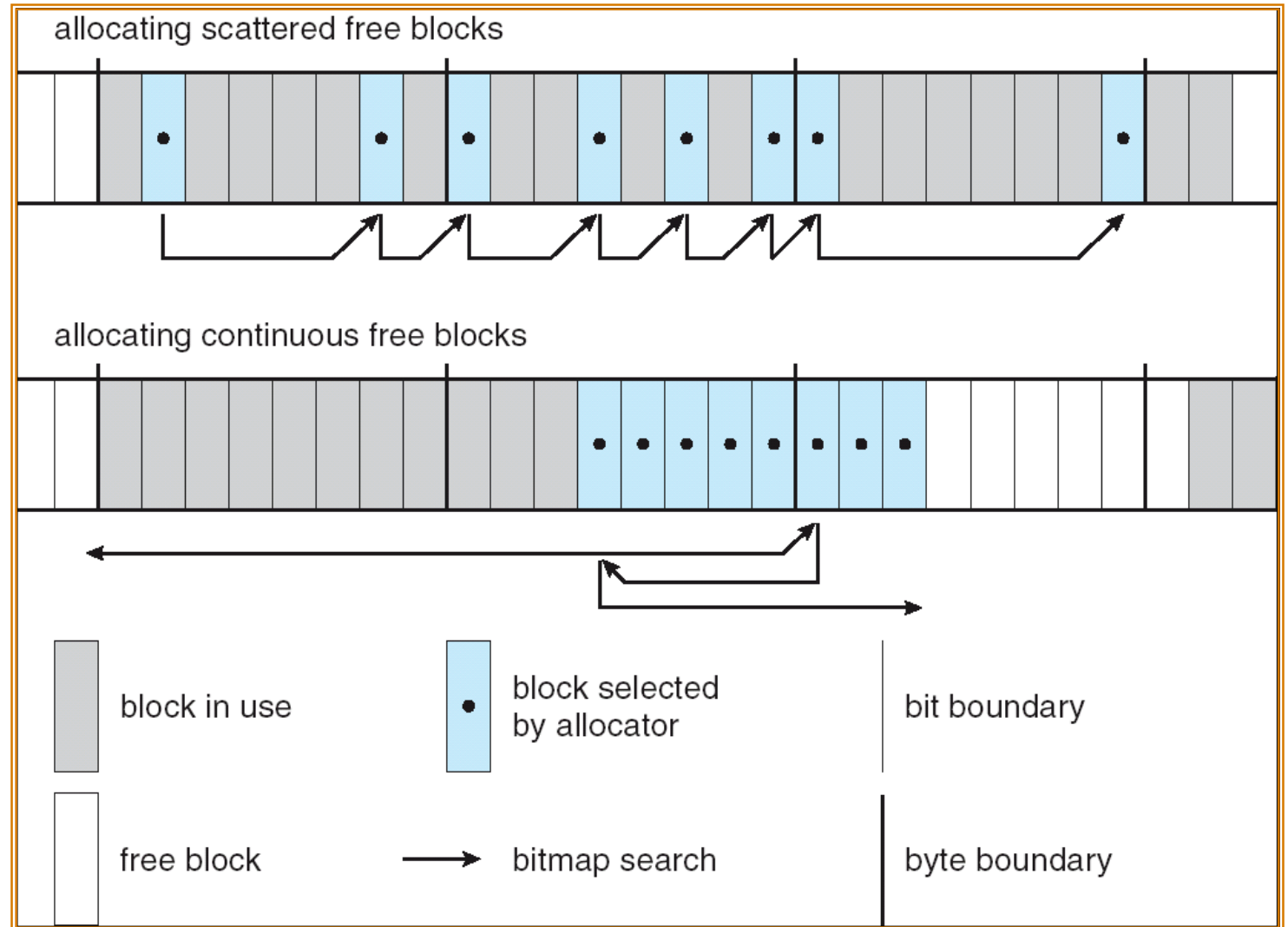- A layer of software to manipulate those objects

# The Linux Ext2fs File System

▸ Ext2fs uses a mechanism similar to that of BSD Fast File System (ffs) for locating data blocks belonging to a specific file

▸ The main differences between ext2fs and ffs concern their disk allocation policies

  ▪ In ffs, the disk is allocated to files in blocks of 8Kb, with blocks being subdivided into fragments of 1Kb to store small files or partially filled blocks at the end of a file

  ▪ Ext2fs does not use fragments; it performs its allocations in smaller units

    ● The default block size on ext2fs is 1Kb, although 2Kb and 4Kb blocks are also supported

  ▪ Ext2fs uses allocation policies designed to place logically adjacent blocks of a file into physically adjacent blocks on disk, so that it can submit an I/O request for several disk blocks as a single operation

# Ext2fs Block-Allocation Policies

# The Linux Proc File System

▶ The **proc** file system does not store data, rather, its contents are computed on demand according to user file I/O requests

▶ **proc** must implement a directory structure, and the file contents within; it must then define a unique and persistent inode number for each directory and files it contains

  ■ It uses this inode number to identify just what operation is required when a user tries to read from a particular file inode or perform a lookup in a particular directory inode

  ■ When data is read from one of these files, **proc** collects the appropriate information, formats it into text form and places it into the requesting process's read buffer

# Input and Output

▸ The Linux device-oriented file system accesses disk storage through two caches:

  ■ Data is cached in the page cache, which is unified with the virtual memory system

  ■ Metadata is cached in the buffer cache, a separate cache indexed by the physical disk block

▸ Linux splits all devices into three classes:

  ■ *block devices* allow random access to completely independent, fixed size blocks of data

  ■ *character devices* include most other devices; they don't need to support the functionality of regular files

  ■ *network devices* are interfaced via the kernel's networking subsystem

# Block Devices

▸ Provide the main interface to all disk devices in a system

▸ The *block buffer* cache serves two main purposes:
  ■ it acts as a pool of buffers for active I/O
  ■ it serves as a cache for completed I/O

▸ The *request manager* manages the reading and writing of buffer contents to and from a block device driver

# Character Devices

▶ A device driver which does not offer random access to fixed blocks of data

▶ A character device driver must register a set of functions which implement the driver's various file I/O operations

▶ The kernel performs almost no preprocessing of a file read or write request to a character device, but simply passes on the request to the device

▶ The main exception to this rule is the special subset of character device drivers which implement terminal devices, for which the kernel maintains a standard interface

# Passing Data Between Processes

▶ The pipe mechanism allows a child process to inherit a communication channel to its parent, data written to one end of the pipe can be read a the other

▶ Shared memory offers an extremely fast way of communicating; any data written by one process to a shared memory region can be read immediately by any other process that has mapped that region into its address space

▶ To obtain synchronization, however, shared memory must be used in conjunction with another Interprocess-communication mechanism

# Shared Memory Object

▶ The shared-memory object acts as a backing store for shared-memory regions in the same way as a file can act as backing store for a memory-mapped memory region

▶ Shared-memory mappings direct page faults to map in pages from a persistent shared-memory object

▶ Shared-memory objects remember their contents even if no processes are currently mapping them into virtual memory

# Security

▸ The *pluggable authentication modules (PAM)* system is available under Linux

▸ PAM is based on a shared library that can be used by any system component that needs to authenticate users

▸ Access control under UNIX systems, including Linux, is performed through the use of unique numeric identifiers (**uid** and **gid**)

▸ Access control is performed by assigning objects a *protections mask*, which specifies which access modes—read, write, or execute—are to be granted to processes with owner, group, or world access