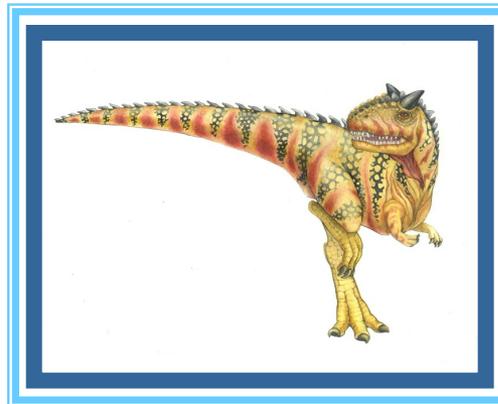
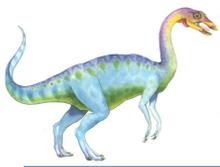


# Chapter 2: System Structures

---



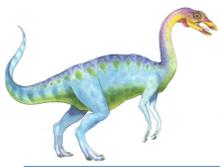


# Chapter 2: Operating-System Structures

---

- Operating System Services
- User and Operating System-Interface
- System Calls
- System Services
- Linkers and Loaders
- Why Applications are Operating System Specific
- Operating-System Design and Implementation
- Operating System Structure
- Building and Booting an Operating System
- Operating System Debugging
- **bpfttrace (dtrace), Systemtap, and Perf**



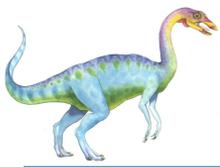


# Objectives

---

- Identify services provided by an operating system
- Illustrate how system calls are used to provide operating system services
- Compare and contrast monolithic, layered, microkernel, modular, and hybrid strategies for designing operating systems
- Illustrate the process for booting an operating system
- Apply tools for monitoring operating system performance
- Design and implement kernel modules for interacting with a Linux kernel



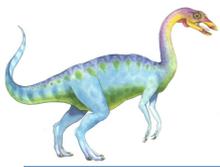


# Operating System Services

---

- Operating systems provide an environment for execution of programs and services to programs and users
- One set of operating-system services provides functions that are helpful to the user:
  - **User interface** - Almost all operating systems have a user interface (**UI**).
    - ▶ Varies between **Command-Line (CLI)**, **Graphics User Interface (GUI)**, **touch-screen**, **Batch**
  - **Program execution** - The system must be able to load a program into memory and to run that program, end execution, either normally or abnormally (indicating error)
  - **I/O operations** - A running program may require I/O, which may involve a file or an I/O device
  - **File-system manipulation** - The file system is of particular interest. Programs need to read and write files and directories, create and delete them, search them, list file information, permission management.

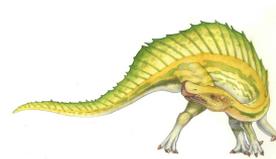


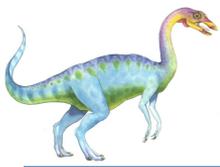


# Operating System Services (Cont.)

---

- **Communications** – Processes may exchange information, on the same computer or between computers over a network
  - ▶ Communications may be via shared memory or through message passing (packets moved by the OS)
- **Error detection** – OS needs to be constantly aware of possible errors
  - ▶ May occur in the CPU and memory hardware, in I/O devices, in user program
  - ▶ For each type of error, OS should take the appropriate action to ensure correct and consistent computing
  - ▶ Debugging facilities can greatly enhance the user's and programmer's abilities to efficiently use the system

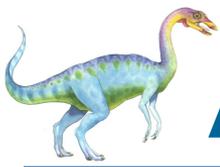




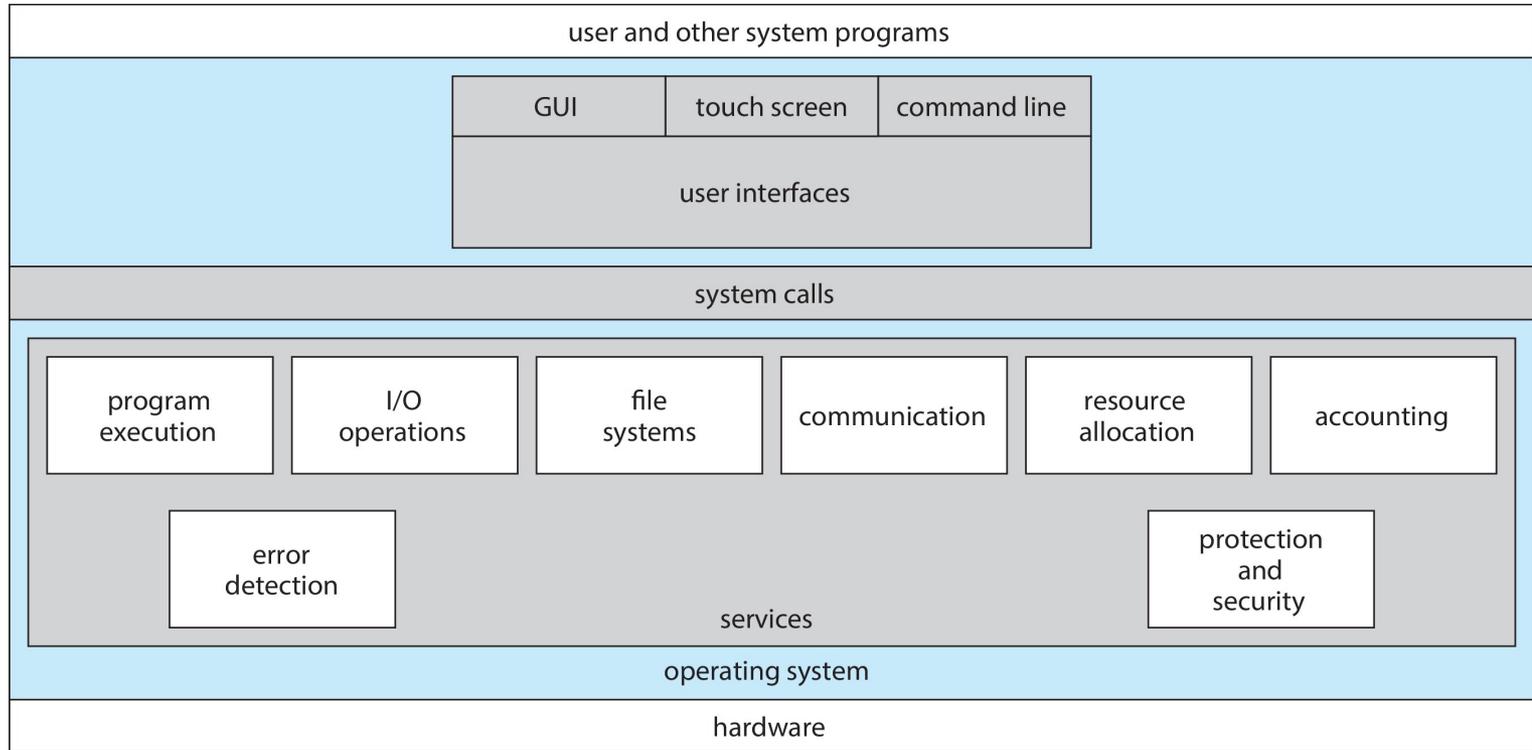
# Operating System Services (Cont.)

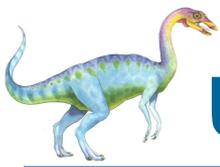
- Another set of OS functions exists for ensuring the efficient operation of the system itself via resource sharing
  - **Resource allocation** - When multiple users or multiple jobs running concurrently, resources must be allocated to each of them
    - ▶ Many types of resources - Some (such as CPU cycles, main memory, and file storage) may have special allocation code, others (such as I/O devices) may have general request and release code
  - **Logging** - To keep track of which users use how much and what kinds of computer resources
  - **Protection and security** - The owners of information stored in a multiuser or networked computer system may want to control use of that information, concurrent processes should not interfere with each other
    - ▶ **Protection** involves ensuring that all access to system resources is controlled
    - ▶ **Security** of the system from outsiders requires user authentication, extends to defending external I/O devices from invalid access attempts
    - ▶ If a system is to be protected and secure, precautions must be instituted throughout it. A chain is only as strong as its weakest link.





# A View of Operating System Services



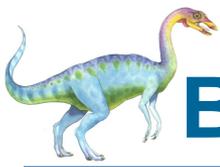


# User Operating System Interface - CLI

---

- CLI or **command interpreter** allows direct command entry
  - ▶ Sometimes implemented in kernel, sometimes by systems program
  - ▶ Sometimes multiple flavors implemented – **shells**
  - ▶ Primarily fetches a command from user and executes it
    - Sometimes commands built-in, sometimes just names of programs
      - » If the latter, adding new features doesn't require shell modification

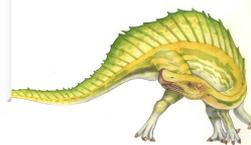


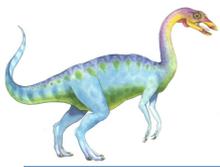


# Bourne Shell Command Interpreter

```
Default
New Info Close Execute Bookmarks

Default Default
PBG-Mac-Pro:~ pbg$ w
15:24 up 56 mins, 2 users, load averages: 1.51 1.53 1.65
USER TTY FROM LOGIN@ IDLE WHAT
pbg console - 14:34 50 -
pbg s000 - 15:05 - w
PBG-Mac-Pro:~ pbg$ iostat 5
          disk0          disk1          disk10          cpu          load average
      KB/t tps MB/s      KB/t tps MB/s      KB/t tps MB/s  us sy id 1m 5m 15m
      33.75 343 11.30      64.31 14 0.88      39.67 0 0.02  11 5 84 1.51 1.53 1.65
      5.27 320 1.65        0.00 0 0.00        0.00 0 0.00   4 2 94 1.39 1.51 1.65
      4.28 329 1.37        0.00 0 0.00        0.00 0 0.00   5 3 92 1.44 1.51 1.65
^C
PBG-Mac-Pro:~ pbg$ ls
Applications          Music                  WebEx
Applications (Parallels)  Pando Packages       config.log
Desktop               Pictures               getsmartdata.txt
Documents             Public                 imp
Downloads             Sites                  log
Dropbox               Thumbs.db              panda-dist
Library               Virtual Machines       prob.txt
Movies                Volumes                scripts
PBG-Mac-Pro:~ pbg$ pwd
/Users/pbg
PBG-Mac-Pro:~ pbg$ ping 192.168.1.1
PING 192.168.1.1 (192.168.1.1): 56 data bytes
64 bytes from 192.168.1.1: icmp_seq=0 ttl=64 time=2.257 ms
64 bytes from 192.168.1.1: icmp_seq=1 ttl=64 time=1.262 ms
^C
--- 192.168.1.1 ping statistics ---
2 packets transmitted, 2 packets received, 0.0% packet loss
round-trip min/avg/max/stddev = 1.262/1.760/2.257/0.498 ms
PBG-Mac-Pro:~ pbg$
```



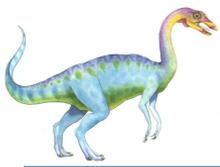


# User Operating System Interface - GUI

---

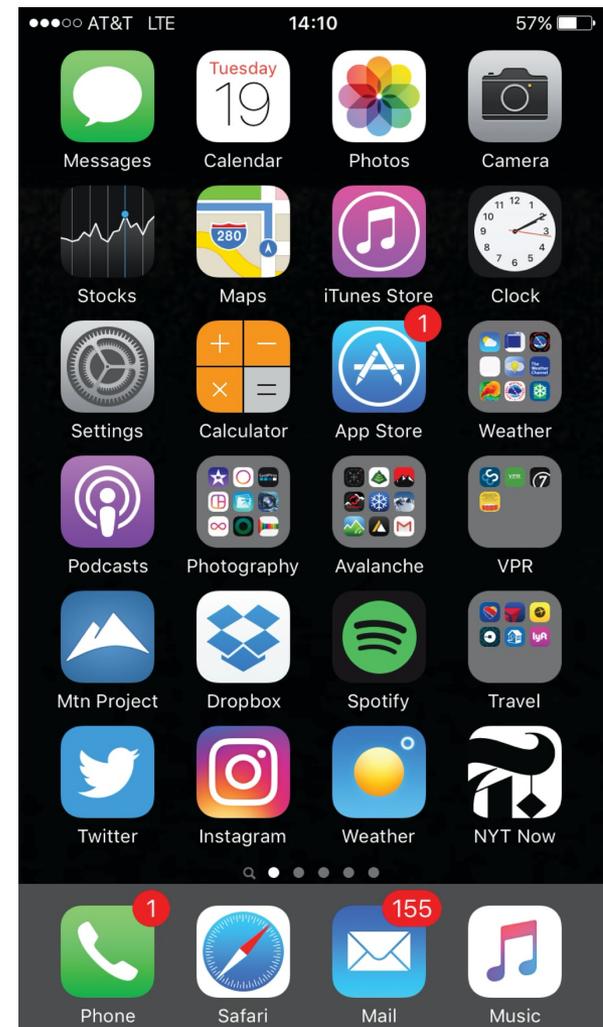
- User-friendly **desktop** metaphor interface
  - Usually mouse, keyboard, and monitor
  - **Icons** represent files, programs, actions, etc
  - Various mouse buttons over objects in the interface cause various actions (provide information, options, execute function, open directory (known as a **folder**))
  - Invented at Xerox PARC
  
- Many systems now include both CLI and GUI interfaces
  - Microsoft Windows is GUI with CLI “command” shell
  - Apple Mac OS X is “Aqua” GUI interface with UNIX kernel underneath and shells available
  - Unix and Linux have CLI with optional GUI interfaces (CDE, KDE, GNOME)





# Touchscreen Interfaces

- Touchscreen devices require new interfaces
  - Mouse not possible or not desired
  - Actions and selection based on gestures
  - Virtual keyboard for text entry
  - Voice Commands



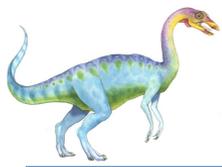


# The Mac OS X GUI

The image displays a collage of various Mac OS X GUI components:

- File Browser:** A window showing a directory listing of files and folders with columns for Name, Size, Date, and Kind.
- Dictionary:** A search window for the dictionary application, showing search results for "dictionary" and a preview of the "Shanter Oxford English Dictionary".
- Presentation Software:** A window showing a slide titled "DSSD HIGH PERFORMANCE PARALLEL FILE SYSTEMS" with a subtitle "CLICK TO EDIT MASTER SUBTITLE STYLE".
- Desktop Environment:** A desktop view showing a dock with various applications, a menu bar at the top, and a background image of a dinosaur.
- System Information:** A window showing system information for "Apple Computer Inc.", including phone, fax, and email details.
- Calculator:** A window showing a calculator interface with a display showing "290.2425".





# System Calls

---

- Programming interface to the services provided by the OS
- Typically written in a high-level language (C or C++)
- Mostly accessed by programs via a high-level **Application Program Interface (API)** rather than direct system call use
- Three most common APIs are Win32 API for Windows, POSIX API for POSIX-based systems (including virtually all versions of UNIX, Linux, and Mac OS X), and Java API for the Java virtual machine (JVM)
- Why use APIs rather than system calls?

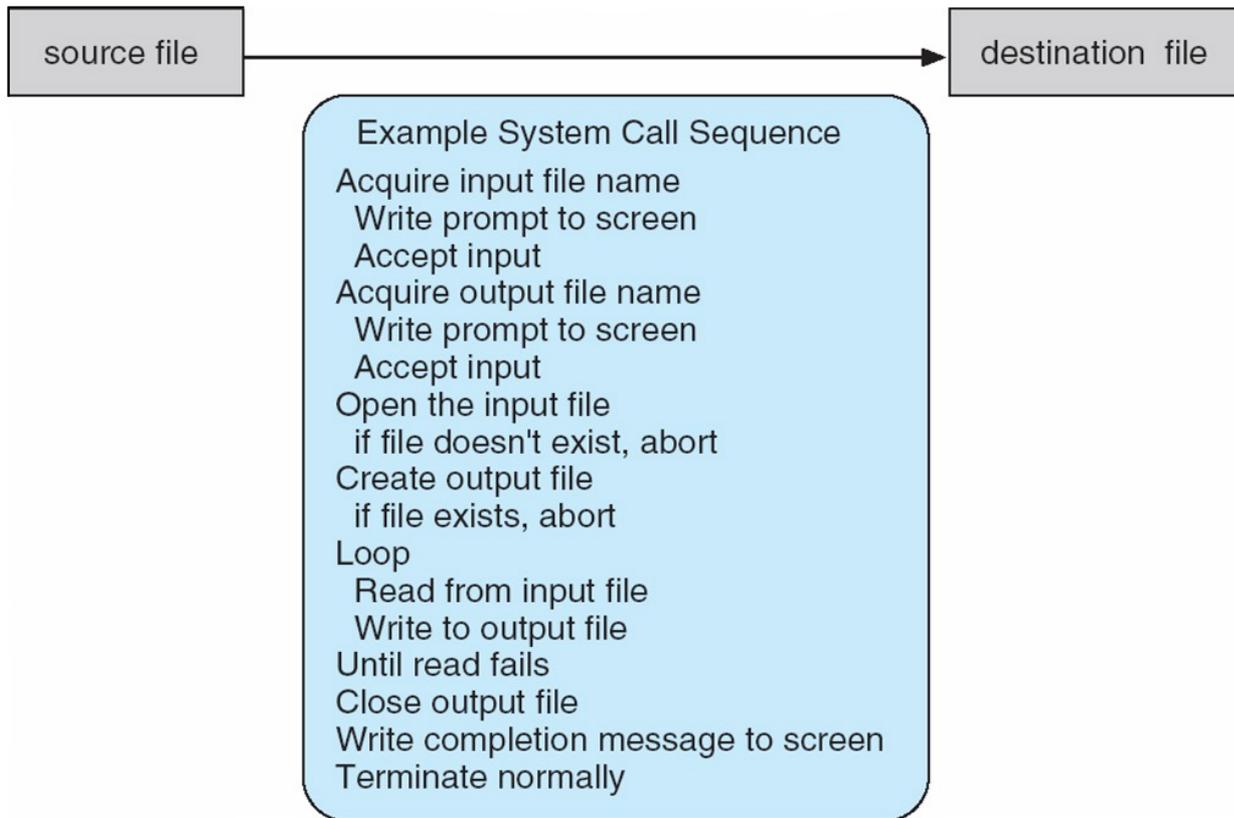
(Note that the system-call names used throughout this text are generic)

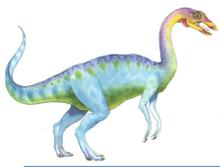




# Example of System Calls

- System call sequence to copy the contents of one file to another file





# Example of Standard API

## EXAMPLE OF STANDARD API

As an example of a standard API, consider the `read()` function that is available in UNIX and Linux systems. The API for this function is obtained from the `man` page by invoking the command

```
man read
```

on the command line. A description of this API appears below:

```
#include <unistd.h>

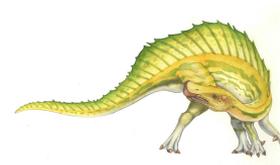
ssize_t  read(int fd, void *buf, size_t count)
```

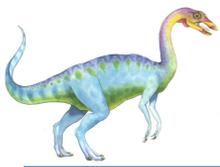
ssize_t	read	(int fd, void *buf, size_t count)
return value	function name	parameters

A program that uses the `read()` function must include the `unistd.h` header file, as this file defines the `ssize_t` and `size_t` data types (among other things). The parameters passed to `read()` are as follows:

- `int fd`—the file descriptor to be read
- `void *buf`—a buffer where the data will be read into
- `size_t count`—the maximum number of bytes to be read into the buffer

On a successful read, the number of bytes read is returned. A return value of 0 indicates end of file. If an error occurs, `read()` returns `-1`.



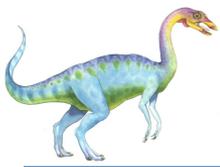


# System Call Implementation

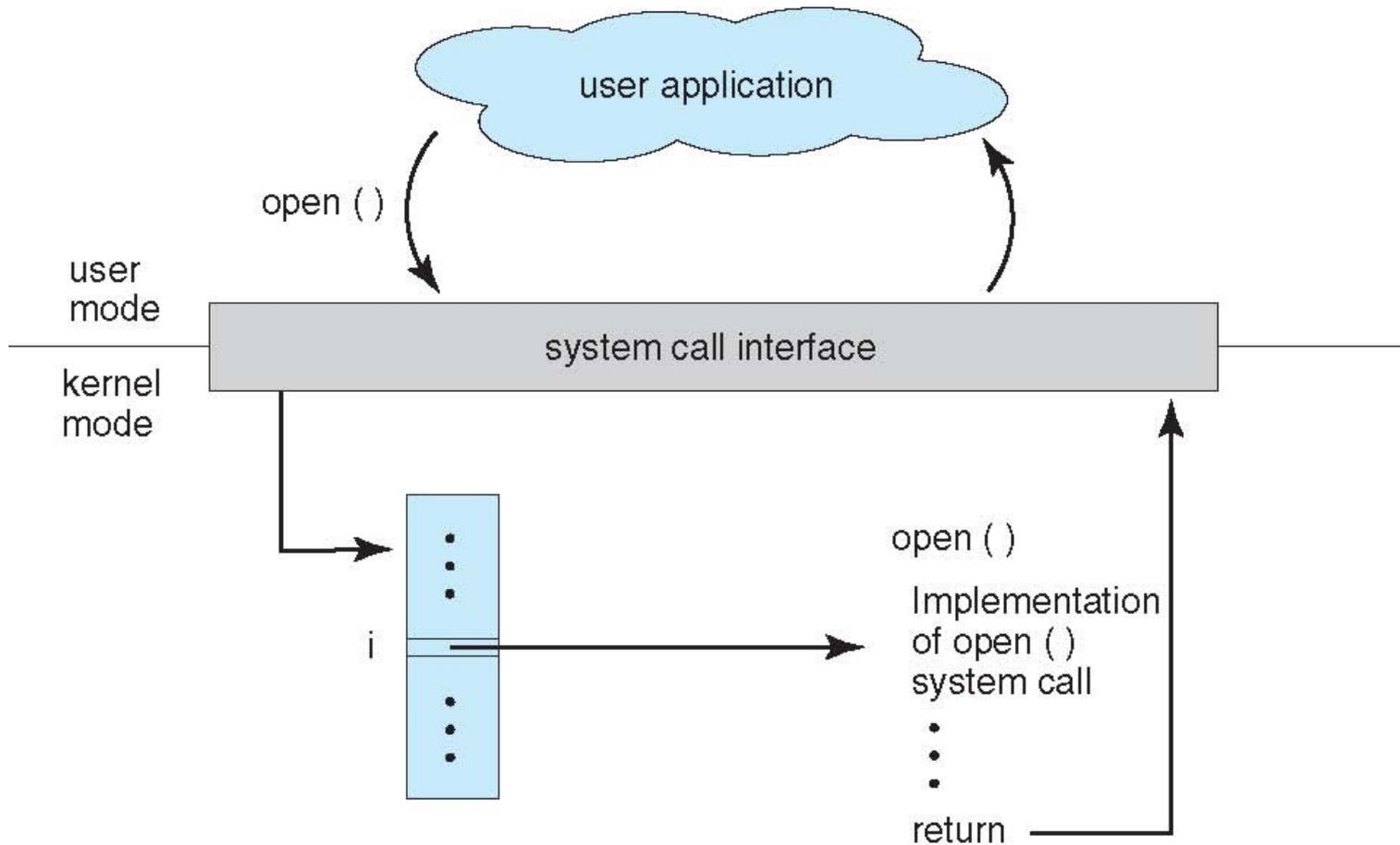
---

- Typically, a number associated with each system call
  - **System-call interface** maintains a table indexed according to these numbers
- The system call interface invokes intended system call in OS kernel and returns status of the system call and any return values
- The caller need know nothing about how the system call is implemented
  - Just needs to obey API and understand what OS will do as a result call
  - Most details of OS interface hidden from programmer by API
    - ▶ Managed by run-time support library (set of functions built into libraries included with compiler)





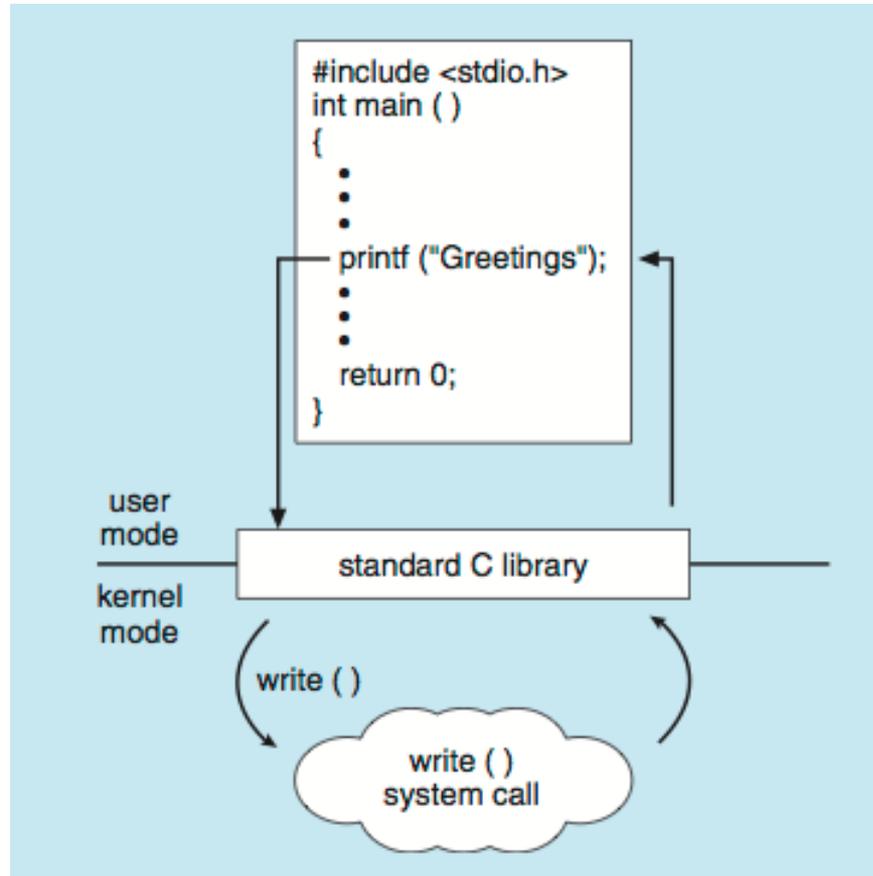
# API – System Call – OS Relationship

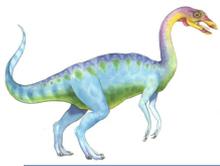




# Standard C Library Example

- C program invoking printf() library call, which calls write() system call

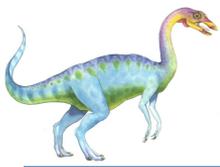




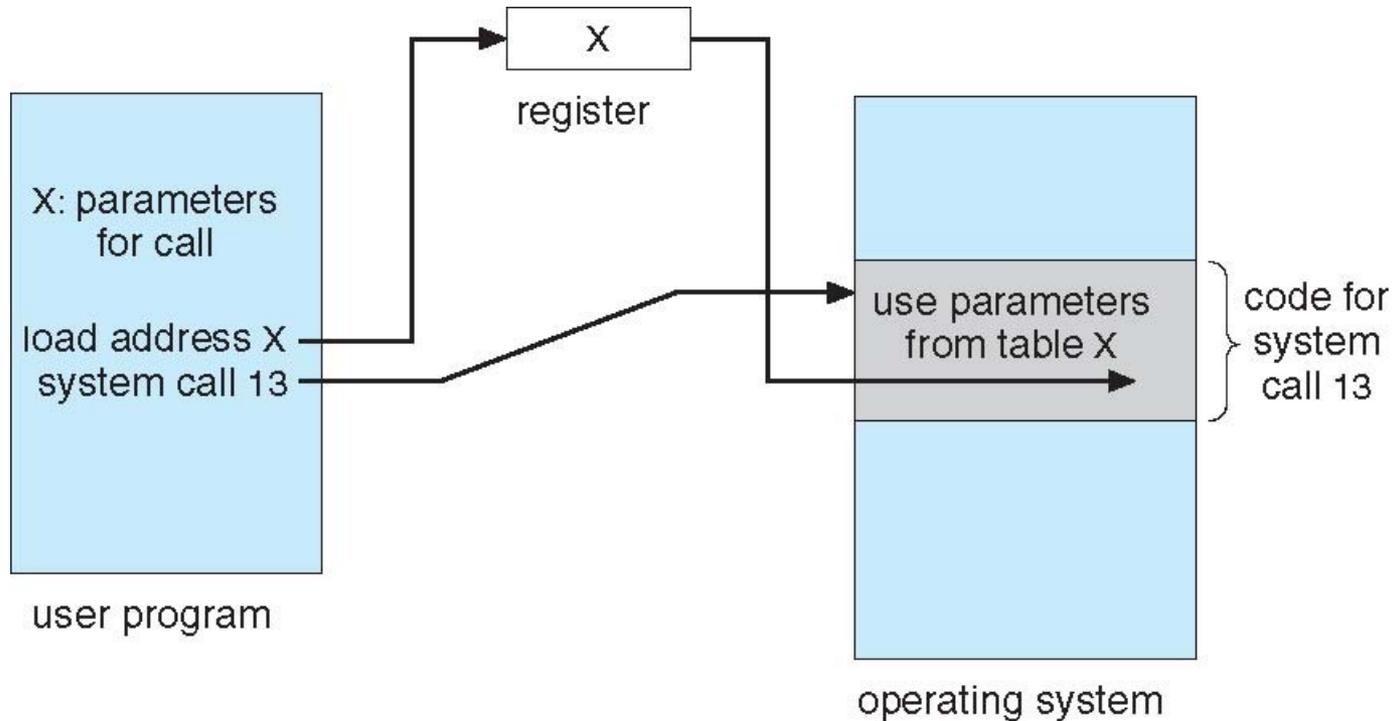
# System Call Parameter Passing

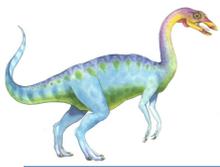
- Often, more information is required than simply identity of desired system call
  - Exact type and amount of information vary according to OS and call
  
- Three general methods used to pass parameters to the OS
  - Simplest: pass the parameters in registers
    - ▶ In some cases, may be more parameters than registers
  - Parameters stored in a block, or table, in memory, and address of block passed as a parameter in a register
    - ▶ This approach taken by Linux and Solaris
  - Parameters placed, or **pushed**, onto the **stack** by the program and **popped** off the stack by the operating system
  - Block and stack methods do not limit the number or length of parameters being passed





# Parameter Passing via Table

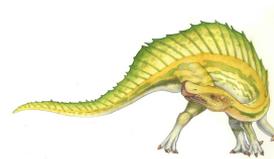


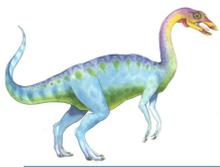


# Types of System Calls

---

- Process control
  - end, abort
  - load, execute
  - create process, terminate process
  - get process attributes, set process attributes
  - wait for time
  - wait event, signal event
  - allocate and free memory
  
  - Dump memory if error
  - **Debugger** for determining **bugs, single step** execution
  - **Locks** for managing access to shared data between processes



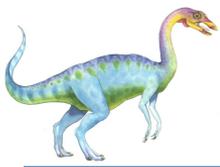


# Types of System Calls

---

- File management
  - create file, delete file
  - open, close file
  - read, write, reposition
  - get and set file attributes
  
- Device management
  - request device, release device
  - read, write, reposition
  - get device attributes, set device attributes
  - logically attach or detach devices



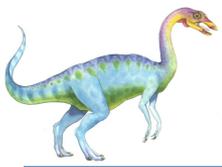


# Types of System Calls (Cont.)

---

- Information maintenance
  - get time or date, set time or date
  - get system data, set system data
  - get and set process, file, or device attributes
  
- Communications
  - create, delete communication connection
  - send, receive messages if **message passing model** to **host name** or **process name**
    - ▶ From **client** to **server**
  - **Shared-memory model** create and gain access to memory regions
  - transfer status information
  - attach and detach remote devices

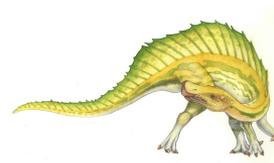




# Types of System Calls (Cont.)

---

- Protection
  - Control access to resources
  - Get and set permissions
  - Allow and deny user access





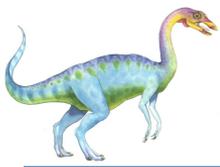
# Examples of Windows and Unix System Calls

## EXAMPLES OF WINDOWS AND UNIX SYSTEM CALLS

The following illustrates various equivalent system calls for Windows and UNIX operating systems.

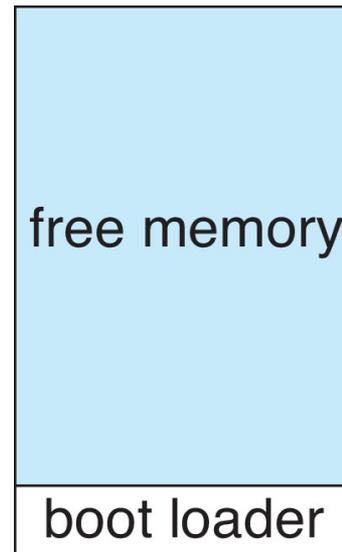
	Windows	Unix
<b>Process control</b>	CreateProcess() ExitProcess() WaitForSingleObject()	fork() exit() wait()
<b>File management</b>	CreateFile() ReadFile() WriteFile() CloseHandle()	open() read() write() close()
<b>Device management</b>	SetConsoleMode() ReadConsole() WriteConsole()	ioctl() read() write()
<b>Information maintenance</b>	GetCurrentProcessID() SetTimer() Sleep()	getpid() alarm() sleep()
<b>Communications</b>	CreatePipe() CreateFileMapping() MapViewOfFile()	pipe() shm_open() mmap()
<b>Protection</b>	SetFileSecurity() InitializeSecurityDescriptor() SetSecurityDescriptorGroup()	chmod() umask() chown()





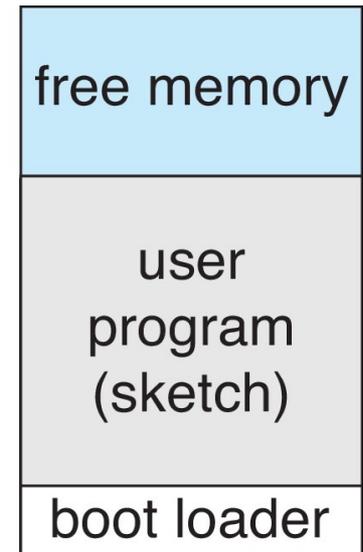
# Example: Arduino

- Single-tasking
- No operating system
- Programs (sketch) loaded via USB into flash memory
- Single memory space
- Boot loader loads program
- <https://www.arduino.cc/>



(a)

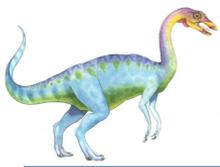
At system startup



(b)

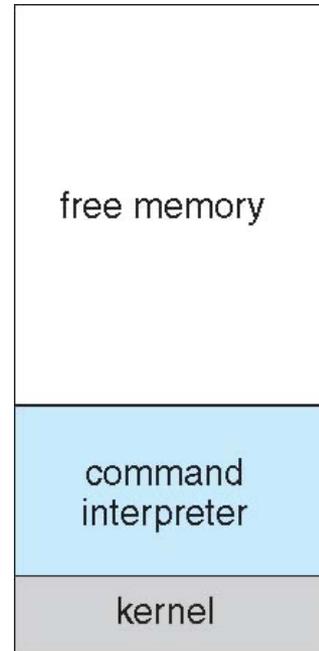
running a program



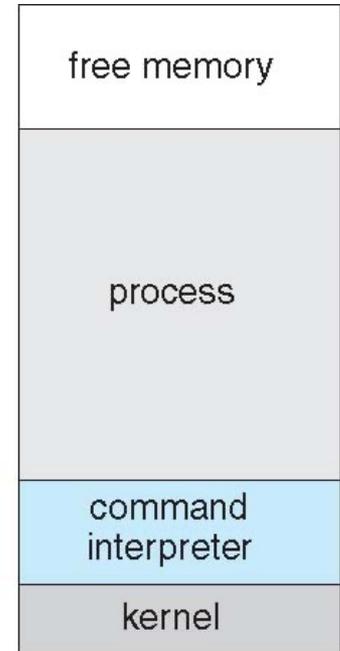


# Example: MS-DOS

- Single-tasking
- Shell invoked when system booted
- Simple method to run program
  - No process created
- Single memory space
- Loads program into memory, overwriting all but the kernel
- Program exit -> shell reloaded



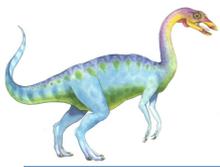
(a)



(b)

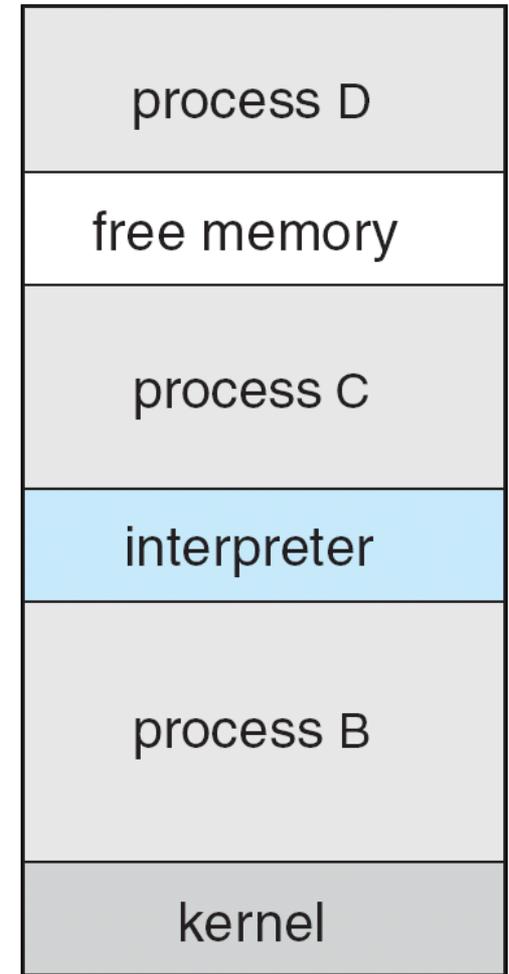
(a) At system startup (b) running a program

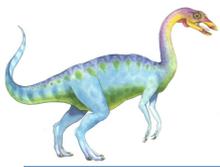




# Example: FreeBSD

- Unix variant
- Multitasking
- User login -> invoke user's choice of shell
- Shell executes `fork()` system call to create process
  - Executes `exec()` to load program into process
  - Shell waits for process to terminate or continues with user commands
- Process exits with
  - code of 0 – no error or
  - code > 0 – error code



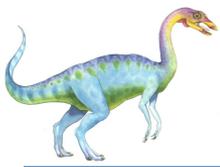


# System Services

---

- System programs provide a convenient environment for program development and execution. They can be divided into:
  - File manipulation
  - Status information sometimes stored in a File modification
  - Programming language support
  - Program loading and execution
  - Communications
  - Background services
  - Application programs
  
- Most users' view of the operation system is defined by system programs, not the actual system calls



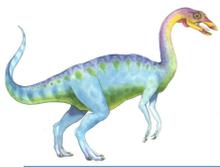


# System Programs

---

- Provide a convenient environment for program development and execution
  - Some of them are simply user interfaces to system calls; others are considerably more complex
- **File management** - Create, delete, copy, rename, print, dump, list, and generally manipulate files and directories
- **Status information**
  - Some ask the system for info - date, time, amount of available memory, disk space, number of users
  - Others provide detailed performance, logging, and debugging information
  - Typically, these programs format and print the output to the terminal or other output devices
  - Some systems implement a **registry** - used to store and retrieve configuration information





# System Programs (Cont.)

---

## ■ File modification

- Text editors to create and modify files
- Special commands to search contents of files or perform transformations of the text

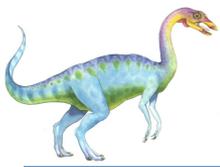
## ■ Programming-language support - Compilers, assemblers, debuggers and interpreters sometimes provided

## ■ Program loading and execution- Absolute loaders, relocatable loaders, linkage editors, and overlay-loaders, debugging systems for higher-level and machine language

## ■ Communications - Provide the mechanism for creating virtual connections among processes, users, and computer systems

- Allow users to send messages to one another's screens, browse web pages, send electronic-mail messages, log in remotely, transfer files from one machine to another





# System Programs (Cont.)

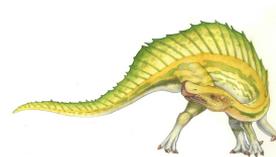
---

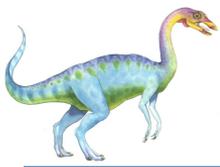
## ■ Background Services

- Launch at boot time
  - ▶ Some for system startup, then terminate
  - ▶ Some from system boot to shutdown
- Provide facilities like disk checking, process scheduling, error logging, printing
- Run in user context not kernel context
- Known as **services**, **subsystems**, **daemons**

## ■ Application programs

- Don't pertain to system
- Run by users
- Not typically considered part of OS
- Launched by command line, mouse click, finger poke



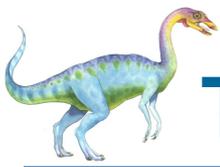


# Linkers and Loaders

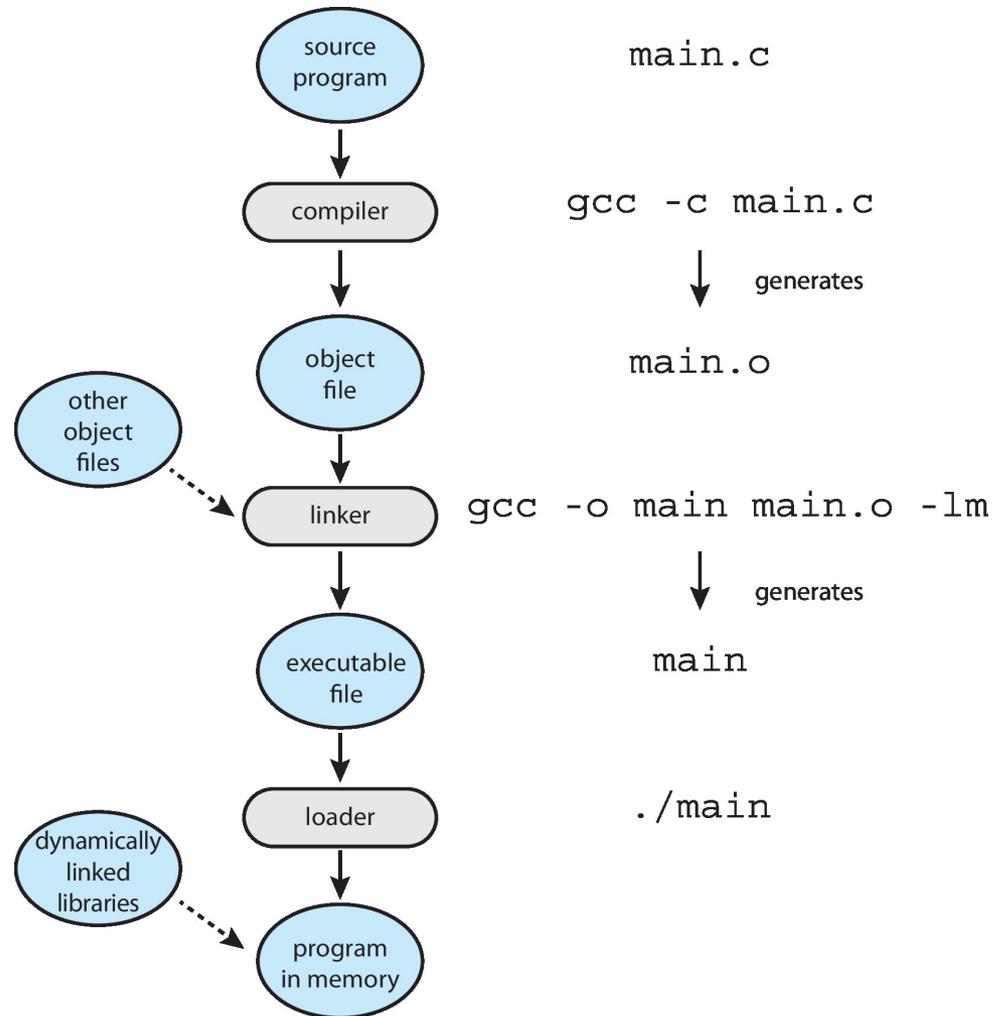
---

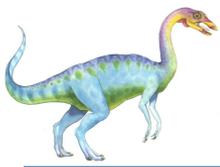
- Source code compiled into object files designed to be loaded into any physical memory location – **relocatable object file**
- **Linker** combines these into single binary **executable** file
  - Also brings in libraries
- Program resides on secondary storage as binary executable
- Must be brought into memory by **loader** to be executed
  - **Relocation** assigns final addresses to program parts and adjusts code and data in program to match those addresses
- Modern general purpose systems don't link libraries into executables
  - Rather, **dynamically linked libraries** (in Windows, **DLLs**) are loaded as needed, shared by all that use the same version of that same library (loaded once)
- Object, executable files have standard formats, so operating system knows how to load and start them





# The Role of the Linker and Loader



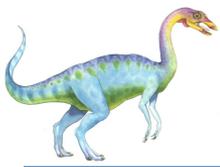


# Why Applications are Operating System Specific

---

- Apps compiled on one system usually not executable on other operating systems
- Each operating system provides its own unique system calls
  - Own file formats, etc
- Apps can be multi-operating system
  - Written in interpreted language like Python, Ruby, and interpreter available on multiple operating systems
  - App written in language that includes a VM containing the running app (like Java)
  - Use standard language (like C), compile separately on each operating system to run on each
- **Application Binary Interface (ABI)** is architecture equivalent of API, defines how different components of binary code can interface for a given operating system on a given architecture, CPU, etc



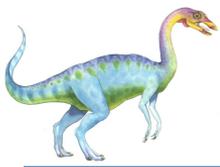


# Operating System Design and Implementation

---

- Design and Implementation of OS not “solvable”, but some approaches have proven successful
- Internal structure of different Operating Systems can vary widely
- Start by defining goals and specifications
- Affected by choice of hardware, type of system
- **User** goals and **System** goals
  - User goals – operating system should be convenient to use, easy to learn, reliable, safe, and fast
  - System goals – operating system should be easy to design, implement, and maintain, as well as flexible, reliable, error-free, and efficient





# Operating System Design and Implementation (Cont.)

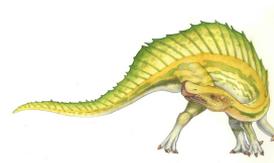
---

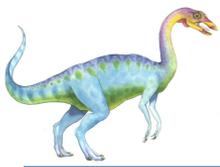
- Important principle to separate

**Policy:** *What* will be done?

**Mechanism:** *How* to do it?

- Mechanisms determine how to do something, policies decide what will be done
  - The separation of policy from mechanism is a very important principle, it allows maximum flexibility if policy decisions are to be changed later
- Specifying and designing OS is highly creative task of **software engineering**

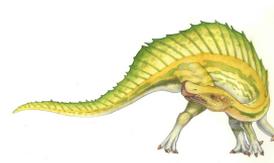


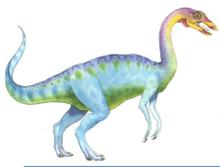


# Implementation

---

- Much variation
  - Early OSes in assembly language
  - Then system programming languages like Algol, PL/1
  - Now C, C++
- Actually usually a mix of languages
  - Lowest levels in assembly
  - Main body in C
  - Systems programs in C, C++, scripting languages like PERL, Python, shell scripts
- More high-level language easier to **port** to other hardware
  - But slower
- **Emulation** can allow an OS to run on non-native hardware





# Operating System Structure

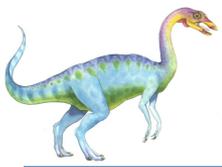
---

- Why careful engineering ?
- Example: LINUX
- September 1991, Linux version 0.01 was released. It had 10,239 lines of code.
- [...]
- 25 January 1999 - Linux 2.2.0 was released (1,800,847 lines of code).
- 4 January 2001 - Linux 2.4.0 was released (3,377,902 lines of code).
- 17 December 2003 - Linux 2.6.0 was released (5,929,913 lines of code).
- 24 December 2008 - Linux 2.6.28 was released (10,195,402 lines of code)
- 3 December 2009 - Linux 2.6.32 was released (12,606,910 lines of code)
- 4 January 2011 - Linux 2.6.37 was released (13,996,612 lines of code).
- 2012, Linux 3.2 was released (14,998,651 lines of code)
- 2013, Linux 3.10 release had 15,803,499 lines of code
- 4.15 (Feb 2018), >20 mio lines of code,
- 5.5 (As of 2020), >27 mio lines of code

<https://phoronix.com/misc/linux-eoy2019/index.html>



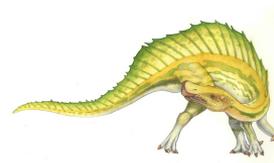


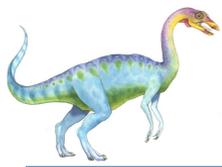


# Operating System Structure

---

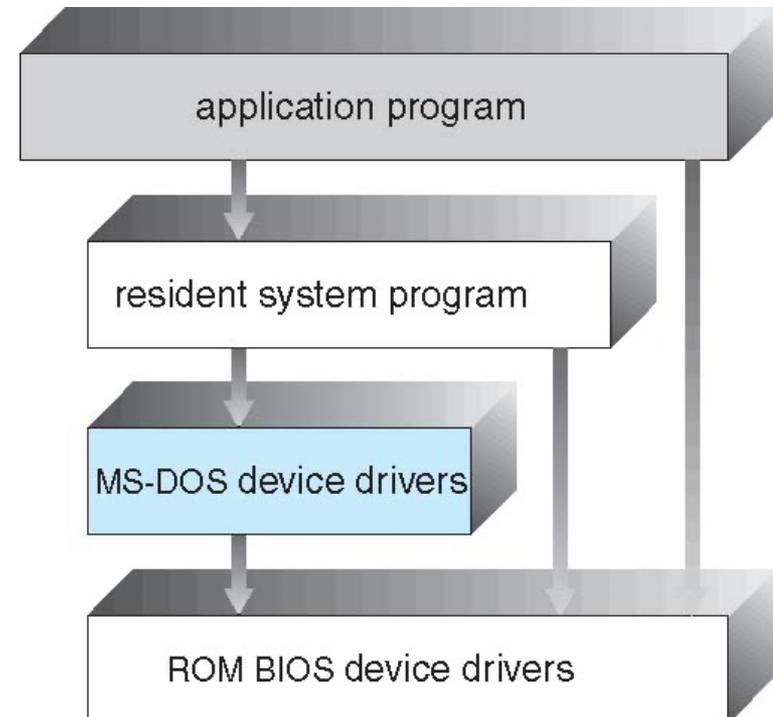
- General-purpose OS is very large program
- Various ways to structure one as follows





# Simple Structure

- I.e. MS-DOS – written to provide the most functionality in the least space
  - Not divided into modules
  - Although MS-DOS has some structure, its interfaces and levels of functionality are not well separated

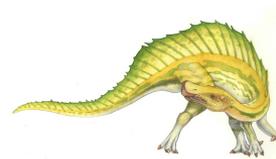


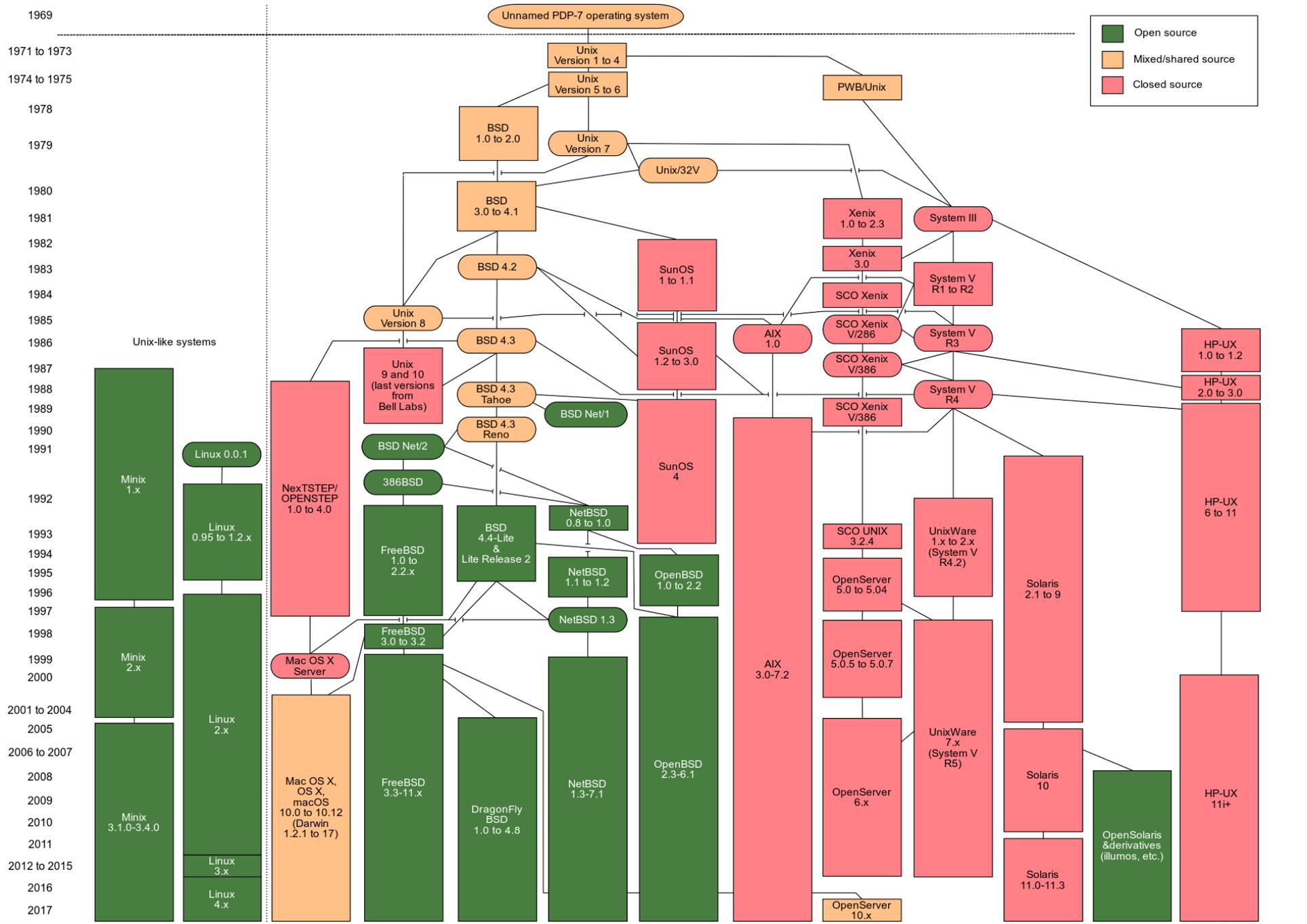


# UNIX

---

- UNIX – limited by hardware functionality, the original UNIX operating system had limited structuring. The UNIX OS consists of two separable parts
  - Systems programs
  - The kernel
    - ▶ Consists of everything below the system-call interface and above the physical hardware
    - ▶ Provides the file system, CPU scheduling, memory management, and other operating-system functions; a large number of functions for one level

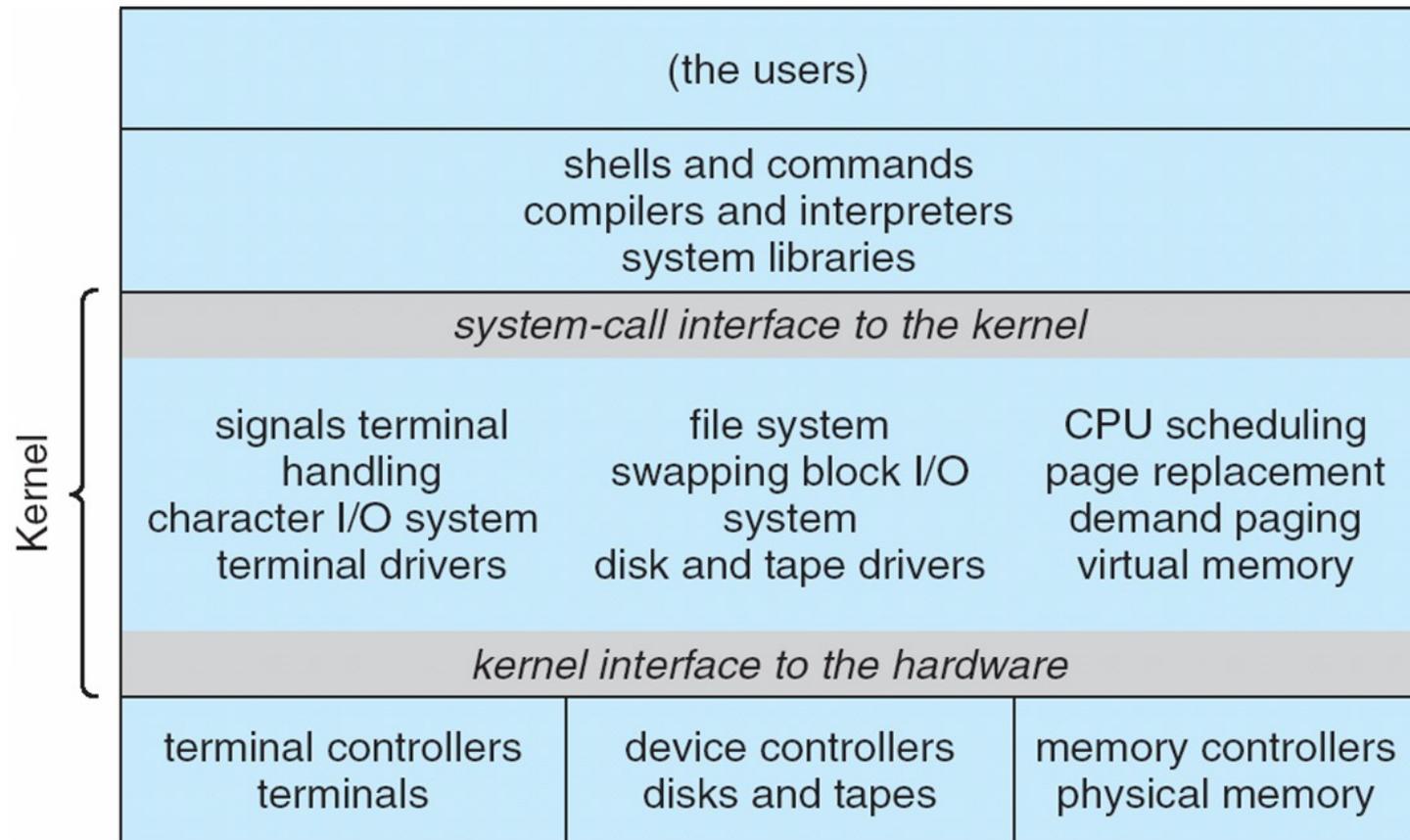


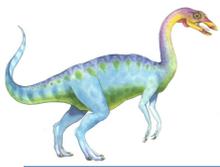




# Traditional UNIX System Structure

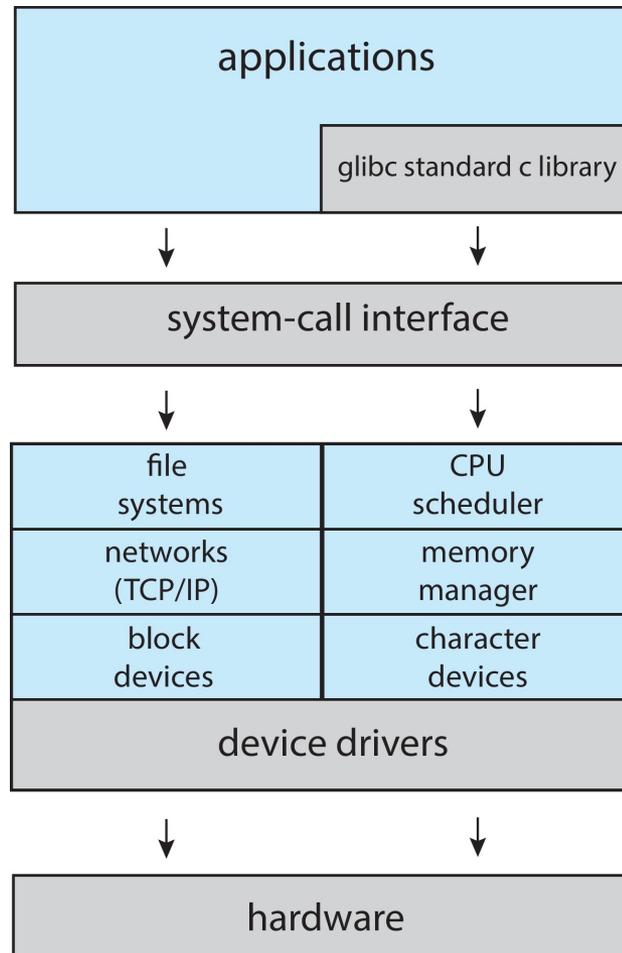
Beyond simple but not fully layered

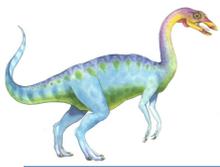




# Linux System Structure

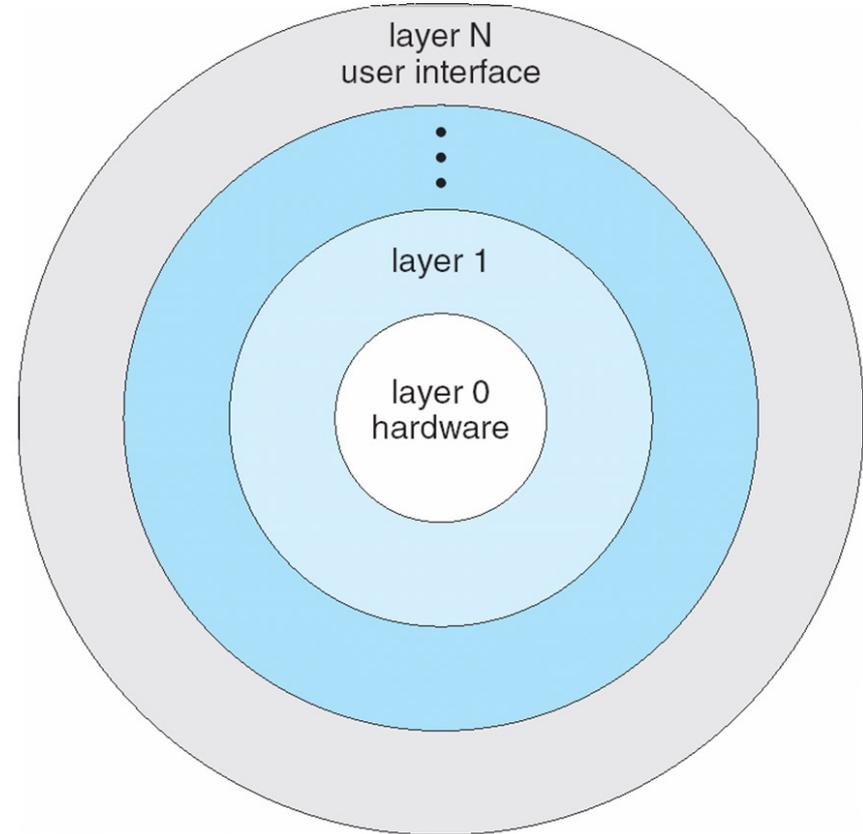
Monolithic plus modular design

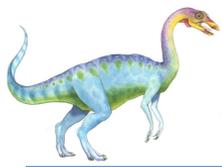




# Layered Approach

- The operating system is divided into a number of layers (levels), each built on top of lower layers. The bottom layer (layer 0), is the hardware; the highest (layer N) is the user interface.
- With modularity, layers are selected such that each uses functions (operations) and services of only lower-level layers





# Microkernels

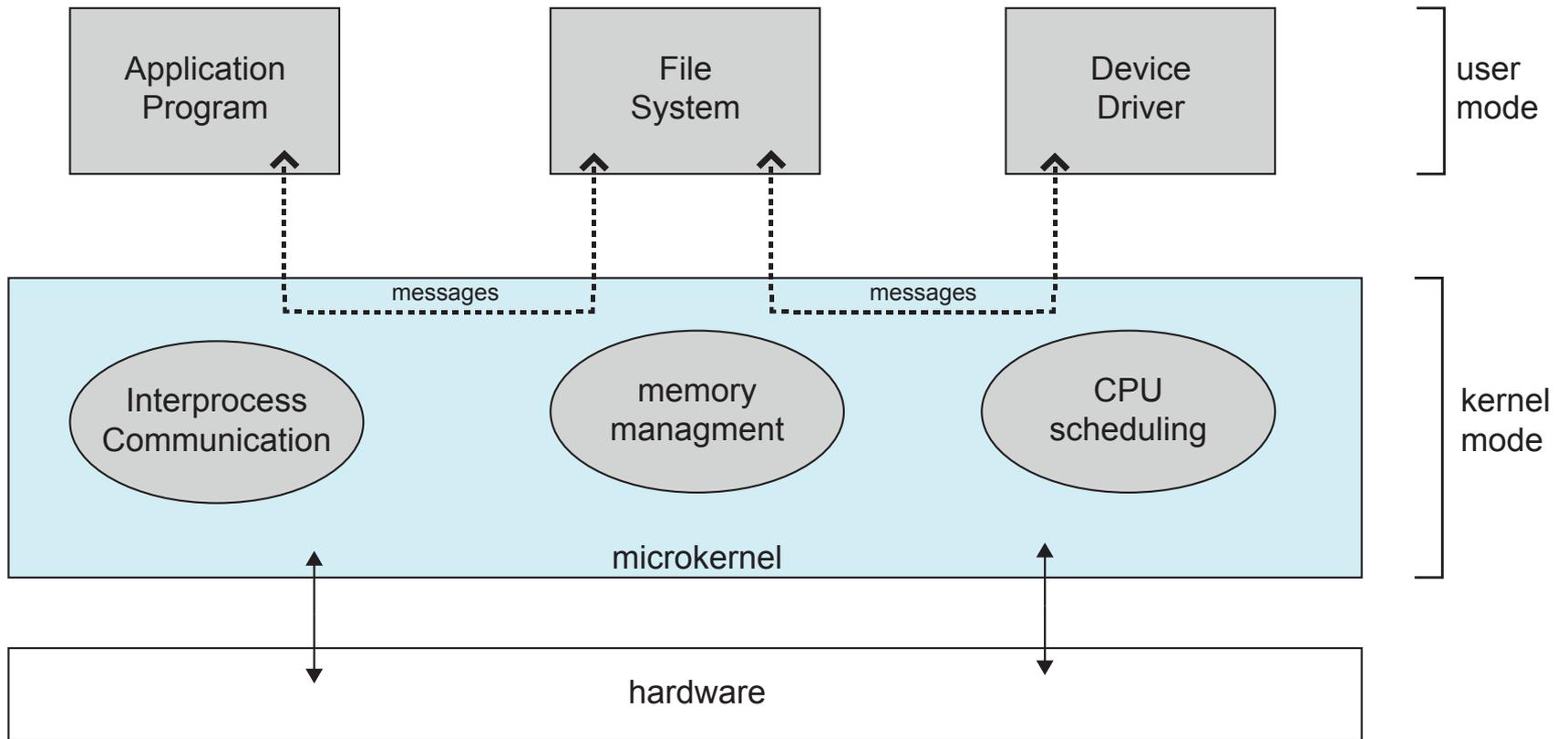
---

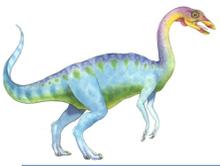
- Moves as much from the kernel into user space
- **Mach** example of **microkernel**
  - Mac OS X kernel (**Darwin**) partly based on Mach
- Communication takes place between user modules using **message passing**
- Benefits:
  - Easier to extend a microkernel
  - Easier to port the operating system to new architectures
  - More reliable (less code is running in kernel mode)
  - More secure
- Detriments:
  - Performance overhead of user space to kernel space communication





# Microkernel System Structure

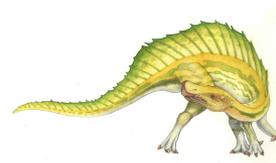


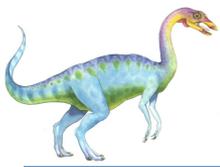


# Modules

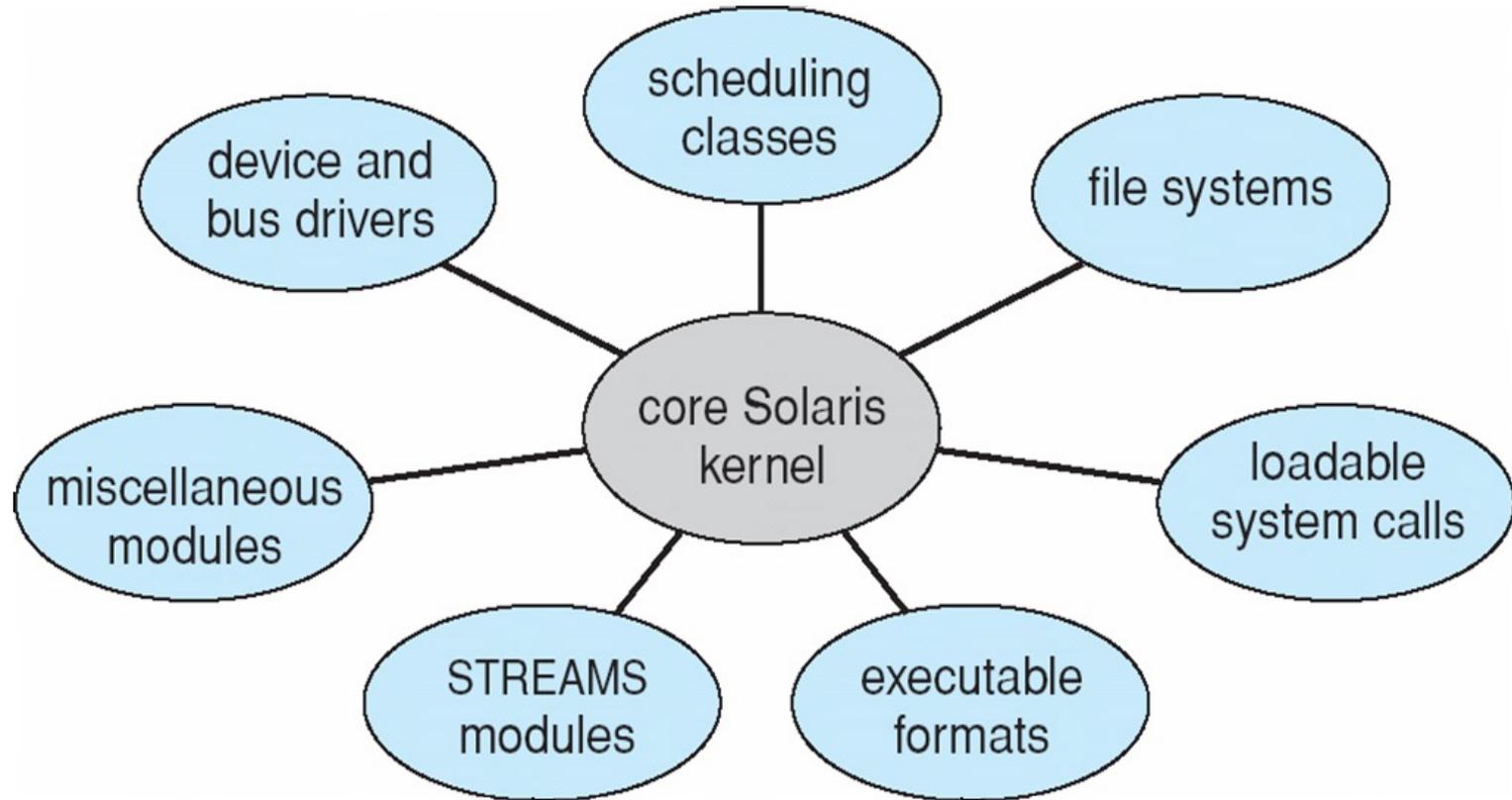
---

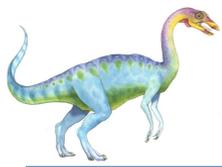
- Most modern operating systems implement **loadable kernel modules**
  - Uses object-oriented approach
  - Each core component is separate
  - Each talks to the others over known interfaces
  - Each is loadable as needed within the kernel
- Overall, similar to layers but with more flexible
  - Linux, Solaris, etc





# Solaris Modular Approach





# Hybrid Systems

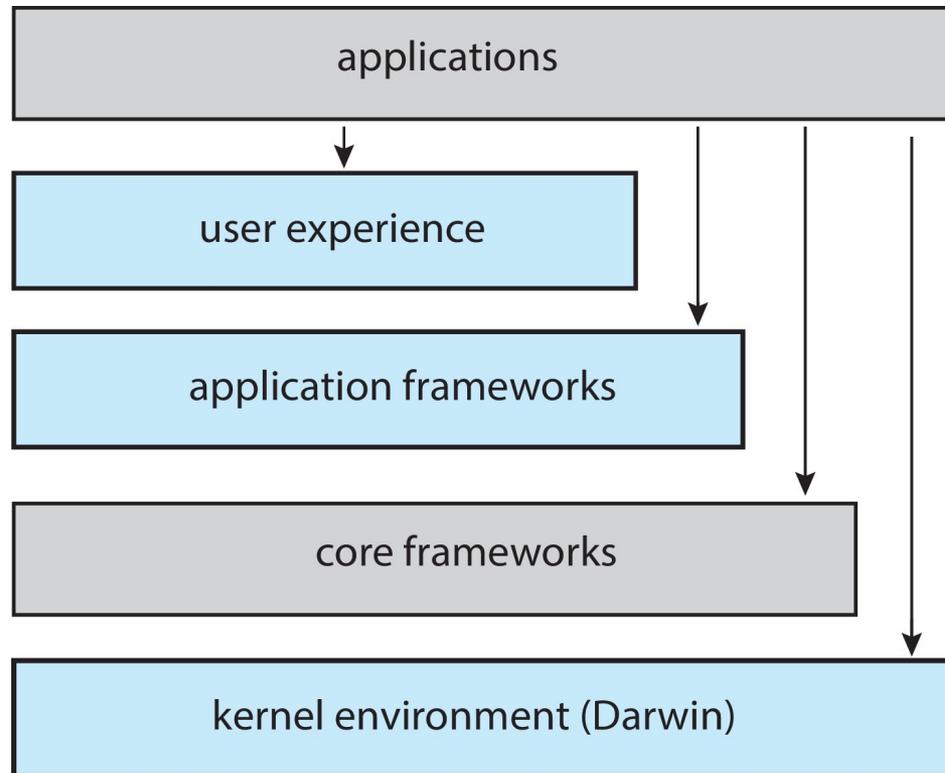
---

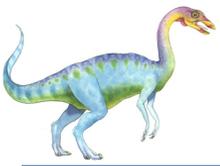
- Most modern operating systems actually not one pure model
  - Hybrid combines multiple approaches to address performance, security, usability needs
  - Linux and Solaris kernels in kernel address space, so monolithic, plus modular for dynamic loading of functionality
  - Windows mostly monolithic, plus microkernel for different subsystem *personalities*
- Apple Mac OS X hybrid, layered, **Aqua** UI plus **Cocoa** programming environment
  - Next slide: kernel consisting of Mach microkernel and BSD Unix parts, plus I/O kit and dynamically loadable modules (called **kernel extensions**)



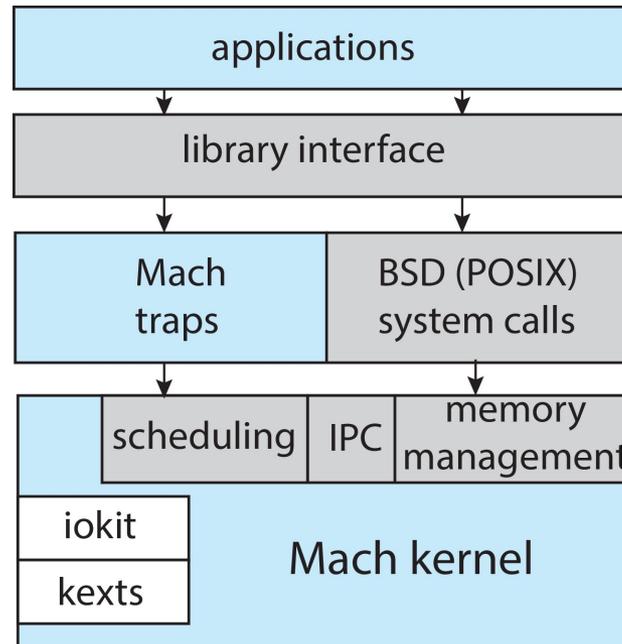


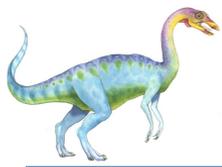
# macOS and iOS Structure



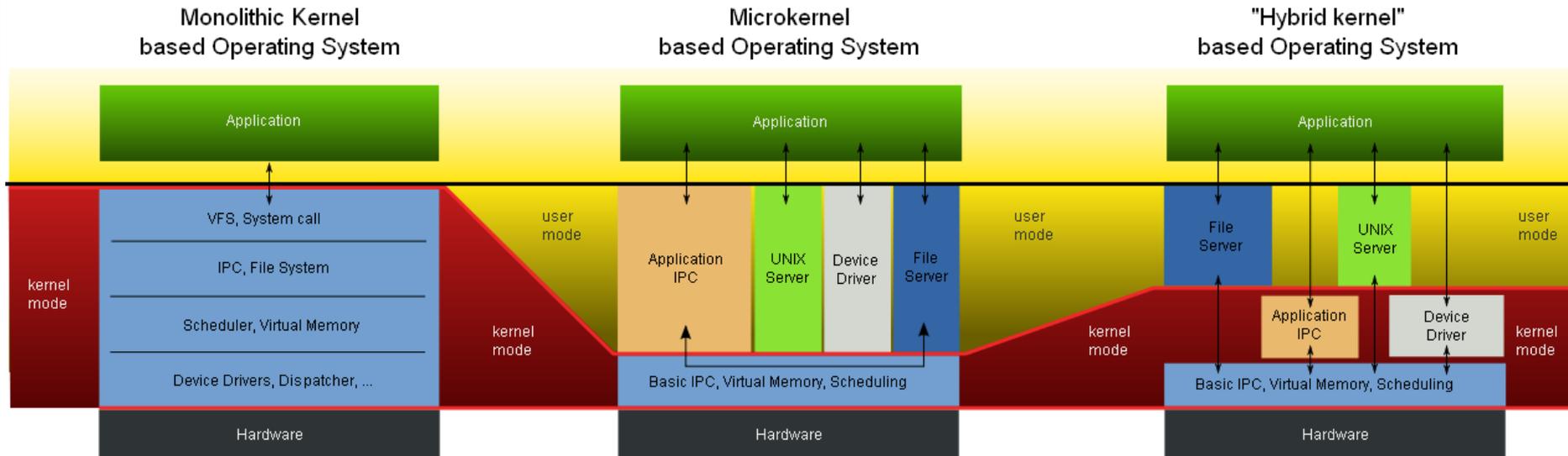


# Darwin





# Hybrid Kernels



- A **hybrid kernel** is a kernel architecture based on combining aspects of microkernel and monolithic kernel architectures.
- Linux Torvalds said of this issue:  
As to the whole 'hybrid kernel' thing — it's just marketing. It's 'oh, those microkernels had good PR, how can we try to get good PR for our working kernel? Oh, I know, let's use a cool name and try to imply that it has all the PR advantages that that other system has'





# iOS

- Apple mobile OS for *iPhone*, *iPad*
  - Structured on Mac OS X, added functionality
  - Does not run OS X applications natively
    - ▶ Also runs on different CPU architecture (ARM vs. Intel)
  - **Cocoa Touch** Objective-C API for developing apps
  - **Media services** layer for graphics, audio, video
  - **Core services** provides cloud computing, databases
  - Core operating system, based on Mac OS X kernel

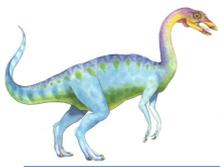
Cocoa Touch

Media Services

Core Services

Core OS

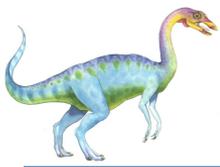




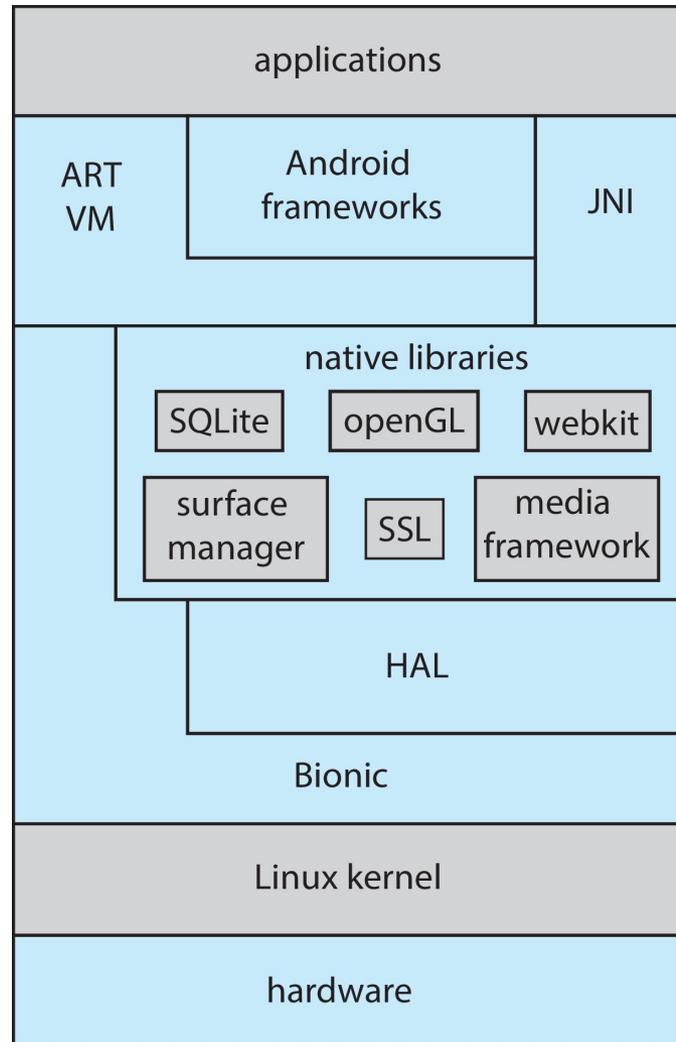
# Android

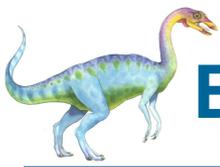
- Developed by Open Handset Alliance (mostly Google)
  - Open Source
- Similar stack to IOS
- Based on Linux kernel but modified
  - Provides process, memory, device-driver management
  - Adds power management
- Runtime environment includes core set of libraries and Dalvik virtual machine
  - Apps developed in Java plus Android API
    - ▶ Java class files compiled to Java bytecode then translated to executable than runs in ~~Dalvik-VM~~ Android Runtime (ART)
- Libraries include frameworks for web browser (webkit), database (SQLite), multimedia, smaller libc





# Android Architecture



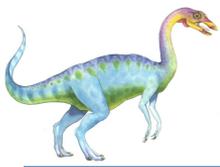


# Building and Booting an Operating System

---

- Operating systems generally designed to run on a class of systems with variety of peripherals
- Commonly, operating system already installed on purchased computer
  - But can build and install some other operating systems
  - If generating an operating system from scratch
    - ▶ Write the operating system source code
    - ▶ Configure the operating system for the system on which it will run
    - ▶ Compile the operating system
    - ▶ Install the operating system
    - ▶ Boot the computer and its new operating system





# Building and Booting Linux

---

- Download Linux source code (<http://www.kernel.org>)
- Configure kernel via “make menuconfig”
- Compile the kernel using “make”
  - Produces `vmlinuz`, the kernel image
  - Compile kernel modules via “make modules”
  - Install kernel modules into `vmlinuz` via “make modules\_install”
  - Install new kernel on the system via “make install”

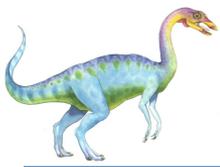




# System Boot

- When power initialized on system, execution starts at a fixed memory location
- Operating system must be made available to hardware so hardware can start it
  - Small piece of code – **bootstrap loader**, **BIOS**, stored in **ROM** or **EEPROM** locates the kernel, loads it into memory, and starts it
  - Sometimes two-step process where **boot block** at fixed location loaded by ROM code, which loads bootstrap loader from disk
  - Modern systems replace BIOS with **Unified Extensible Firmware Interface (UEFI)**
- Common bootstrap loader, **GRUB**, allows selection of kernel from multiple disks, versions, kernel options
- Kernel loads and system is then **running**
- Boot loaders frequently allow various boot states, such as single user mode





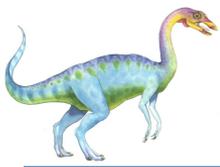
# Operating-System Debugging

---

- **Debugging** is finding and fixing errors, or **bugs**
- OSes generate **log files** containing error information
- Failure of an application can generate **core dump** file capturing memory of the process
- Operating system failure can generate **crash dump** file containing kernel memory
- Beyond crashes, performance tuning can optimize system performance
  - Sometimes using **trace listings** of activities, recorded for analysis
  - **Profiling** is periodic sampling of instruction pointer to look for statistical trends

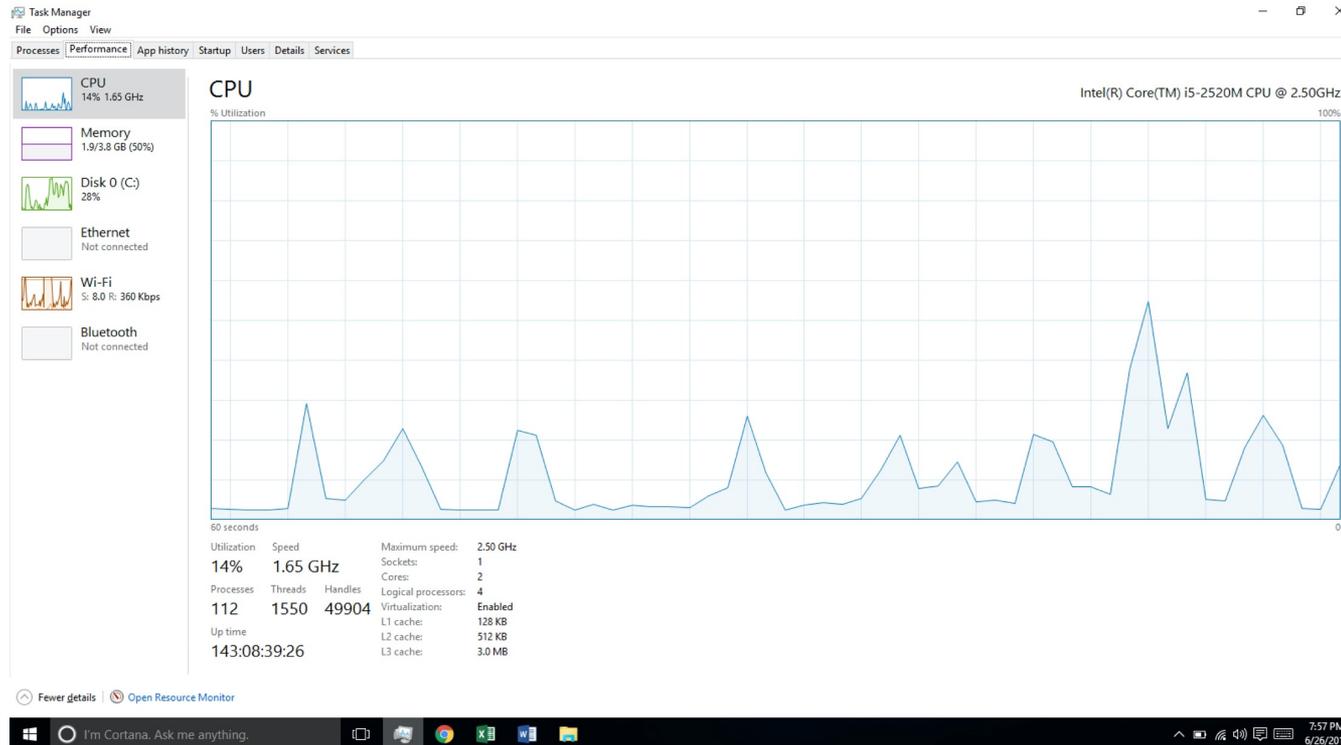
Kernighan's Law: "Debugging is twice as hard as writing the code in the first place. Therefore, if you write the code as cleverly as possible, you are, by definition, not smart enough to debug it."

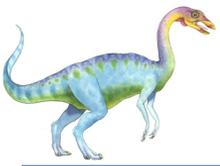




# Performance Tuning

- Improve performance by removing bottlenecks
- OS must provide means of computing and displaying measures of system behavior
- For example, “top” program or Windows Task Manager

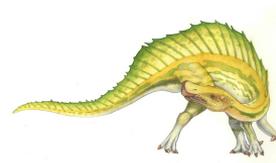




# Tracing

---

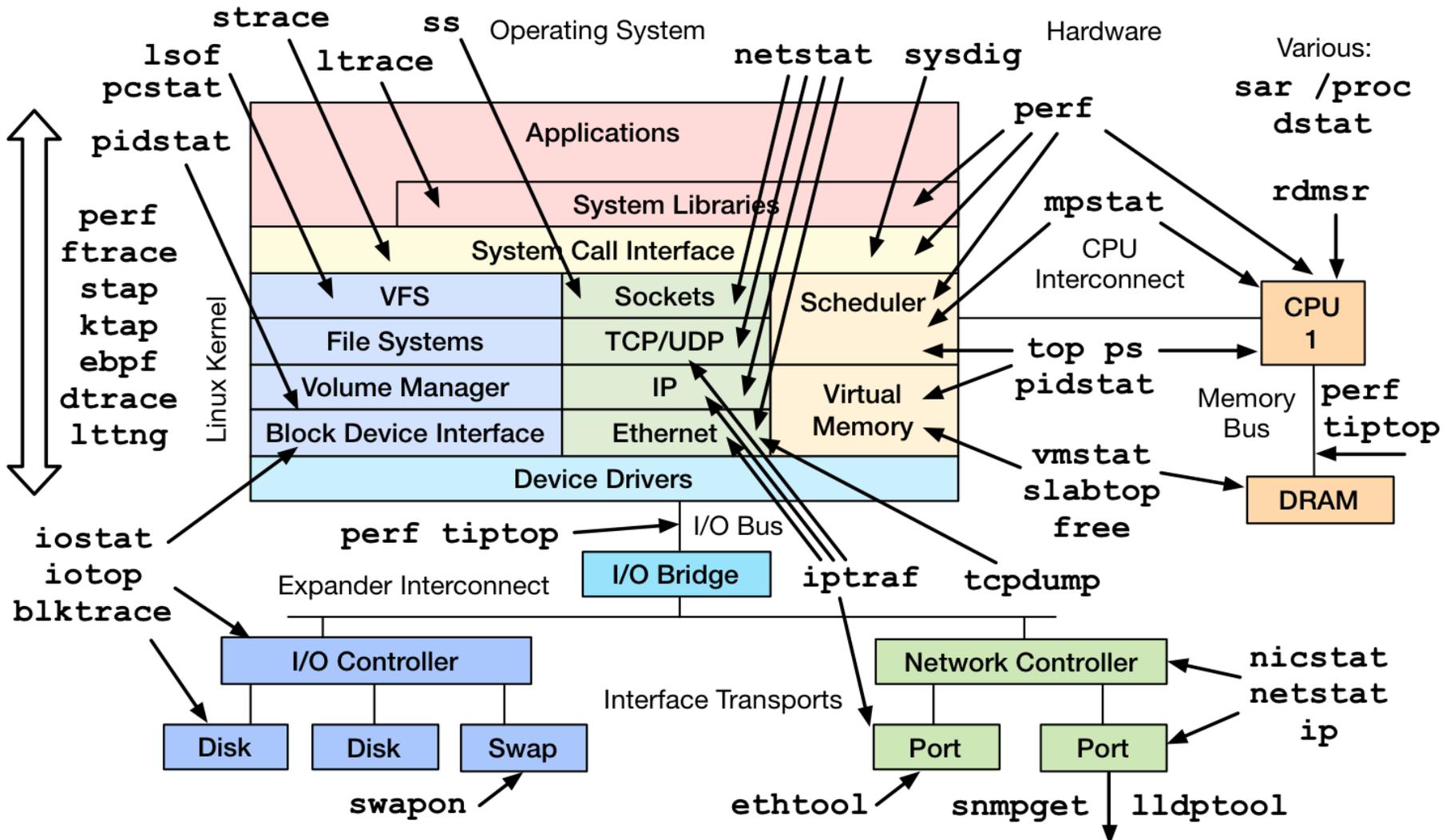
- Collects data for a specific event, such as steps involved in a system call invocation
- Tools include
  - strace – trace system calls invoked by a process
  - gdb – source-level debugger
  - perf – collection of Linux performance tools
  - tcpdump – collects network packets

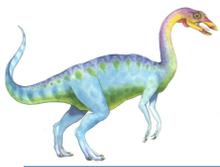




# Linux

## Linux Performance Observability Tools



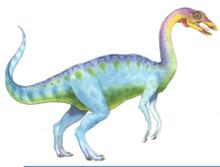


# Overview of DTrace

---

- DTrace tool in Solaris, FreeBSD, Mac OS X, Linux allows live instrumentation on production systems
- Can be used to troubleshoot performance and logic problems in user applications
- Uses dynamic binary instrumentation
- Inserts instrumentation code in running processes
- Specialized C-like language, D
- Terminology
  - **Probes** : points of instrumentation: probes “fire” when code is executed, capturing state data and sending it to consumers of those probes
  - **Providers**: make probes available (technically providers are loadable kernel modules)
- Meant to be used on production systems



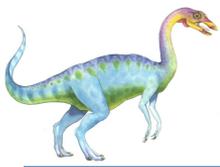


# DTrace Usage Overview

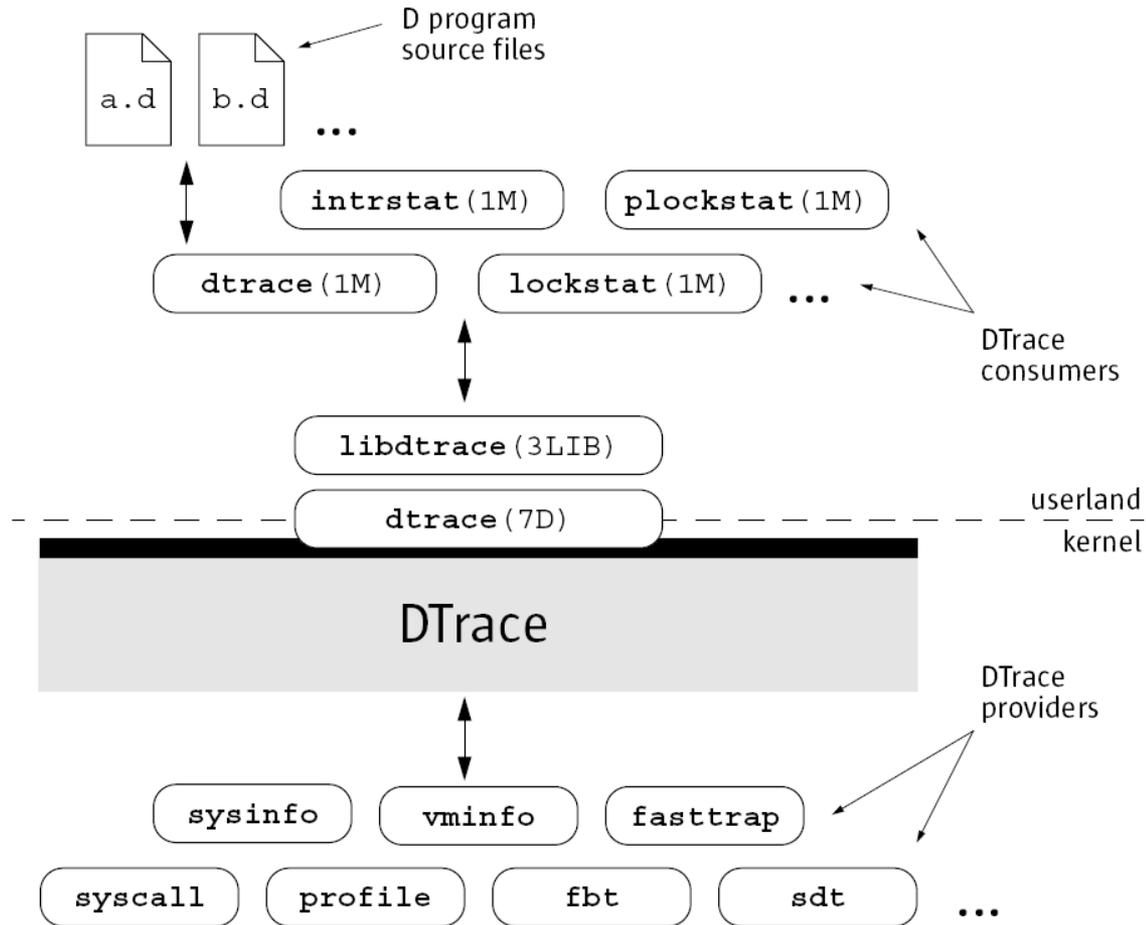
---

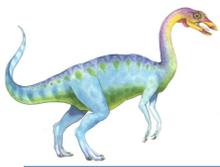
- Users write D programs that collect information at runtime
- Users invoke `dtrace` to insert instrumentation code in kernel and user processes
  - Security mechanism ensures code can be inserted only by authorized users
- When events occur at runtime, user's D code is executed by the DTrace providers, which causes
  - Information to be recorded
  - Data to be printed
  - Whatever else users defines in their D programs!
- Many, many providers exist for getting lots of different data from running programs





# DTrace Architecture





# DTrace – Probe Clause

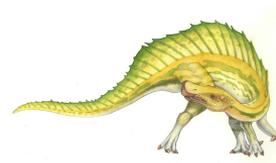
---

- Probe clause:

```
probe-description  
/predicate/  
{  
    action statement  
}
```

- Probe description:

```
provider:module:function:name
```





# Example Toy D Program

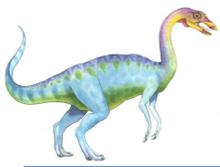
## ■ D code

```
/*
  Count off and report the
  number of seconds elapsed
*/
dtrace:::BEGIN
{
  i = 0;
}
profile:::tick-1sec
{
  i = i + 1;
  trace(i);
}
dtrace:::END
{
  trace(i);
}
```

## ■ Output

```
# dtrace -s counter.d
dtrace: script 'counter.d' matched 3
probes
CPU ID      FUNCTION:NAME
0   25499  :tick-1sec    1
0   25499  :tick-1sec    2
0   25499  :tick-1sec    3
0   25499  :tick-1sec    4
0   25499  :tick-1sec    5
0   25499  :tick-1sec    6
^C
0   2      :END          6
#
```





# More Realistic Program

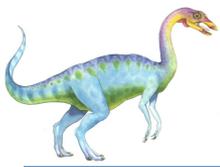
## ■ D code to time `read()` and `write()` syscalls

```
syscall::read:entry,  
syscall::write:entry  
/pid == $1/  
{  
    ts[probefunc] = timestamp;  
}  
syscall::read:return,  
syscall::write:return  
/pid == $1 && ts[probefunc] != 0/  
{  
    printf("%d nsecs", timestamp -  
        ts[probefunc]);  
}
```

## ■ Output

```
# dtrace -s rvertime.d `pgrep -n ksh`  
dtrace: script 'rvertime.d' matched 4  
probes  
  
CPU ID FUNCTION:NAME  
0    33 read:return    22644 nsecs  
0    33 read:return    3382 nsecs  
0    35 write:return   25952 nsecs  
0    33 read:return   916875239 nsecs  
0    35 write:return   27320 nsecs  
0    33 read:return    9022 nsecs  
0    33 read:return    3776 nsecs  
0    35 write:return   17164 nsecs  
...
```

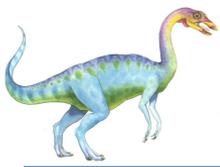




# Solaris 10 dtrace Following System Call

```
# ./all.d 'pgrep xclock' XEventsQueued
dtrace: script './all.d' matched 52377 probes
CPU FUNCTION
0 -> XEventsQueued U
0 -> _XEventsQueued U
0 -> _X11TransBytesReadable U
0 <- _X11TransBytesReadable U
0 -> _X11TransSocketBytesReadable U
0 <- _X11TransSocketBytesreadable U
0 -> ioctl U
0 -> ioctl K
0 -> getf K
0 -> set_active_fd K
0 <- set_active_fd K
0 <- getf K
0 -> get_umatamodel K
0 <- get_umatamodel K
...
0 -> releasef K
0 -> clear_active_fd K
0 <- clear_active_fd K
0 -> cv_broadcast K
0 <- cv_broadcast K
0 <- releasef K
0 <- ioctl K
0 <- ioctl U
0 <- _XEventsQueued U
0 <- XEventsQueued U
```





# Some Example DTrace Providers

---

- `syscall`: makes a probe available at the entry to and return from every system call
- `vminfo`: makes a probe available on VM activity (page out, page faults, etc)
- `profile`: makes a probe available that can run every X milliseconds
- `fpuinfo`: makes a probe available when hardware floating point operations are emulated in software
- Users can also create their own providers by using the DTrace API
  - Ex: provide probes before/after a request is serviced in a web server or database server

-> `dtrace` Live demo on a Solaris and a macOS X system





# BCC

- Debugging interactions between user-level and kernel code nearly impossible without toolset that understands both and an instrument their actions
- BCC (BPF Compiler Collection) is a rich toolkit providing tracing features for Linux
  - See also the original DTrace
- For example, disksnoop.py traces disk I/O activity

TIME(s)	T	BYTES	LAT(ms)
1946.29186700	R	8	0.27
1946.33965000	R	8	0.26
1948.34585000	W	8192	0.96
1950.43251000	R	4096	0.56
1951.74121000	R	4096	0.35

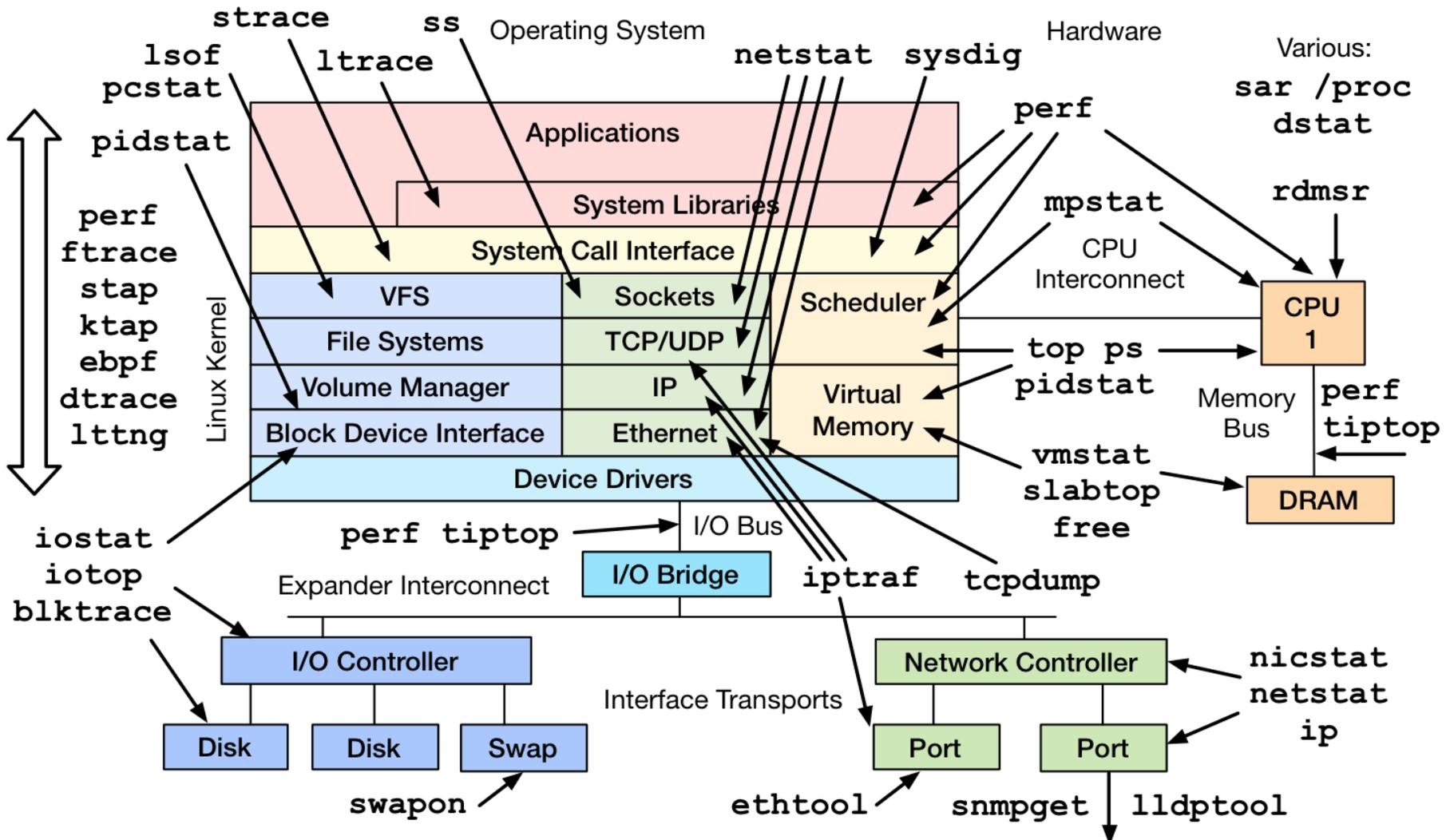
- Many other tools (next slide)

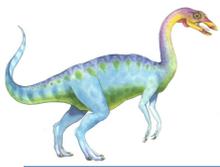




# Linux

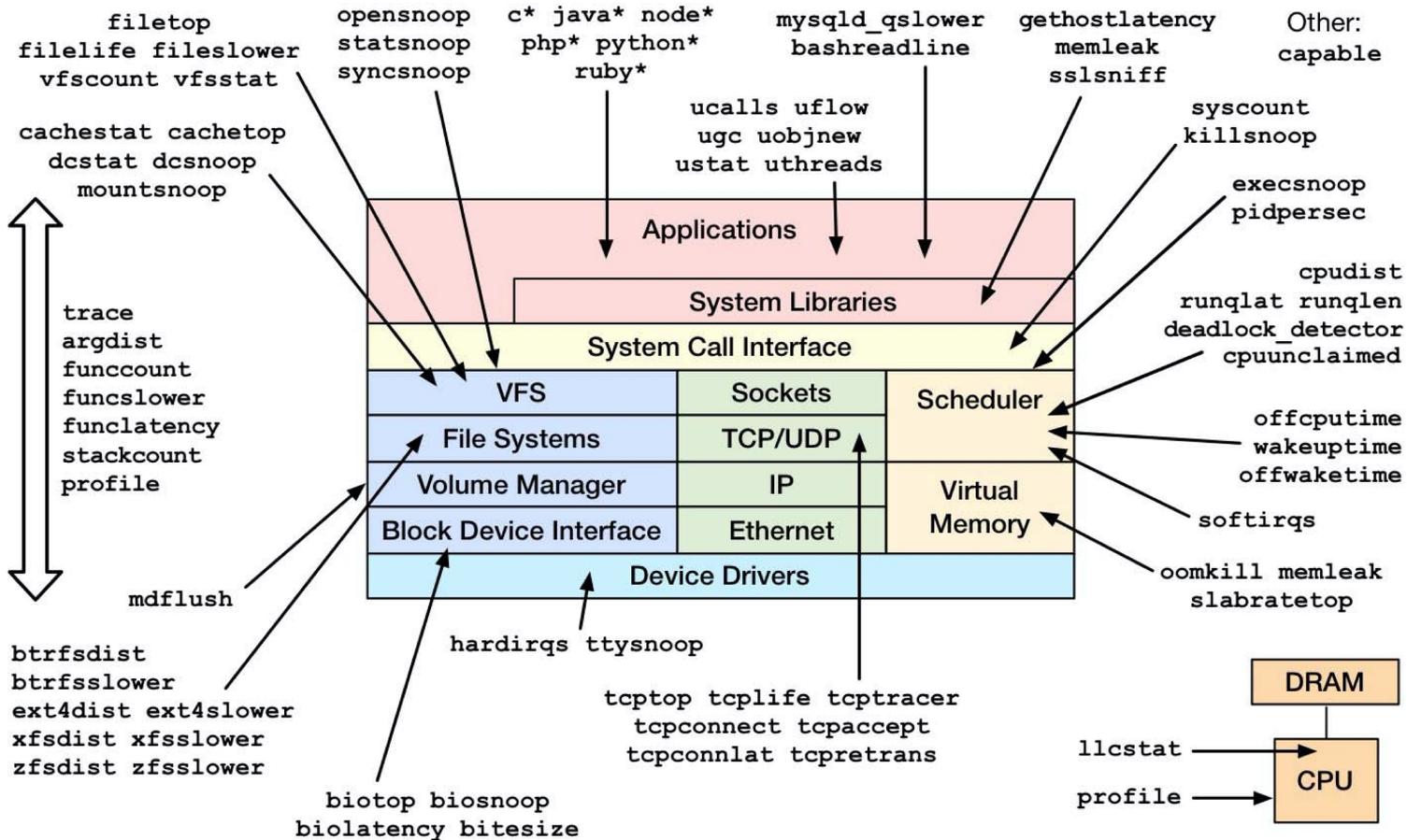
## Linux Performance Observability Tools





# Linux bcc/BPF Tracing Tools

## Linux bcc/BPF Tracing Tools

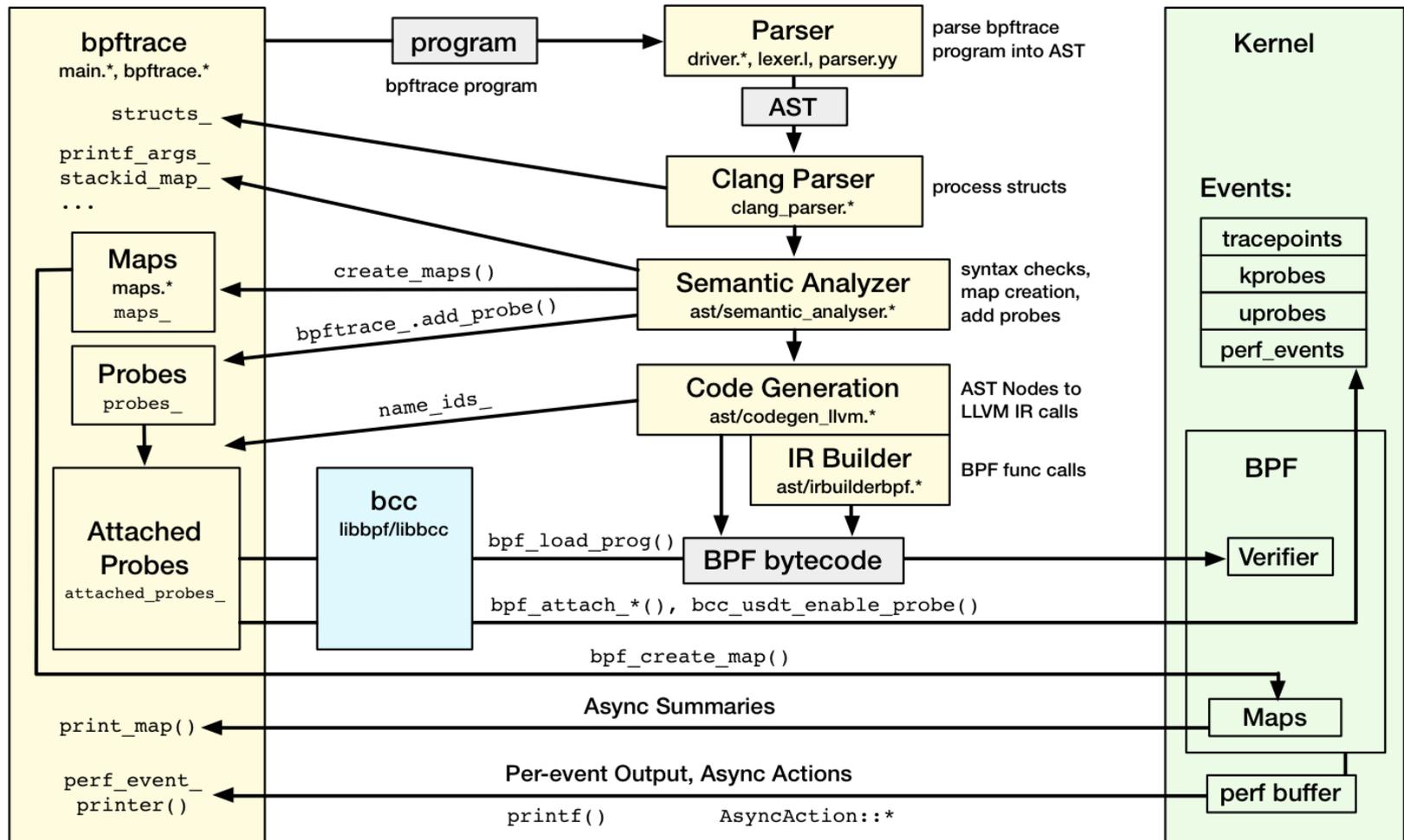


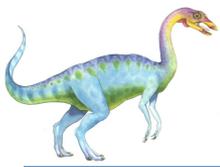
<https://github.com/iovisor/bcc#tools> 2017





# bpfftrace Internals





# Kernel Instrumentation in Linux

---

For example : Perf ( <http://www.brendangregg.com/perf.html> )

- **Hardware Events:** low-level processor activity based on CPU performance counters. For example, CPU cycles, instructions retired, memory stall cycles, level 2 cache misses, etc.
- **Software Events:** low level events based on kernel counters. For example, CPU migrations, minor faults, major faults, etc.
- **Tracepoint Events:** kernel-level events based on the “ftrace” framework. For example, system calls, TCP events, file system I/O, disk I/O, etc. These are grouped into libraries of tracepoints; eg, "sock:" for socket events, "sched:" for CPU scheduler events.
- Details about the events can be collected, including timestamps, the code path that led to it, and other specific details. The capabilities of perf\_events are enormous, and you're likely to only ever use a fraction.
- A 2017 comparasion of tools for gathering of information about the running Linux system:

<https://sourceware.org/systemtap/wiki/SystemtapDtraceComparison>

