# IF2230
# I/O Systems

# Chapter 12: I/O Systems

- I/O Hardware
- Application I/O Interface
- Kernel I/O Subsystem
- Transforming I/O Requests to Hardware Operations
- Streams
- Performance

# Objectives

▸ Explore the structure of an operating system's I/O subsystem

▸ Discuss the principles of I/O hardware and its complexity

▸ Provide details of the performance aspects of I/O hardware and software

# Overview

- I/O management is a major component of operating system design and operation
  - Important aspect of computer operation
  - I/O devices vary greatly
  - Various methods to control them
  - Performance management
  - New types of devices frequent
- Ports, busses, device controllers connect to various devices
- **Device drivers** encapsulate device details
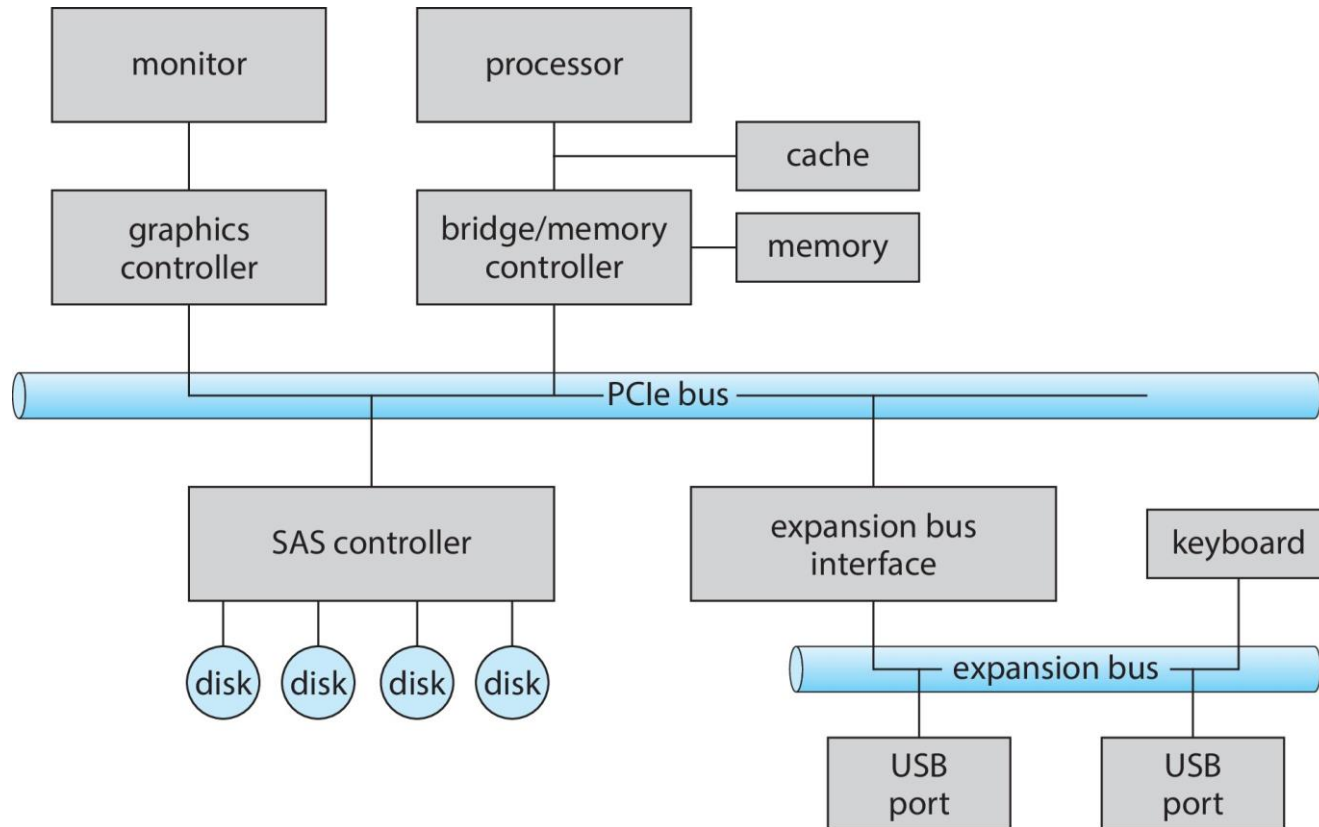  - Present uniform device-access interface to I/O subsystem

# I/O Hardware

- Incredible variety of I/O devices
- Common concepts
  - **Port**
  - **Bus** (**daisy chain** or shared direct access)
    - **PCI** bus common in PCs and servers, PCI Express (**PCIe**)
    - **expansion bus** connects relatively slow devices
    - **Serial-attached SCSI** (**SAS**) common disk interface
  - **Controller** (**host adapter**)
- I/O instructions control devices
- Devices have addresses, used by
  - Direct I/O instructions
  - **Memory-mapped** I/O

# A Typical PC Bus Structure

# Device I/O Port Locations on PCs (partial)

| I/O address range (hexadecimal) | device |
|---|---|
| 000–00F | DMA controller |
| 020–021 | interrupt controller |
| 040–043 | timer |
| 200–20F | game controller |
| 2F8–2FF | serial port (secondary) |
| 320–32F | hard-disk controller |
| 378–37F | parallel port |
| 3D0–3DF | graphics controller |
| 3F0–3F7 | diskette-drive controller |
| 3F8–3FF | serial port (primary) |

# Polling

- For each byte of I/O
    1. Read busy bit from status register until 0
    2. Host sets read or write bit and if write copies data into data-out register
    3. Host sets command-ready bit
    4. Controller sets busy bit, executes transfer
    5. Controller clears busy bit, error bit, command-ready bit when transfer done
- Step 1 is **busy-wait** cycle to wait for I/O from device
    - Reasonable if device is fast
    - But inefficient if device slow
    - CPU switches to other tasks?
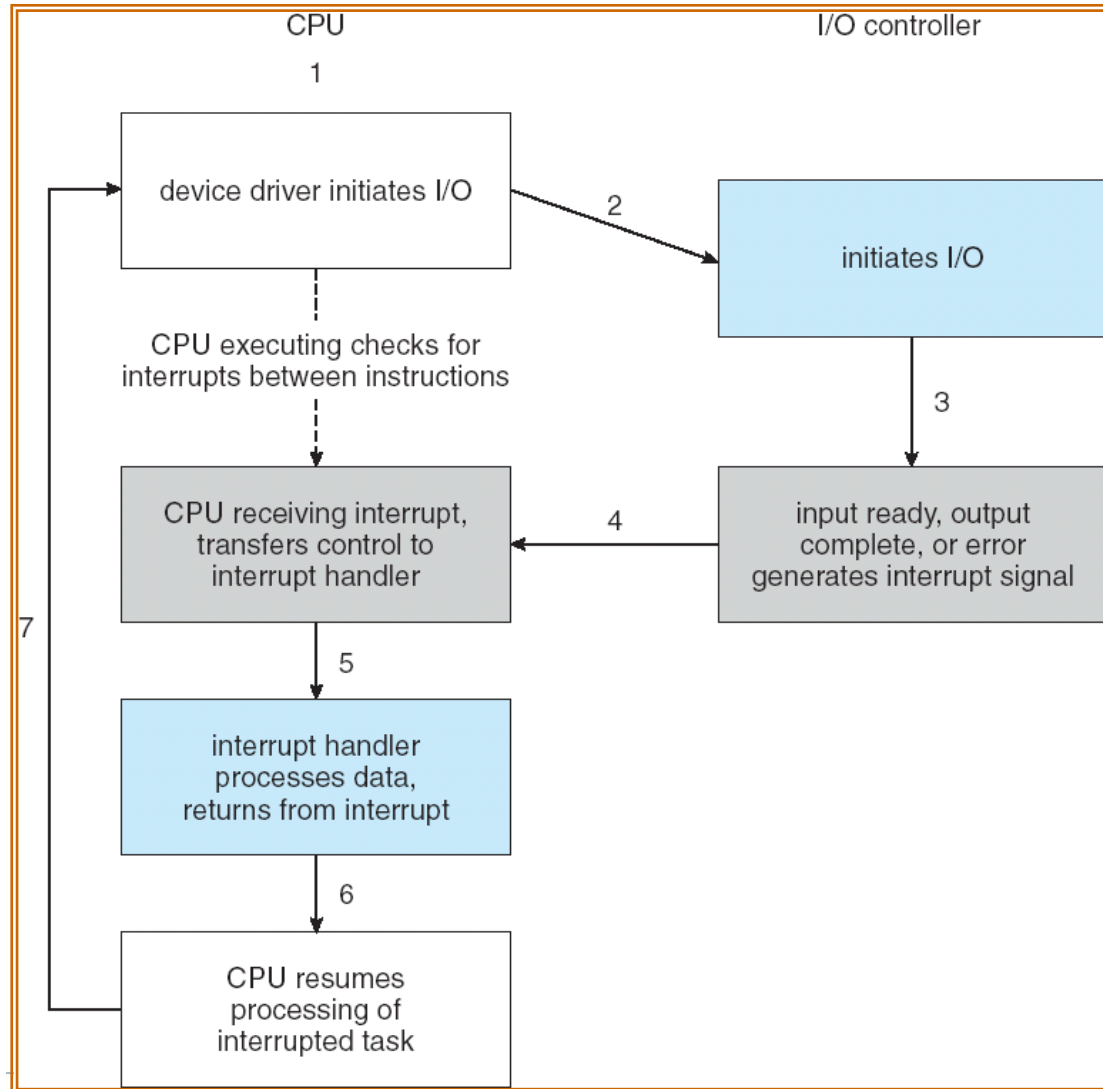        - But if miss a cycle data overwritten / lost

# Interrupts

- Polling can happen in 3 instruction cycles
  - Read status, logical-and to extract status bit, branch if not zero
  - How to be more efficient if non-zero infrequently?
- CPU **Interrupt-request line** triggered by I/O device
  - Checked by processor after each instruction
- **Interrupt handler** receives interrupts
  - **Maskable** to ignore or delay some interrupts
- **Interrupt vector** to dispatch interrupt to correct handler
  - Context switch at start and end
  - Based on priority
  - Some **nonmaskable**
  - Interrupt chaining if more than one device at same interrupt number

# Interrupt-Driven I/O Cycle

# Interrupts (Cont.)

‣ Interrupt mechanism also used for **exceptions**

  ‣ Terminate process, crash system due to hardware error

‣ Page fault executes when memory access error

‣ System call executes via **trap** to trigger kernel to execute request

‣ Multi-CPU systems can process interrupts concurrently

  ‣ If operating system designed to handle it

‣ Used for time-sensitive processing, frequent, must be fast

# Latency

▸ Stressing interrupt management because even single-user systems manage hundreds or interrupts per second and servers hundreds of thousands

▸ For example, a quiet macOS desktop generated 23,000 interrupts over 10 seconds

```
Fri Nov 25 13:55:59                                    0:00:10
                        SCHEDULER       INTERRUPTS
------------------------------------------------------
total_samples                  13           22998

delays <  10 usecs             12           16243
delays <  20 usecs              1            5312
delays <  30 usecs              0             473
delays <  40 usecs              0             590
delays <  50 usecs              0              61
delays <  60 usecs              0             317
delays <  70 usecs              0               2
delays <  80 usecs              0               0
delays <  90 usecs              0               0
delays < 100 usecs              0               0
total   < 100 usecs            13           22998
```
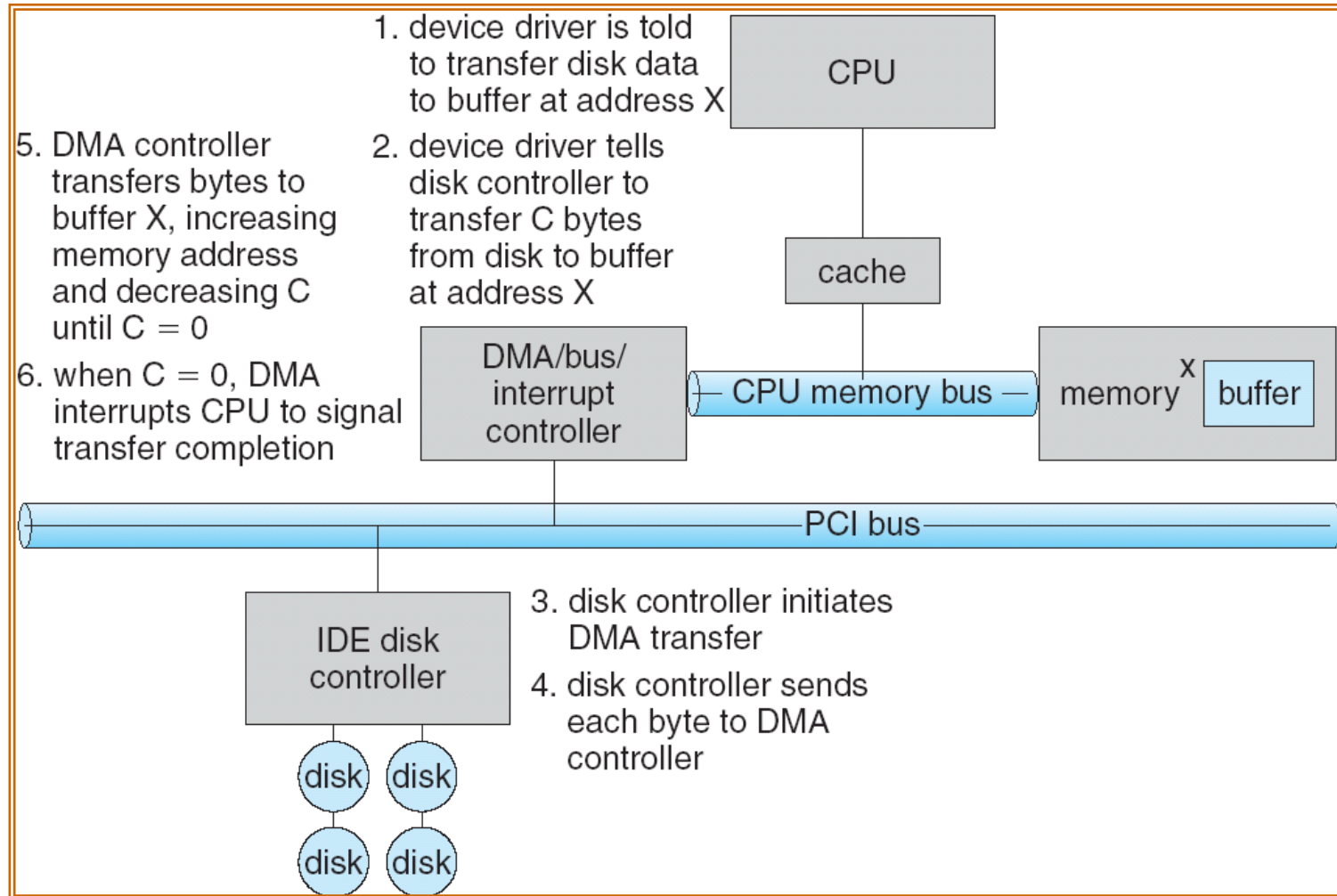
# Intel Pentium Processor Event-Vector Table

| vector number | description |
| --- | --- |
| 0 | divide error |
| 1 | debug exception |
| 2 | null interrupt |
| 3 | breakpoint |
| 4 | INTO-detected overflow |
| 5 | bound range exception |
| 6 | invalid opcode |
| 7 | device not available |
| 8 | double fault |
| 9 | coprocessor segment overrun (reserved) |
| 10 | invalid task state segment |
| 11 | segment not present |
| 12 | stack fault |
| 13 | general protection |
| 14 | page fault |
| 15 | (Intel reserved, do not use) |
| 16 | floating-point error |
| 17 | alignment check |
| 18 | machine check |
| 19–31 | (Intel reserved, do not use) |
| 32–255 | maskable interrupts |

# Direct Memory Access

▸ Used to avoid **programmed I/O** (one byte at a time) for large data movement

▸ Requires **DMA** controller

▸ Bypasses CPU to transfer data directly between I/O device and memory

▸ OS writes DMA command block into memory

  ▸ Source and destination addresses
  ▸ Read or write mode
  ▸ Count of bytes
  ▸ Writes location of command block to DMA controller
  ▸ Bus mastering of DMA controller – grabs bus from CPU
    ▸ **Cycle stealing** from CPU but still much more efficient
  ▸ When done, interrupts to signal completion

▸ Version that is aware of virtual addresses can be even more efficient - **DVMA**

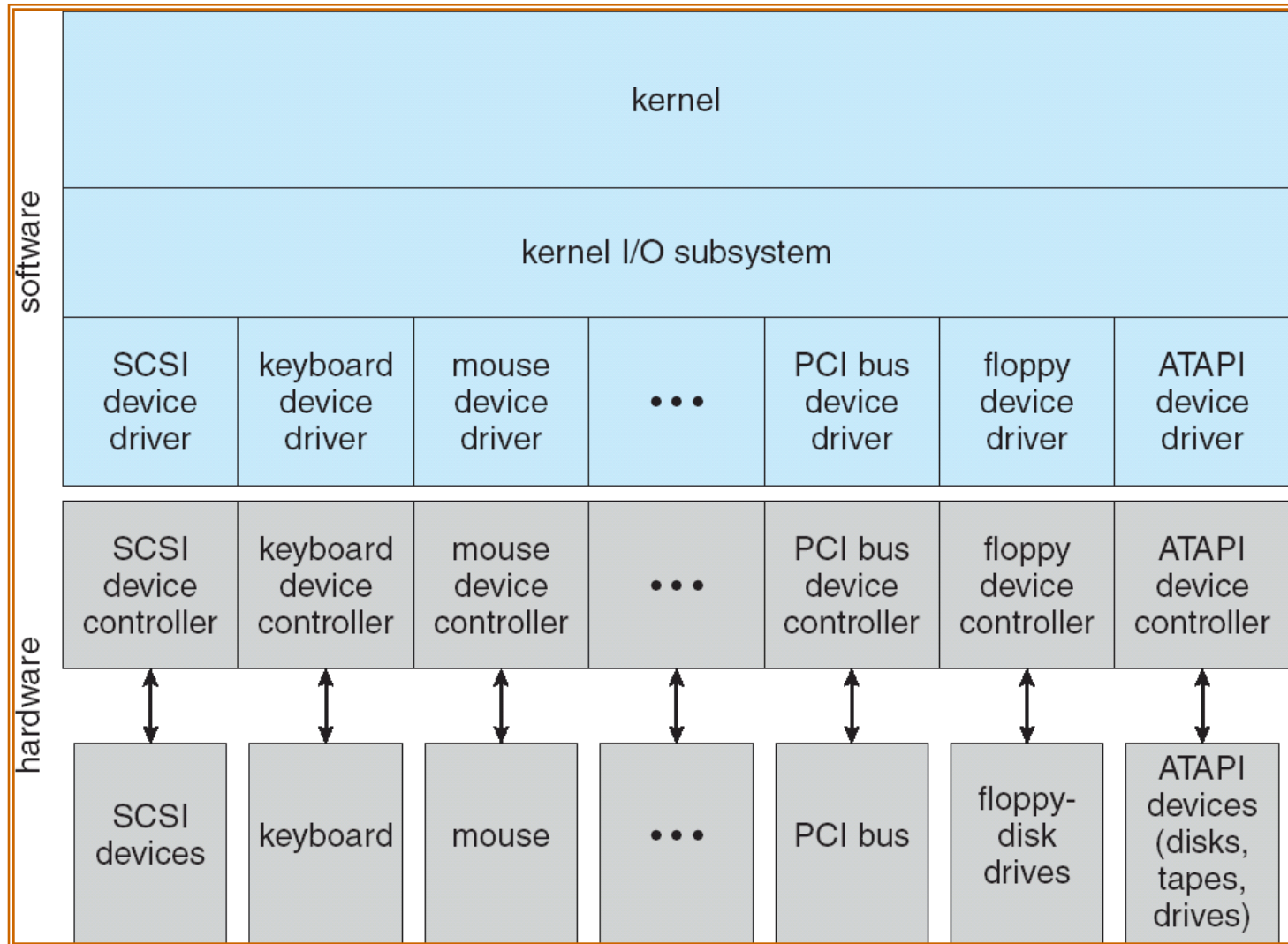# Six Step Process to Perform DMA Transfer

# Application I/O Interface

▸ I/O system calls encapsulate device behaviors in generic classes

▸ Device-driver layer hides differences among I/O controllers from kernel

▸ Devices vary in many dimensions
  ▸ **Character-stream** or **block**
  ▸ **Sequential or random-access**
  ▸ **Sharable or dedicated**
  ▸ **Speed of operation**
  ▸ **read-write, read only,** or **write only**

▸

# A Kernel I/O Structure

# Characteristics of I/O Devices

| aspect | variation | example |
|---|---|---|
| data-transfer mode | character<br>block | terminal<br>disk |
| access method | sequential<br>random | modem<br>CD-ROM |
| transfer schedule | synchronous<br>asynchronous | tape<br>keyboard |
| sharing | dedicated<br>sharable | tape<br>keyboard |
| device speed | latency<br>seek time<br>transfer rate<br>delay between operations | |
| I/O direction | read only<br>write only<br>read–write | CD-ROM<br>graphics controller<br>disk |

# Characteristics of I/O Devices (Cont.)

▸ Subtleties of devices handled by device drivers

▸ Broadly I/O devices can be grouped by the OS into

  ▸ Block I/O

  ▸ Character I/O (Stream)

  ▸ Memory-mapped file access

  ▸ Network sockets

▸ For direct manipulation of I/O device specific characteristics, usually an escape / back door

  ▸ Unix `ioctl()` call to send arbitrary bits to a device control register and data to device data register

▸ UNIX and Linux use tuple of "major" and "minor" device numbers to identify type and instance of devices (here major 8

```
brw-rw----  1 root disk 8, 0 Mar 16 09:18 /dev/sda
brw-rw----  1 root disk 8, 1 Mar 16 09:18 /dev/sda1
brw-rw----  1 root disk 8, 2 Mar 16 09:18 /dev/sda2
brw-rw----  1 root disk 8, 3 Mar 16 09:18 /dev/sda3
```

▸

# Block and Character Devices

- Block devices include disk drives
  - Commands include read, write, seek
  - **Raw I/O**, **direct I/O**, or file-system access
  - Memory-mapped file access possible
    - File mapped to virtual memory and clusters brought via demand paging
  - DMA
- Character devices include keyboards, mice, serial ports
  - Commands include `get()`, `put()`
  - Libraries layered on top allow line editing

# Network Devices

- Varying enough from block and character to have own interface

- Linux, Unix, Windows and many others include **socket** interface

    - Separates network protocol from network operation

    - Includes `select()` functionality

- Approaches vary widely (pipes, FIFOs, streams, queues, mailboxes)

# Clocks and Timers

▸ Provide current time, elapsed time, timer

▸ **Programmable interval timer** used for timings, periodic interrupts

▸ `ioctl` (on UNIX) covers odd aspects of I/O such as clocks and timers

# Blocking and Nonblocking I/O

▸ **Blocking** - process suspended until I/O completed
  ▸ Easy to use and understand
  ▸ Insufficient for some needs

▸ **Nonblocking** - I/O call returns as much as available
  ▸ User interface, data copy (buffered I/O)
  ▸ Implemented via multi-threading
  ▸ Returns quickly with count of bytes read or written

▸ **Asynchronous** - process runs while I/O executes
  ▸ Difficult to use
  ▸ I/O subsystem signals process when I/O completed

▸

# Two I/O Methods



(a)

(b)

# Kernel I/O Subsystem

▶ Scheduling
  ▶ Some I/O request ordering via per-device queue
  ▶ Some OSs try fairness

▶ Buffering - store data in memory while transferring between devices
  ▶ To cope with device speed mismatch
  ▶ To cope with device transfer size mismatch
  ▶ To maintain "copy semantics"

▶

# Device-status Table

# Common PC and Data-center I/O devices and Interface Speeds

# Kernel I/O Subsystem

- **Caching** - fast memory holding copy of data
  - Always just a copy
  - Key to performance

- **Spooling** - hold output for a device
  - If device can serve only one request at a time
  - i.e., Printing

- **Device reservation** - provides exclusive access to a device
  - System calls for allocation and deallocation
  - Watch out for deadlock

# Error Handling

▸ OS can recover from disk read, device unavailable, transient write failures

▸ Most return an error number or code when I/O request fails

▸ System error logs hold problem reports

# I/O Protection

▸ User process may accidentally or purposefully attempt to disrupt normal operation via illegal I/O instructions

  ▸ All I/O instructions defined to be privileged

  ▸ I/O must be performed via system calls

    ▸ Memory-mapped and I/O port memory locations must be protected too

# Use of a System Call to Perform I/O

# Kernel Data Structures
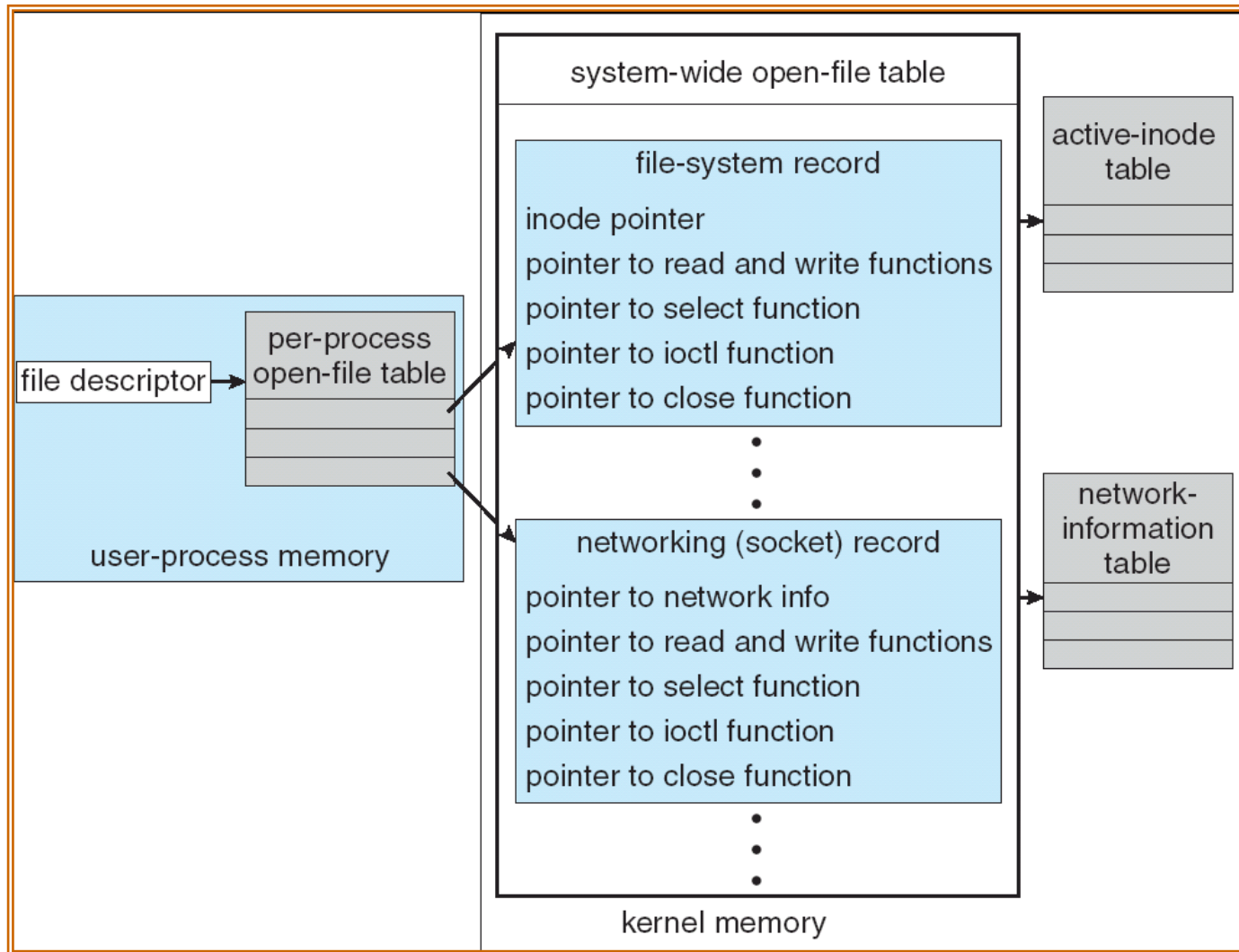
▸ Kernel keeps state info for I/O components, including open file tables, network connections, character device state

▸ Many, many complex data structures to track buffers, memory allocation, "dirty" blocks

▸ Some use object-oriented methods and message passing to implement I/O
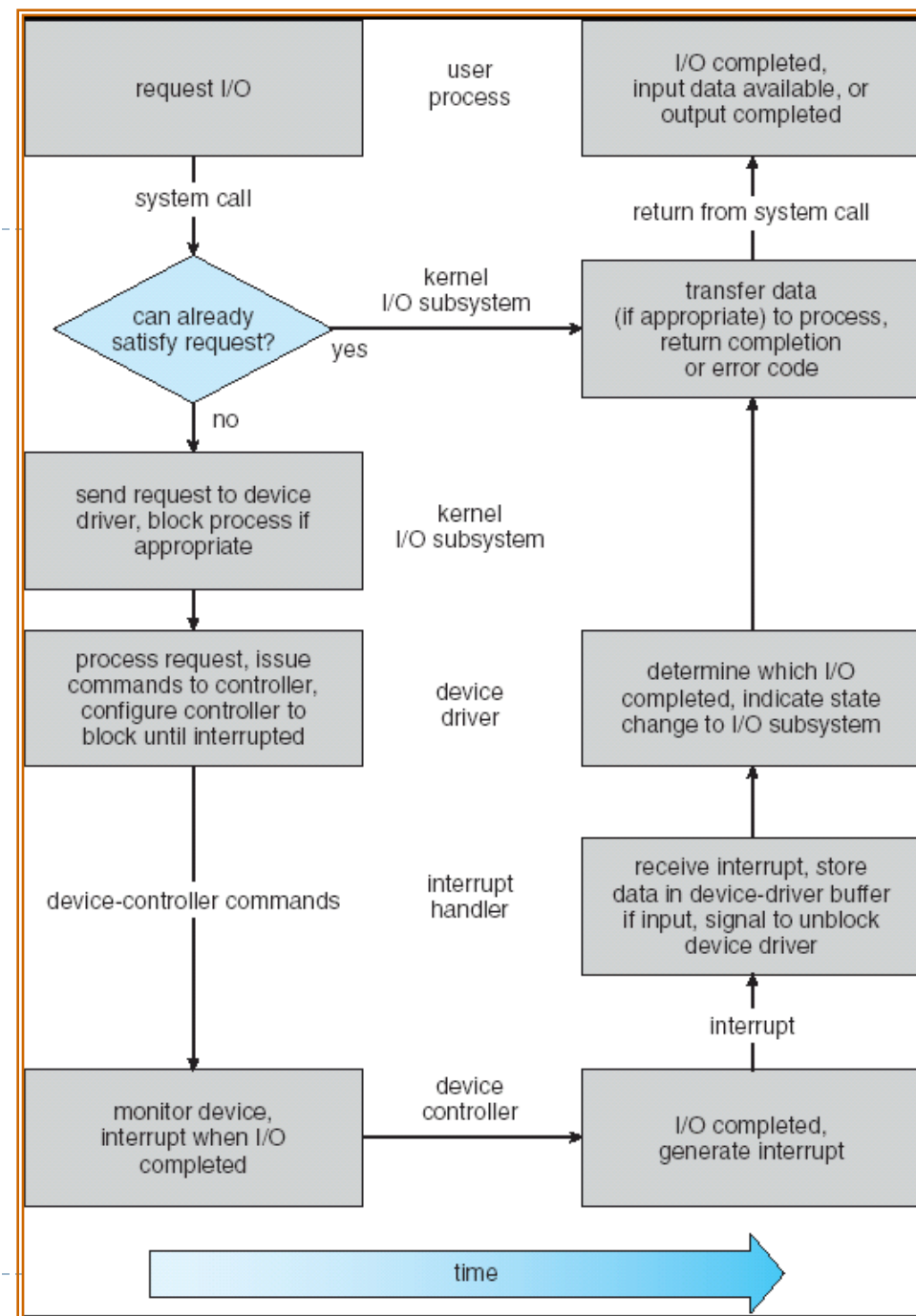
# UNIX I/O Kernel Structure

# I/O Requests to Hardware Operations

‣ Consider reading a file from disk for a process:

  ‣ Determine device holding file

  ‣ Translate name to device representation

  ‣ Physically read data from disk into buffer

  ‣ Make data available to requesting process

  ‣ Return control to process

# Life Cycle of An I/O Request
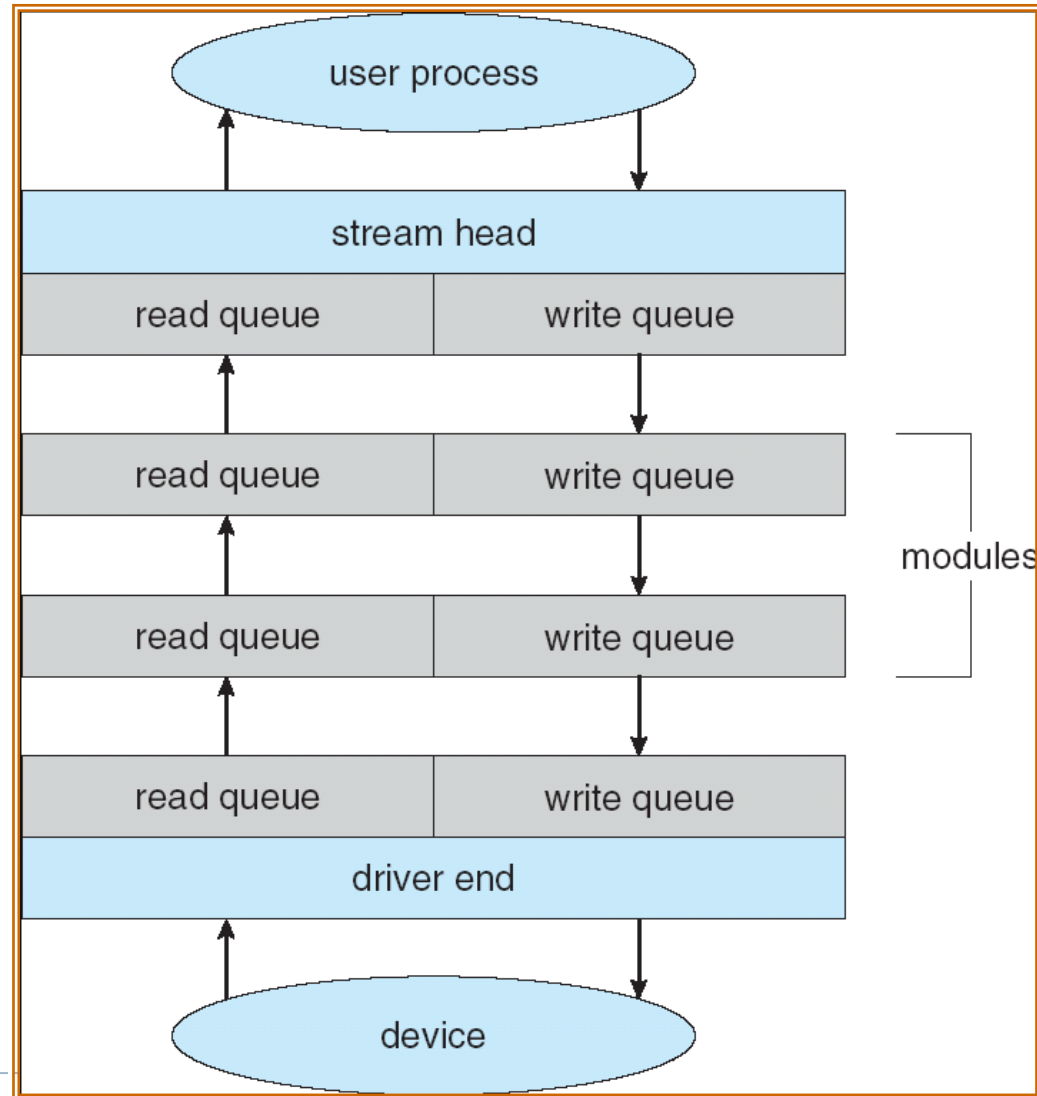
# STREAMS

▸ **STREAM** – a full-duplex communication channel between a user-level process and a device in Unix System V and beyond

▸ A STREAM consists of:

 - STREAM head interfaces with the user process

 - driver end interfaces with the device

 - zero or more STREAM modules between them.

▸ Each module contains a **read  queue** and a **write queue**

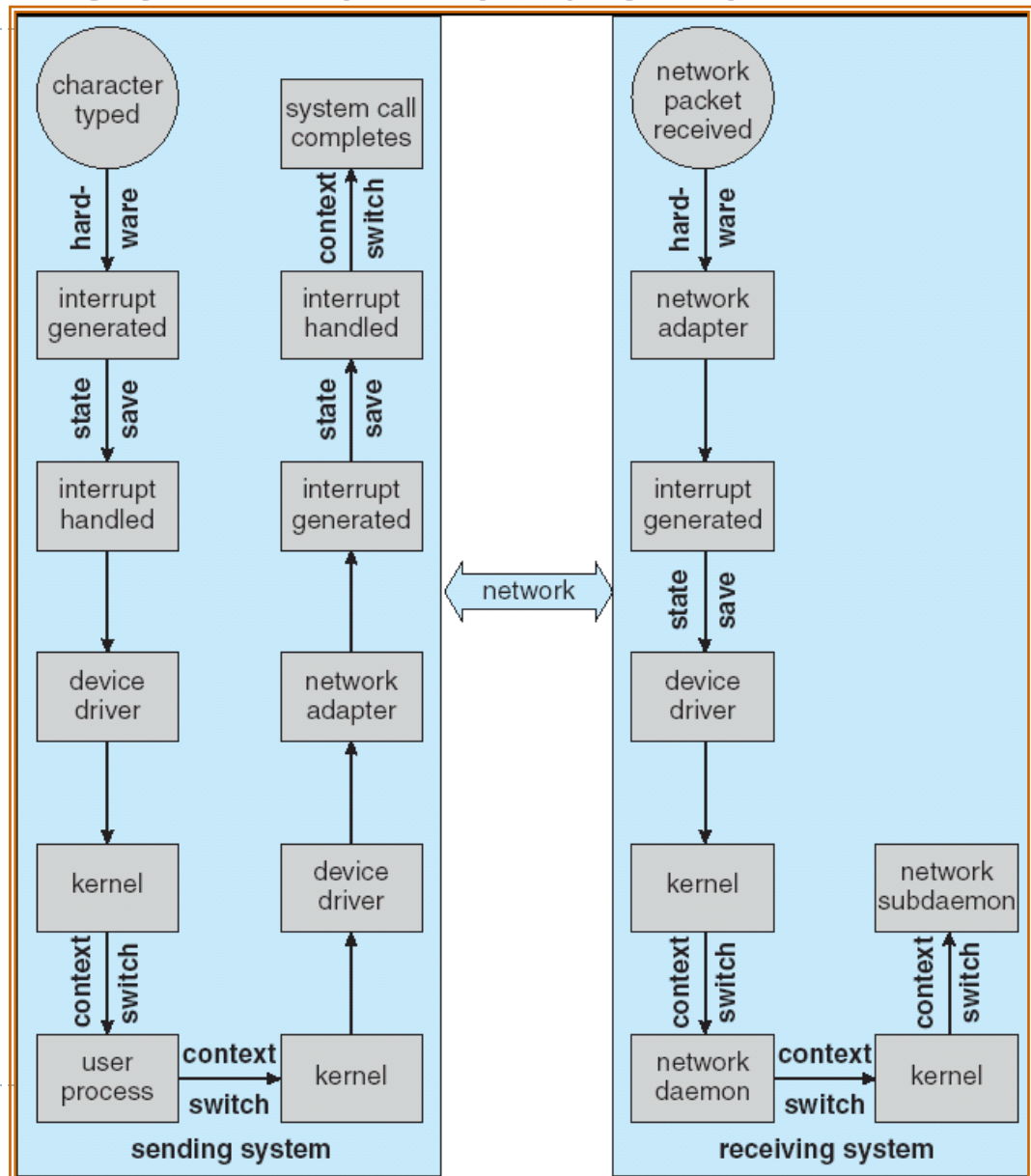▸ Message passing is used to communicate between queues

# The STREAMS Structure

# Performance

▸ I/O a major factor in system performance:

  ▸ Demands CPU to execute device driver, kernel I/O code

  ▸ Context switches due to interrupts

  ▸ Data copying

  ▸ Network traffic especially stressful
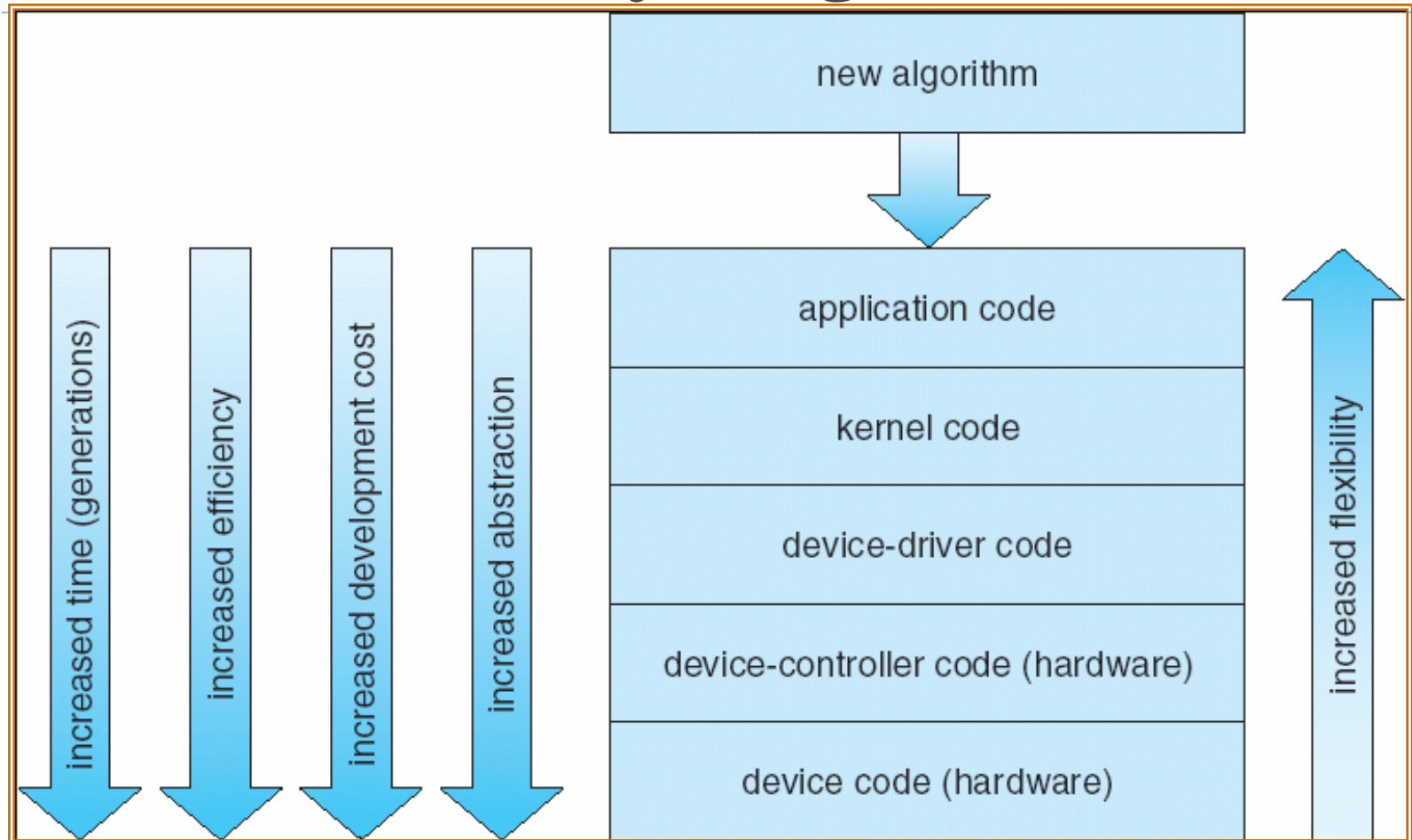
# Intercomputer Communications

# Improving Performance

‣ Reduce number of context switches

‣ Reduce data copying

‣ Reduce interrupts by using large transfers, smart controllers, polling

‣ Use DMA

‣ Balance CPU, memory, bus, and I/O performance for highest throughput

# Device-Functionality Progression

# I/O Performance of Storage (and Network Latency)