

Memory Mapped I/O

Certain memory addresses correspond to registers in I/O devices and not normal memory.

Control Register: Indicates if it is okay to read/write data register
Data Register: Contains I/O data

| Register | Location | Contains |
|---------------------|------------|--|
| Receiver Control | 0xffff0000 | Lowest two bits: Interrupt Enable Bit, Ready Bit |
| Receiver Data | 0xffff0004 | Received data stored at lowest byte |
| Transmitter Control | 0xffff0008 | Lowest two bits: Interrupt Enable Bit, Ready Bit |
| Transmitter Data | 0xffff000c | Transmitted data stored at lowest byte |

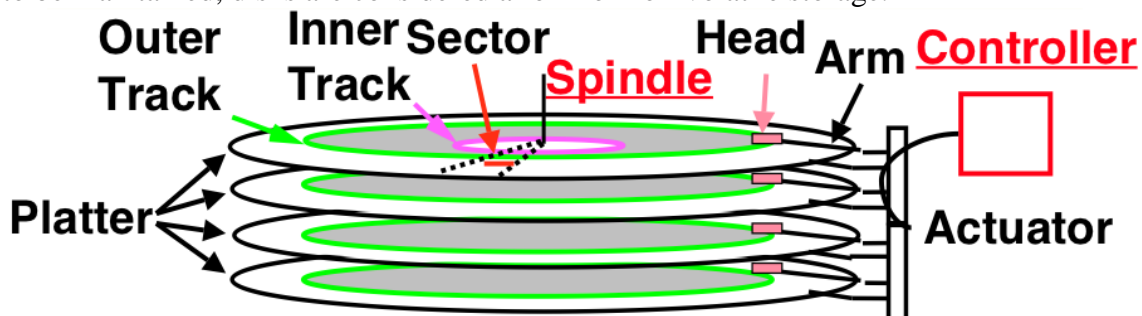
Describe MIPS code to read a byte from the receiver and immediately send it to the transmitter.

Polling and Interrupts

| Operation | Definition | Pro/Con | Good For |
|-------------------|--|---------|----------|
| Polling | Periodically check to see if the device is ready to transfer data. | | |
| Interrupts | Device makes an asynchronous request for data transfer. | | |

Disk Organization

Magnetic disks are one of the most common types of I/O devices. Bits are encoded by the controlling the polarity of magnetic fields on some sort of substrate. Since the magnetic fields do not require power to be maintained, disks are considered a form of non-volatile storage.



Based on notes by Aaron Staley, which were in turn based on notes by David Jacobs.

An additional term not shown here is that the collection of corresponding tracks across all the platters is called a cylinder.

There are two ways to address disks. Logical addressing treats the disk drive as one big array of blocks. Physical addressing uses a (cylinder, sector, platter) tuple to specify a blocks physical position in the disk drive.

Disk Performance

Disk Latency = Seek Time + Rotation Time + Transfer Time + Controller Overhead

RAID

Big disks are expensive (and dangerous). We can use an array of smaller disks to simulate the behavior of one larger disk with a more reasonable cost.

| | |
|---------|--|
| RAID 0 | No redundancy, just multiple disks |
| RAID 1 | Mirroring for redundancy, doubles read bandwidth |
| RAID 2 | Bit-level striping, increases bandwidth further |
| RAID 3 | Parity Disks, allows recovery from a single disk failure |
| RAID 4 | Block-level striping with Parity Disk, increases bandwidth |
| RAID 5+ | Striped Parity, reduces wear and tear |

Disk Exercises

We have a 7200 RPM drive with 3 ms seek time and a 20 MB/sec transfer rate once the head is in place. Assume that the controller overhead is negligible.

- 1) What is the overall throughput reading 1 MeBi of contiguous data on a random track?
- 2) What is the overall throughput reading 1MeBi of data spread randomly across the drive as 8 KiBi files? (this is why disk fragmentation is bad)

Performance Metrics

In order to get any meaningful definition of performance, we need to develop a quantitative metric that we all can agree on. This is harder than it sounds. We briefly talked about these when discussing pipelining.

Response Time, Execution Time, Latency - the time it takes to complete one task

Throughput, Bandwidth - tasks completed per time unit

Megahertz Myth

A processor's performance is determined by more than just the clock speed.

CPU time = Instruction Count * CPI * clock period

Based on notes by Aaron Staley, which were in turn based on notes by David Jacobs.

Exercises

You are the lead developer of a new video game at AE, Inc. The graphics are quite sexy, but the frame rates (performance) are horrible. Doubly unfortunately, you have to show it off at a shareholder meeting tomorrow. What do you do?

You need to render your latest and greatest über-133t animation. If your rendering software contains the following mix of instructions, which processor is the best choice?

| Operation | Frequency |
|-----------|-----------|
| ALU | 30% |
| Load | 30% |
| Store | 20% |
| Branch | 20% |

| A's CPI | B's CPI | C's CPI |
|---------|---------|---------|
| 1 | 1 | 1 |
| 3 | 5 | 3 |
| 2 | 3 | 4 |
| 3 | 2 | 2 |

What if the processors had different clock speeds? Assume A is a 1 Ghz processor, B is a 1.5 Ghz processor, and C is a 750 Mhz processor.

But wait, these processors are made by different manufacturers, and use different instruction sets. So the renderer (for the different architectures) takes a different number of instructions on each. Which is best if your main loop on A averages 1000 instructions; on B it averages 800 instructions; and on C it averages 1200 instructions?

Parallel Computing

Parallel computing refers more to multicore and multiprocessor machines. This is sometimes also called "supercomputing." Since the processors are physically closer together, there is a potential for much faster communications between them. However, synchronizing the processors can prove a difficult problem.

Amdahl's Law

The potential speedup from parallelization is limited by the amount a program can be parallelized. Let s be the fraction of the work that must be done sequentially and P be the number of processors. Then,

$$\text{Speedup}(P) \leq 1/s$$

Exercise

What are the contributing factors to Amdahl's law? Why isn't it an equality?