

Solutions are due at the beginning of class on the due date and must be **typed** and neatly organized. Late homeworks will not be accepted. You are permitted to discuss these problems with your classmates; however, your work and answers must be your own. All questions should be directed to the teaching assistant.

### Problem 1 : Block mapping.

The Linux ext2 file system inode structure uses 12 direct blocks as well as one single indirect, one double indirect, and one triple indirect block. Assume blocks are 4 KB in size, and block pointers are unsigned 32-bit integer values. What is the maximum file size accessible:

- (a) Using only the direct blocks?
- (b) Using the direct and single indirect blocks?
- (c) Using the direct, single, and double indirect blocks?
- (d) Using the direct, single, double, and triple indirect blocks?
- (e) Repeat (d), assuming blocks are only 1 KB in size. (Block size is chosen when a partition is formatted.)
- (f) Imagine adding a “quadruple indirect” block. What is the maximum file size (for 4 KB blocks) using the direct, single, double, triple, and quadruple indirect blocks?
- (g) Why are quadruple-indirect blocks irrelevant given the inode structure, block size, and block pointer size assumed in this problem?

An interesting artifact of the original ISO/ANSI C specification of the `fseek()` and `ftell()` system calls is that compliant operating systems were unable to address anywhere near the file size of (d). This is because these functions use the *long* type (a signed 32 bit integer) to specify the file position indicator (in bytes) for a stream.

- (h) Because of the use of the *long* type, what is the maximum file size that can be handled by the ISO/ANSI C `fseek()` and `ftell()` system calls (and, therefore, the maximum file size that can be handled by compliant operating systems)?

### Problem 2 : Extents.

Another technique to map data blocks from the inode is to use extent lists. An extent list is similar to the existing block pointers (see previous problem) except each entry is a vector:

- Pointer to the first block of the extent;
- Length of the extent (number of blocks following the first block).

Imagine replacing the 12 direct blocks in the ext2 inode with 6 extent descriptors, each 64 bits long (32 bits for the block pointer and 32 bits for the extent length). Assume a 512-byte block size.

- (a) What is the largest file addressible using these 6 extent descriptors?

- (b) What is the smallest file possible using all 6 descriptors?
- (c) Are indirect blocks still necessary when extent lists are used? Why or why not?

### **Problem 3 : Polling, interrupts, and DMA.**

Consider a system with the following components: a 500 MHz CPU, a hard drive that is capable of streaming data at 20 MB/s. Assume an infinite buffer in the operating system and that the I/O bus is capable of sustaining the drive's maximum streaming bandwidth.

The drive can use three different methods for notification and data transfer. These methods are:

- Polling and Programmed I/O: each polling operation consumes 500 cycles. Data is transferred in 8 word chunks (4 bytes/word). Polling must be done often enough that no data is lost.
  - Interrupts and Programmed I/O: each interrupt consumes 500 cycles. Again no transfer can be missed and data is transferred in 8 word chunks.
  - Interrupts and DMA: It takes 1500 cycles to initiate a DMA transfer and another 500 cycles to process the interrupt upon DMA completion. Assume the DMA engine transfers data in 4 KB chunks.
- (a) For each of the above scenerios calculate the fraction of the CPU consumed given the assumption that the disk is transferring 100% of the time.
  - (b) Assume now that the disk is only active 10% of the time. Again calculate the fraction of CPU consumed for each scenerio. Hint: think carefully about how often the CPU must poll in scenerio 1.
  - (c) For each of the above scenerios give the relative advantages and disadvantages of that scheme over the others.

### **Problem 4 : Those that do not learn from the mistakes of history are doomed to repeat them.**

Back in the days when Bill Gates was only moderately rich and “high-end” hard disk capacities were perhaps tens of megabytes, the INT 13 interface (BIOS software interrupt 0x13) was used by operating systems to send commands to hard disks. The INT 13 interface uses 24 bits to specify 512-byte sector locations based on physical disk geometry:

- 10 bits for cylinder number (up to 1024 cylinders),
  - 8 bits for head number (up to 256 heads/cylinder),
  - 6 bits for sector number (up to 63 sectors/track: sectors are numbered 1–63, omitting zero).
- (a) What is the largest disk capacity (in sectors and GB) addressable through the INT 13 interface? Note: use  $1 \text{ GB} = 2^{30} \text{ B}$ .

In a similar manner, the ATA/ATAPI-5 (what used to be IDE) interface defines 28 bits for addressing sectors based on the physical disk geometry:

- 16 bits for cylinder number (up to 65,536 cylinders),

- 4 bits for head number (up to 16 heads/cylinder),
- 8 bits for sector number (up to 255 sectors/track; zero is not used).

Notice that these bit fields vary significantly in size from the INT 13 specification. This was caused by poor communication and cooperation among the standards bodies, and resulted (around 1994) in the first of many embarrassing limitations on disk addressing. When combined, each of the cylinder/head/sector addresses are truncated:

Standard	Bits for cylinder number	Bits for head number	Bits for sector number	Total bits for address
INT 13	<b>10</b>	8	<b>6</b>	24
ATA	16	<b>4</b>	8	28
Combined	<b>10</b>	<b>4</b>	<b>6</b>	20

- (b) What is the largest disk capacity addressable through the INT 13 interface using ATA hard disks? (Recall that sector zero is unused in both interfaces.)

A short-term fix for these problems is logical block addressing (LBA) on the ATA bus combined with the “INT 13 extensions” described below. With LBA addressing, the BIOS no longer sends a physical cylinder/head/sector address. Instead, all available address bits are combined and used to send a logical address, with the logical-to-physical conversion done by the disk firmware.

- (c) What is the largest disk capacity addressable using LBA over ATA?
- (d) Maxtor Corporation released the DiamondMax 80 in August 2000. This ATA disk contains 160,086,528 sectors (76.335 GB). Assuming disk capacity growth continues present trends (doubling every 12 months), in what month of what year will ATA disks no longer be addressable using LBA over ATA?

Now isn’t that interesting? In 1998 the ATA standards committee (NCITS Technical Committee T13) began evaluating options for increasing the LBA address to a 48- or 64-bit value. 48-bit LBA extensions were officially added to the developing ATA-6 standard in October 2000; however, the ATA-6 standard has not yet been standardized.

- (e) What is the largest disk capacity addressable using 48-bit LBA over ATA-6?
- (f) Assuming disk capacity growth continues present trends as in (d), in what year will 48-bit LBA no longer be sufficient?

To overcome the BIOS limitation of (a), the INT 13 interface was changed to use the “INT 13 extensions.” Instead of the OS passing a cylinder/head/sector address through the INT 13 call, the OS passes a physical memory address pointing to a 64-bit logical block address.

- (g) What is the largest disk capacity addressable using the INT 13 extensions?
- (h) Assuming disk capacity growth continues present trends as in (d), in what year will the INT 13 extensions no longer be sufficient?

## Problem 5 : SCSI busses.

### Part I: Asynchronous SCSI

Note that the  $\text{ACK}/$  and  $\text{REQ}/$  signals are active low (this means that when they are held low they are active). Assume that the time between  $\text{REQ}/$  asserted and  $\text{ACK}/$  asserted is 0 and that data is transferred in 1 byte chunks.

- (a) Given the above timing diagram (Figure 1) for a SCSI asynchronous data transmit calculate the maximum possible bus bandwidth (MB/s).

Hint: it may be useful to extend the diagram to incorporate two data phases.

### Part II: Synchronous SCSI

SCSI synchronous data transfers are a little different from their asynchronous counterparts. For a data out phase the target sends a number of  $\text{REQ}$  pulses at a fixed frequency determined by the synchronous transfer period. The maximum number of  $\text{REQ}$  pulses without receiving an  $\text{ACK}$  is called the  $\text{REQ}/\text{ACK}$  offset. The target must then wait until it receives an  $\text{ACK}$  before it can send more  $\text{REQ}$ s. Finally, the  $\text{ACK}$  pulses come along with the data from the initiator at the same defined frequency. With the arrival of the  $\text{ACK}$  pulse the number of outstanding pulses has dropped below the offset and the target responds by sending data continually at the defined frequency.

- (b) What are the tradeoffs between using a large and small  $\text{REQ}/\text{ACK}$  offset?

### Part III: Asynchronous vs. synchronous SCSI

This question will investigate whether asynchronous or synchronous SCSI transfers are faster and why. To do this we must look at the propagation delays of the SCSI cable, the turn-around times in the controller chip, and the nature of the protocol.

The following information is based on SCSI-1:

- Cable propagation delays have been measured to be about 1.7 ns/foot.
  - Typical turn-around times are about 40 ns. (The turn-around time is the amount of time the SCSI chip takes to change an output in response to an input.)
  - The bus is 8 bits wide.
- (c) Remember that asynchronous SCSI mode uses a handshake before each data byte can be sent. First  $\text{REQ}$  goes true, then the  $\text{ACK}$  is permitted to go true, then  $\text{REQ}$  is permitted to go false, and finally  $\text{ACK}$  is permitted to go false (each of these signal must travel the full length of the cable and perform a turn-around before the next signal can proceed). For the following length cables calculate the maximum data rate:
- 1 foot
  - 6 meters (19.6 feet)
  - 25 meters (82 feet)
- (d) In synchronous mode the sender is permitted to transmit data without waiting for the  $\text{ACK}$  to come back and  $\text{REQ}$ s are transmitted at a fixed transfer rate. As a result only one propagation delay and turnaround are necessary (in the best case). For SCSI-1 the frequency at which data can be sent is 5 MHz. What is the maximum data rate? What is the maximum cable length given this data rate?

### **Problem 6 : I/O system design.**

You have been given a number of individual components with which you have been asked to build a large storage system. However, there are some tradeoffs between the price and performance of the components and your company wants to achieve the greatest performance possible at the lowest cost.

Given the following components with their performance information:

- 3 GHz CPU
- 16 byte wide memory with a 50 ns cycle time
- 33 MHz/32-bit PCI bus with up to 8 SCSI controllers
- SCSI controller adds 0.2 ms overhead per I/O, 1 bus/controller, cost \$350
- SCSI bus 7 disks/bus, bus @ 320 MB/s
- OS with 10,000 CPU cycles per disk I/O
- large disk: 146 GB
- small disk: 36 GB
- both disks are 10,000 RPM, 6 ms avg seek, stream at 25 MB/s, and cost \$4/GB
- avg I/O size of 8 KB

Make the following assumptions when answering the questions below:

- Every disk I/O requires an average seek and average rotational delay (of half a rotation).
  - All devices are used at 100% capacity and that the workload is evenly distributed between all the disks.
- (a) Give the breakdown of the max I/O operations per second (IOPS) for each component (CPU, memory, I/O bus, a SCSI controller, a disk).
- (b) In a fully configured system (a system with the max number of disks controllers) which components are under-utilized? Which components will be a bottleneck?
- (c) For a 2 TB database what is the optimal configuration of the system? Find highest performance first and then lowest cost for that performance level (in terms of disks and controllers).
- (d) Given the configuration in Part (c), you are told you can upgrade one type of component. Given the following choices which is the best component to upgrade:
- new disks (same sizes) but 15,000 RPM, 4 ms avg seek, stream at 30 MB/s
  - new PCI I/O backplane 66 MHz/64-bit
  - new 1000 MIPS processor

## Target Asynchronous Send

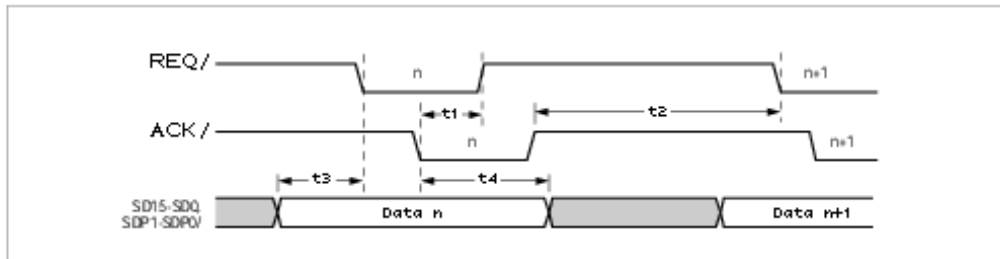


Figure 6-40: Target Asynchronous Send Waveforms

Table 6-40: Target Asynchronous Send Timings

Parameter	Symbol	Min	Max	Units
REQ/ deasserted from ACK/ asserted	$t_1$	10	-	ns
REQ/ asserted from ACK/ deasserted	$t_2$	10	-	ns
Data setup to REQ/ asserted	$t_3$	55	-	ns
Data hold from ACK/ asserted	$t_4$	20	-	ns

Figure 1: SCSI Asynchronous Data Out Timing