

Computer Systems Notes Zoe Farnum

1 Program Structure and Execution

1.1 Information Storage

Computers store information as a series of bits. These bits can be interpreted by users in either source binary, comfortable decimal, or compatible hexadecimal.

Computers also have a default word size, i.e. the largest continuous block of memory the computer can access. Currently, most computers are 32 bit, however more are becoming 64.

Going with word size, each data type also has a typical size in memory:

C Declaration	32 Bit	64 Bit
char	1	1
short int	2	2
int	4	4
long int	4	8
long long int	8	8
char*	4	4
float	4	4
double	8	8

Table 1: Size of C Data Types

Besides the bits themselves, the order also matters, which brings up the distinction between little-endian and big-endian.¹ Big Endian has the highest place values put in the lowest memory location, while Little Endian has the lowest place values put in the lowest memory location.

1.2 Integer Arithmetic

Depending on the type of numbers involved in addition, we can get strange or unexpected behavior:

If we're dealing with large numbers, ones that are near to the word size in length, we have to be concerned about overflow. Overflow occurs when the full integer result cannot fit within the word data type.

Unsigned integer addition results in a something that resembles modular addition. If we have signed integers, we need to now concern ourselves with the negative numbers as well. Negative overflow often results in a positive number due to the definition of two's complement.²

¹Name stolen from Gallium's Thaw.
²Computers represent negative numbers by essentially inverting the bits. Normal binary adds each place, this system subtracts each place from the highest set bit. For a 4 bit word size, the number 0111 would actually be -8 instead of 7.

CSCI 2400 1 Hsu

Computer Systems Notes Zoe Farnum

There are a number of different commands that assembly uses in order to provide all the functionality available in C code. These are summarized in the table below. We have to delve slightly deeper in the jump command however. The jump command depends on operational flags, and allows for sophisticated code.

2.1 Procedures

A procedure call involves taking data from one part of the program to another. This is managed by the program stack.

The program stack has two registers assigned to it, the stack pointer `%esp` and the frame pointer `%ebp`.

By convention, certain registers are kept private. Usually the registers `%eax`, `%edx`, and `%ecx` are classified as callee-save registers, that is the registers that are written to by the function caller. On the other hand, the registers `%edx`, `%esi`, and `%edi` are classified as callee-save registers. That is the registers needed by the called function.

2.2 Arrays

Arrays are a common tool, and are fairly simplistic in nature. They work by allocating the required memory for the specified data type, and referring to it in order.

Structures are treated (to an extent) the exact same as arrays. When a structure is created, the required memory for the data types fills contiguous space in memory.

Figure 1: IA32 Stack Structure

CSCI 2400 2 Hsu

Computer Systems Notes

3 Processor Architecture

3.1 Y86 Instruction Set

In this instruction set some commands are split up into several instructions.

3.2 Sequential Implementations vs. Pipelined

In a sequential implementation, all cycles over one data operation can start until the old one has finished.

Figure 2: Sequential Implementation

Conversely in a pipelined implementation, we split up the begin operations before old ones have finished.

4 Optimization

Optimization is an art.

We can use the metric *cycles per element* (CPE) to compare programs.

First step in code optimization is to reduce the number of instructions referred to as code motion.

After that's complete, we remove unnecessary memory accesses.

Followed by loop unrolling, the strategy that involves concrete number of steps and turning them into a set of instructions that can be unrolled in certain cases.

CSCI 2400

Computer Systems Notes Zoe Farnum

halt	stop the program
nop	no operation
movl	register ← register
movl	immediate ← register
movl	register ← memory
movl	memory ← register
opl	integer operation
jnl	jump
call	call function
retl	return
pushl	push onto stack
popl	pop from stack

Table 3: Y86 Instruction Set

Figure 3: Sequential Implementation

5 Memory Hierarchy

5.1 RAM

RAM comes in two forms, static and dynamic. Static RAM is much more expensive, but way faster.

5.1.1 SRAM

Each bit is stored in a bistable memory cell. It can only be in one position of another, never both, or neither. It will also keep its state indefinitely as long as it's kept powered.

5.1.2 DRAM

Each bit is stored as a charge on a capacitor. They lose power relatively quickly, and are much cheaper.

5.2 Disk Storage

Disks have several platters that spin at a fixed rate. The capacity of a disk is determined by:

$$\text{Capacity} = \frac{\text{bits}}{\text{sector}} \times \frac{\text{sectors}}{\text{track}} \times \frac{\text{tracks}}{\text{surface}} \times \frac{\text{platters}}{\text{platter}} \times \text{disk}$$

There is an actuator arm responsible for reading and writing from and to the disk.

Recording Density \times Track Density \times Area Density

or

$$\text{Capacity} = \frac{\text{bits}}{\text{sector}} \times \frac{\text{sectors}}{\text{track}} \times \frac{\text{tracks}}{\text{surface}} \times \frac{\text{platters}}{\text{platter}} \times \text{disk}$$

CSCI 2400 4 Hsu

Computer Systems Notes Zoe Farnum

5.4.4 Fully Associative Caches

A fully associative cache has $E = C/B$. Set selection is trivial as there is only one set, however indexing is tricky given that there are many similar tags. These caches are not effective at a large scale.

Figure 4: Sequential Implementation

5.4.2 Direct Mapped Caches

A cache with one line per set is a direct-mapped cache. When we have a cache miss, three steps occur:

Step One: Set Selection. The correct set is located using s .

Step Two: Line Matching. Determine if any line already contains the data. If one does, we have a cache hit.

Step Three: Line Replacement. If we have a cache miss, we need to replace a line with the new data.

5.4.3 Set Associative Caches

Cache's miss arise when we have direct mapped caches, where there is one line per set. We can relax this constraint a little, and establish each set to have at least 2 lines. A set with $1 < E < C/B$ is called an E -way set associative cache. This is a much more sophisticated caching method, and we need to search each line separately instead of each set line before. The replacement policies differ from machine to machine, and it can be tricky determining what replacement method to utilize.

CSCI 2400 5 Hsu

Computer Systems Notes

6 Linking

7 Virtual Memory

All Assembly code is executed using virtual memory. At the actual hardware address.

The MMU is in charge of converting these physical addresses and vice-versa. This can lead to problems however with the MMU.

7.1 Page Tables

To help solve this, most machines use fixed length page virtual page number as equal to the virtual address divided.

Page tables also map main memory to these fixed size page table records a virtual page \rightarrow hardware page mapping.

A page table is an array of page table entries (PTE) that point to physical pages.

The virtual address space $\{0, 1, 2, \dots, N-1\}$, with N physical pages are $\{0, 1, 2, \dots, M-1\}$.

Physical pages are $P = 2^p$ bytes in size and have been partitioned. Underscored have not been created by the VM system in memory. Underscored are not.

There is a valid bit on each row that indicates whether or not in memory.

CSCI 2400

Computer Systems Notes Zoe Farnum

Figure 5: TLB, page table, and cache for real memory systems. All values in the TLB, page table, and cache are in hexadecimal notation.

CSCI 2400 7 Hsu