Seminar 3

Petru Eles petel@ida.liu.se tf: 281396

Some few slides are based on material from the course book as well as from the book "Computer Systems: A Programmer's Perspective" by Bryant & O'Hallaron

# **Binary Representation**

- Information is represented as a sequence of binary digits: Bits
- What the actual bits represent depends on the context:
  - □ Numerical value (integer, floating point, fixed point)
  - □ Sequence of characters (text)
  - **Executable instruction**

# **Binary Representation**

- Information is represented as a sequence of binary digits: Bits
- What the actual bits represent depends on the context:
  - Numerical value (integer, floating point, fixed point)
  - □ Sequence of characters (text)
  - Executable instruction
- Depending on the context, operations performed are:
  - **Logical computation (context: logic)** 
    - 1: true Operations: And, Or, Exclusive-Or (Xor), Not
    - 0: false
  - Numerical Computation (context: numbers)
    - 1 <u>Operations</u>: Addition, Subtraction, Multiplication, 0 Division

# **Logical Computation: Boolean Algebra**



# **Logical Computation: Boolean Algebra**



It applies similarly to bit vectors (operations apply bitwise):

11100101	11100101	11100101	
& 01101101	01101101	<b>^ 01101101</b>	~ 01101101
01100101	11101101	10001000	10010010

#### **Arithmetical Computation**

#### Adding two one bit numbers



#### How is this Done in Computers?

- Logic values are represented by voltage levels:
  - □ High Voltage (e.g. 3.3V):1
  - □ Low Voltage (e.g. 0V): 0

At the output of a circuit we can have the following signal; this circuit produces the sequence 0, 1, 0 (or false, true, false):



# The Basic Building Block: The Transistor



H: high voltage level (1, true) L: low voltage level (0, false)

# The Basic Building Block: The Transistor



Observe!

This implements logic Not from V<sub>in</sub> to V<sub>out</sub>!



Such a circuit is called a *Not gate* (also *inverter*):



H: high voltage level (1, true) L: low voltage level (0, false)

# **Gates for Boolean Operations**

- Gates are electronic devices that perform Boolean operations.
  - **Gates are built as small electronic circuits based on transistors;**
  - Gates are the basic building blocks out of which VLSI (very Large Scale Integration) circuits are built; today computers are implemented as VLSI circuits, with up to billions of transistors on a chip.



# **Gates for Boolean Operations**

- Gates are electronic devices that perform Boolean operations.
  - **Gates are built as small electronic circuits based on transistors;**
  - Gates are the basic building blocks out of which VLSI (very Large Scale Integration) circuits are built; today computers are implemented as VLSI circuits, with up to billions of transistors on a chip.

Any logical function can be implemented as a combination of such gates.

• The one bit adder:



Sum = A Xor B C<sub>out</sub> = A And B

This is just a truth table capturing a logical function; thus, it can be implemented with a combination of logical gates!

• The one bit adder:

Α	B	Sum	Cout
0	0	0	0
0	1	1	0
1	0	1	0
1	1	0	1

Sum = A Xor B C<sub>out</sub> = A And B

This is just a truth table capturing a logical function; thus, it can be implemented with a combination of logical gates! Here is the circuit:



This is called a *Half Adder* (does not consider input carry).

The Full Adder (adds two bits and input carry):



The Full Adder (adds two bits and input carry):



• A four bits adder: adds two four bits numbers and an input carry:



 By further cascading full adders, one can build 8, 16, 32, 64, ... bit adders.

 In a similar way, circuits for other arithmetic operations can be implemented.

#### How to Store a BIT?

 Circuits like those shown in the previous slides are called "combinatorial": they produce an output that only depends on the input; the output is maintained as long as that input is applied.

What about a circuit that is able to store a bit? You can write 1 or 0 to the circuit, and the output will keep the value also after the input has disappeared.

# **Flip-Flops**

- The circuit below has two inputs. One (called S) for setting it to 1 (H), the other (called R) for setting to 0 (L).
  - When there arrives an input 1 to S, the output becomes 1; it will stay 1, until there comes an input 1 to R.
  - Once an input 1 arrived to R, the output switches 0, and stays so until an 1 arrives to S.















# **Flip-Flops**

One flip-flop can store one bit. Using groups of several flip-flops, arbitrary long sequences of bits can be stored. This is a basic technique to store data in computers e.g. in registers.

# Let's Go Over to Computers

• We have seen how data (logical and numerical) is represented in a computer.

We have seen that it is possible to construct circuits that are able to operate on data and perform logical and arithmetical operations.

• We have seen that circuits can be built which are able to store data.

## Let's Go Over to Computers

• We have seen how data (logical and numerical) is represented in a computer.

We have seen that it is possible to construct circuits that are able to operate on data and perform logical and arithmetical operations.

• We have seen that circuits can be built which are able to store data.

#### Now, let's see how a computer is built and works!

# What is a Computer/Computer-System?

A computer is a data processing machine which is operated <u>automatically</u> under the control of a list of <u>instructions</u> (called a program) stored in its main memory.

 Computers today are extremely complex and are built of many interconnected components; in addition to actual data processing, they have to perform tasks, such as communicate with other computers and devices, to interact with the user and the environment, etc. Therefore, we speak about *Computer Systems*.













#### **Input/Output Devices**

They connect the computer to the external world. Connection is via controllers/adaptors.



# **How Does a Computer Work?**

- All computers in use, simple or complicated, big or small, cheap or expensive work according to the same basic concept, known as the <u>von Neumann</u> <u>architecture</u>:
  - Data and instructions are both stored in the main memory (stored program concept);
  - The content of the memory is addressable by location (without regard to what is stored in that location);
  - Instructions are executed sequentially (from one instruction to the next, in order of their location in memory) unless the order is explicitly modified.
#### CPU



#### The basic organization (architecture):

- Central processing unit (CPU) contains:
  - Control unit (CU) that coordinates the execution of instructions;
  - Arithmetic/logic unit (ALU) that performs arithmetic and logic operations;
  - □ A set of registers.
- Main memory.



### **Register Organization**

- The set of registers within the CPU represents the top level of the memory hierarchy inside the computer system:
  - <u>User visible registers</u>: can be accessed by programs, for data storing.
  - <u>Control and Status registers</u>: used by the Control Unit to control the operation of the CPU; not directly accessible by the programmer.

#### CPU



#### **User Visible Registers**

A set of registers which can be used without restrictions as operands for any operation and as address registers; these are so called general-purpose registers.

#### CPU



#### **Control and Status Registers**

- Program Counter (PC): holds the address of the instruction to be fetched and executed.
- Instruction Register (IR): holds the last instruction fetched.
- Program Status Word (PSW): Condition
  Code Flags + other bits defining the status of the CPU.

#### CPU



#### Arithmetic Logic Unit (ALU)

Performs arithmetic and logic operations.
 There might be several of them in a CPU.
 ALUs are different, depending on the data type they operate on: integer ALU, floating point ALU, etc.

#### CPU



### Arithmetic Logic Unit (ALU)

Performs arithmetic and logic operations.
 There might be several of them in a CPU.
 ALUs are different, depending on the data type they operate on: integer ALU, floating point ALU, etc.

### **Control Unit**

- The control unit generates the appropriate signals such that all other components of the CPU and the computer system, together, execute the current instruction.
- The current instruction to execute is stored in the instruction register (IR); it is the instruction whose memory address is stored in the program counter (PC)

- A CPU can only execute *machine instructions*,
- Each computer has a set of specific machine instructions which its CPU is able to recognize and execute.

- A CPU can only execute *machine instructions*,
- Each computer has a set of specific machine instructions which its CPU is able to recognize and execute.

A machine instruction is represented as a sequence of bits (binary digits). These bits are organized into *fields* that define:

- A CPU can only execute *machine instructions*,
- Each computer has a set of specific machine instructions which its CPU is able to recognize and execute.



- A machine instruction is represented as a sequence of bits (binary digits). These bits are organized into *fields* that define:
  - □ What has to be done (the operation code).

- A CPU can only execute *machine instructions*,
- Each computer has a set of specific machine instructions which its CPU is able to recognize and execute.



- A machine instruction is represented as a sequence of bits (binary digits). These bits are organized into *fields* that define:
  - □ What has to be done (the operation code).
  - **To whom the operation applies (source operands).**

- A CPU can only execute *machine instructions*,
- Each computer has a set of specific machine instructions which its CPU is able to recognize and execute.



- A machine instruction is represented as a sequence of bits (binary digits). These bits are organized into *fields* that define:
  - □ What has to be done (the operation code).
  - **To whom the operation applies (source operands).**
  - Where does the result go (destination operand); in this example CPU it is assumed that the result of the operation is stored in the same place where the second operand was stored; no additional field is needed.

- A CPU can only execute *machine instructions*,
- Each computer has a set of specific machine instructions which its CPU is able to recognize and execute.

The number of bits, number and length of the fields and their order is particular to each computer; this defines the *instruction format* of that computer.

# **Types of Machine Instructions**

- Machine instructions are of four types:
  - **Data transfer between memory and CPU registers**
  - □ Arithmetic and logic operations
  - Program control (test and branch); these are those instructions that change the flow of instruction execution by *jumping* to an instruction *different* from the instruction following the current one in memory.
  - □ I/O transfer

You see, there are very simple things a machine instruction does! But many machine instructions, together, perform the big thing!

### **Instruction Execution**

#### Let's imagine you write in a program the following instruction:

Z := (Y + X) \* 3;

The instruction will be executed by the CPU as a sequence of four machine instructions!

### **Instruction Execution**

Let's imagine you write in a program the following instruction: Z := (Y + X) \* 3;

Men whic tion	nory address at ch the instruc- /data is stored	Content of the memory	
Move value of Y to Reg 2	00001000	0000101110001010 Move addr of Y Reg 2	
Add value of X to Reg 2 (result kept in Reg 2)	00001001	0001101110000010 Add addr of X Reg 2	Instr
Multiply Reg 2 with 3 (result kept in Reg 2)	00001010	0010100000011010 Mul value "3" Reg 2	inctions
Store Reg 2 at address of Z	00001011	0001001110010010 Move addr of Z Reg 2	
Value of X: 11 Value of Y: 3 Final value of Z: 42	01110000 01110001 01110010	00000000000001011 ← X 0000000000000011 ← Y 0000000000101010 ← Z	ר על גענו 10 אוג גענו 10 אוג



















# **Compilers**

**High Level Language** Z := (Y + X) \* 3; (e.g. C, C++, Java)

We have written in our program:

What the computer executes is:



**Machine instructions** for the particular processor that runs the program.

# **Compilers**

High Level Language (e.g. C, C++, Java)

We have written in our program:

What the computer executes is:



Machine instructions for the particular processor that runs the program.

Who brings us from our program to the machine instructions?



A compiler is a program that translates programs written in a high level language into machine code to be executed on a certain processor.

# **The Machine Cycle**

From the previous example you have seen that many things have to be done to execute a simple machine instruction;

- Fetch instruction
- Decode instruction
- Execute instruction

# **The Machine Cycle**

From the previous example you have seen that many things have to be done to execute a simple machine instruction;

- Fetch instruction
- **Decode** instruction
- Execute instruction
  Execute instruction
  Execute instruction

# The Machine Cycle

From the previous example you have seen that many things have to be done to execute a simple machine instruction;

- Fetch instruction
- **Decode instruction**
- **Execute instruction**

Fetch operand(s)Execute instruction

#### Machine Cycle

Each instruction is performed as a sequence of steps; the steps corresponding to the execution of one instruction are referred together as a *machine cycle*.

The number and nature of steps in the machine cycle differ from processor to processor.



# The Quest for Speed

Running faster (more instructions per time unit) has been a permanent goal of computer designers.

Two main factors contribute to high performance of modern processors:

- 1. Fast circuit technology: smaller and faster switching transistors, allowing the processor to run at higher frequency.
- 2. Architectural features such as:
  - **Smart memory hierarchies**
  - **D** Pipelining

Superscalar architectures are executed in parallel. 

# **Memory System**

One of the most crucial aspects in designing efficient computer architectures is the memory system.

#### What do we need?

We need memory to fit very large programs and to work at a speed comparable to that of the microprocessors.

- Main problem:
  - **Processors are working at a high clock rate and they need large memories;**
  - Memories are much slower than microprocessors; but for executing a single instruction you need several memory accesses (fetch the instruction and operands); it doesn't help that the processor is fast, if the memory is orders of magnitude slower.

# **The CPU-Memory Gap**



- Fast memories are more expensive per byte and cannot be very large (main memory is much smaller than SSD or Disk)
- It is possible to build memory structures that are as fast as the CPU, but they are very expensive and small.

- Fast memories are more expensive per byte and cannot be very large (main memory is much smaller than SSD or Disk)
- It is possible to build memory structures that are as fast as the CPU, but they are very expensive and small.



#### The good news:

It is possible to build a composite memory system which combines *small, fast memories* (from the top of the hierarchy) and *large slow memories* (from the middle and bottom of the hierarchy) and which <u>behaves (most of the time) like a large fast memory</u>.

### How can this work?

#### The good news:

It is possible to build a composite memory system which combines *small, fast memories* (from the top of the hierarchy) and *large slow memories* (from the middle and bottom of the hierarchy) and which <u>behaves (most of the time) like a large fast memory</u>.

### How can this work?

### The answer is: Locality
## **The Principle of Locality**

- During execution of a program, memory references by the processor, for both instructions and data, tend to cluster: once an area of the program is entered, there are repeated references to a small set of instructions (loop, subroutine) and data (components of a data structure, local variables or parameters on the stack).
  - <u>Temporal locality</u> (locality in time): If an item is referenced, it will tend to be referenced again soon.
  - Spacial locality (locality in space): If an item is referenced, items whose addresses are close by will tend to be referenced soon.

### **Cache Memory**

A cache memory is a small, very fast memory that retains copies of recently used information (instructions and data). It operates transparently to the programmer, automatically deciding which values to keep and which to overwrite.

- Due to the property of locality, most of the time, the instruction or data required by the CPU will be available in the top cache. If not, it will be loaded from the lower level cache; once loaded the information will be written into the top level cache and replace some existing one, in order to make space for the new information.
- Which information is replaced when new one has to be written?
  - Some information is overwritten that has, for a long time, not been used by the CPU (and, thus, is less likely to be needed in the future)
- The above procedure is repeated at each level of the hierarchy.

### **The Quest for Speed**

Running faster (more instructions per time unit) has been a permanent goal of computer designers.

Two main factors contribute to high performance of modern processors:

- 3. Fast circuit technology: smaller and faster switching transistors, allowing the processor to run at higher frequency.
- 4. Architectural features such as:
  - Smart memory hierarchies
  - Pipelining
  - Superscalar architectures

Several instructions are executed in parallel.

# **Pipelining**

#### We remember the machine cycle



# **Pipelining**



# Pipelining



#### **Superscalar Architectures**

- You can imagine a superscalar processor as composed of several pipelines running together.
  - As opposed to simple pipelined computers, superscalars fetch several instructions and produce several results simultaneously



### **The Quest for Speed**

Running faster (more instructions per time unit) has been a permanent goal of computer designers.

Two main factors contribute to high performance of modern processors:

- 5. Fast circuit technology: smaller and faster switching transistors, allowing the processor to run at higher frequency.
  - 6. Architectural features such as:
    - Smart memory hierarchies
    - **D** Pipelining
    - Superscalar architectures

That one has been a primary source of performance improvement all over the years.
Processors running at higher and higher frequencies allowed for a continuous increase in speed. That doesn't work any more!!!

#### **The Power Wall**

We have reached the limit due to the temperature produced by the high power consumption! Further increase of the frequency is impossible!



This is the main challenge today!

New ways have to be explored in order to deliver performance!

### **Multicore Chips**

Multicore chips: Several processors on the same chip.

- This is the only way to increase chip performance without excessive increase in power consumption:
  - Instead of increasing processor frequency, use several processors and run them in parallel, each at lower frequency.

### **Intel Core Duo**

#### Composed of two Pentium M superscalar processors.



# Intel Core i7

#### Composed of four x86 SMT (simultaneous multithreading) processors.



### **ARM11 MPcore**

Composed of four ARM11 processor cores.



# Intel's Single-Chip Cloud Computer (SCC)

#### Composed of 48 P54C Pentium cores



#### Where are We

