Computer Systems and Architecture Revision Notes

James Brown May 3, 2017

Contents

1	Intr	$\operatorname{roduction}$	1
2	Fun 2.1	damentals of Computer Organisation The ven Neumann Architecture and Everyting Programs	1
	2.1	The von Neumann Architecture and Executing Programs	1
		2.1.1 von Neumann Architecture	1
		2.1.2 The Clock Cycle	2
	2.2	2.1.3 Executing Programs	2
	2.2	Harvard Architecture	4
	2.3	Case Studies	4
		2.3.1 MIPS	4
		2.3.2 Intel x86	4
3	Inst	cruction Sets and Assembly Language	5
	3.1	Types of Instruction Set	5
	3.2	Types of MIPS Instructions	6
	3.3	MIPS Register Conventions	6
	3.4	Machine Code	6
	3.5	Further MIPS instructions	8
4	CP	U Microarchitecture	8
	4.1	MIPS Microarchitecture	8
		4.1.1 Instruction Fetch	8
		4.1.2 Instruction Execution	8
		4.1.3 Branches and Jumps	10
	4.2	Integrating the Datapath	11
	4.3	CPU Control	12
5	Dio	ital Logic	12
J		Digital Logic	12
	0.1	5.1.1 Transistors	12
		5.1.2 Decoders	13
		5.1.3 Multiplexers	14
		5.1.4 A Simple Adder	14
	- 0	5.1.5 A Simple ALU	15
	5.2	Number Representation	15
		5.2.1 Positive Integers	15
		5.2.2 Positive Real Numbers	16
		5.2.3 Negative Numbers and Two's Complement	16
		5.2.4 Floating Point Numbers	17
	5.3	Clocked Logic	18
6	I/O	and Peripherals	18
7	Imr	proving Performance	18
-	7.1	Caches and Virtual Memory	18
	7.2	Pipelining and Branch Prediction	18

1 Introduction

These are notes I have written in preparation of the 2017 Computer Systems and Architecture exam. This year the module was run by Iain Styles (I.B.Styles@cs.bham.ac.uk). This is the module did not cover networks and they are not examinable - as such I will not be writing about them here.

2 Fundamentals of Computer Organisation

Computer programs consist of **instructions** and **data** which are identical in appearance, but they are logically distinct. Programs have a ordered set of instructions which are executed sequentially, unless it's otherwise stated. Programs also have data which is there to be manipulated by the instructions which are run. In the computers memory these will both have the same physical representation, but are not the same as each other. Because of this, when storing instructions and data they must be kept logically separate. That is to say they must be stored in different regions of memory, and not interspersed with each other for example.

We may want to describe the computers architecture at a variety of levels of abstraction:

- Level 5: High Level Languages. These are largely independent of the physical machine, occasionally regarded as part of the architecture.
- Level 4: Assembly Language. Programming in terms of the machine's basic operations.
- Level 3: Operating System. Common services and management functions.
- Level 2: Instruction Set. The basic operations that the machine can execute.
- Level 1: Microarchitecture. The distinct functional units that are required to implement the instruction set, and their organisation.
- Level 0: Digital Logic. The implementation of the functional units in terms of basic logic operations.
- Level -1: Physical Device. The implementation of the logic using basic electronic components such as transistors, and the physical substrate on which these are constructed.

2.1 The von Neumann Architecture and Executing Programs

2.1.1 von Neumann Architecture

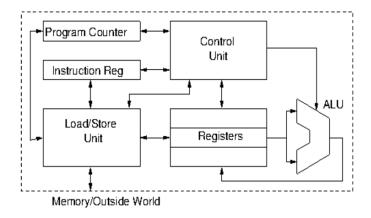


Figure 1: The von Neumann Architecture

Many modern computers are built on the (or use slightly modified) von Neumann architectures and can be considered it's heart.

We consider the main memory as being logically - but not necessarily physically - separate from the CPU. The main memory holds all of the instructions and data that make up a program(s). As stated earlier, instructions and data are stored in distinct locations within the main memory so that they are easily distinguished. Within main memory, instructions are stored sequentially so that you can determine the flow of the program implicitly from their order. Lastly, main memory is also a volatile storage method, meaning that all data is lost once the power is cut.

The **Load/Store unit** is used as the interface between the CPU and the outside world. It issues and receives requests to transfer instructions and data between the CPU and the memory via the bus.

The **registers** are small amounts of local, fast access storage the hold data that is currently in use. Data is passed to the registers by the load/store unit. Each register can hold one 'word' of data. The main registers are used purely to hold data - instructions are dealt with separately.

The **instruction register** hold the current instruction that is being executed so that it can be used by the control unit to configure the ALU. Only one instruction is active at any one time, unless the design features special techniques for performance improvements that rely on multiple instructions being executed simultaneously.

The **ALU** (arithmetic and logic unit) is the 'engine' of the computer. It performs all the computations and comparisons. It also reads data from registers and writes the results of calculations back into the registers.

The **program counter** is a special register that contains the memory location (address) of the next instruction which shall be executed - a bookmark in essence. In the normal execution of a program, the program counter is incremented after each instruction to point to the next memory location. Some instructions may change the value of the program counter in order to change the order of execution.

Additionally, we may want to add some other components just outside of the CPU which can be just as important. Due to the fact the main memory is remote from the CPU, access times can be slow. We may use an intermediate layer of memory known as **cache** that is smaller but much faster to access. This would be used to hold portions of programs that are likely to be used again shortly. We may also want to use stable, **long-term memories** such as disks or DVDs. Typically these all will require the use of an **Input-Output (IO) controller** which handles peripheral devices such as disk drives, mice and keyboards. It may do this through an extension of the memory addressing protocol or via an interrupt based protocol. Memory, peripherals and the CPU all communicate with each other via the **bus** which carries data around and allows, for example, data to be transferred from disk into main memory. This bus consists of a set of physical wires plus a protocol (there are many available, such as PCI, ISA, IDE, SCSI) that is implemented by the **bus controller**. The bus controller determines which subsystems can communicate. It should be noted that only one piece of data can be on the bus at any time.

2.1.2 The Clock Cycle

The vast, vast majority of computer systems are *synchronous* - meaning their activities are synchronised by an external clock signal in the form of a **square-wave electric pulse**. This speed is frequently quoted as a measure of CPU performance but to say it is analogous to performance would be slightly false. There are many other factors and different architectures cannot be compared on the basis of their clock rate alone. The time between two different pulses is related to the frequency of the CPU by t = 1/f and is called a cycle time.

For the most part computer systems must be synchronus as variability in manufacturing means that it's not possible to know exactly how long it will take for a particular operation to complete. The clock cycle is chosen to be slightly longer than the longest delay in the system which ensure the machine is in a well defined state when the next set of operations start which is triggered by the **rising edge** of the clock pulse.

2.1.3 Executing Programs

Computer programs, in their most basic form, are just sequential series' of instructions. In von Neumann architecture, the execution of these instructions is governed by the **instruction execution cycle**. The instruction execution cycle is triggered by the clock cycle, but has several stages

within it which are triggered by successive clock pulses. One complete instruction cycle usually takes several clock cycles to execute - exact numbers depend on the type of instruction and the details of the particular machine. Fetch data from memory for example may take several clock cycles to execute, and it may take several cycles before the data is safely loaded into a register. Others may complete in a single clock cycle such as the addition of the contents of two values stored in registers.

Most architectures follow the same basic set of stages in the **Fetch-Decode-Execute Cycle**. In can roughly be broken down into 8 steps in an idealised version:

Fetch:

- 1. Inspect the program counter to find the address of the next instruction
- 2. Load the next instruction from memory into the instruction register
- 3. Update the program counter to point at the next instruction

Decode:

- 4. Determine the type of instruction fetched
- 5. If the instruction requires data from memory, determine its address (usually embedded in the instruction

Execute:

- 6. Fetch and required data from memory into one of the CPU registers
- 7. Execute the instruction
- 8. Return to step 1 for the next instruction

Starting a program doesn't fit neatly into this simple model and before we enter the cycle we also need to take a few actions to make sure that things are ready. Firstly, we need to load the program from disk into main memory. The instructions and data needed by the program will each occupy a block of memory, which is allocated by the operating system, and the memory address of the first instruction is called the **entry point**. When first started, the entry point is loaded into the program counter which then becomes the starting point of the cycle.

Fetch. Once we have a valid instruction location in the program counter (PC), we can begin the cycle. An important note is that at this point in time all we have is the memory address - not the actual instruction - we still need to actual fetch it, hence the title. At the start of the next clock cycle, the CPU issues a request via the load/store unit to the memory by sending the memory address and a request to read from the memory via the bus. Later in time, the instruction will be received from memory by the load/store unit and then stored in the instruction register (IR). Depending on the relative speed of the clock cycle and the memory, it could take several cycles before the instruction is ready in the IR. Once the request has been made, the value of the PC is changed to point to the next instruction - which usually just involves simply incrementing the PC. This may however be modified by some instructions such as branch or jump.

Decode. Now we have the instruction in the IR, we can begin to act upon it in the CPU. The type of instruction is determined by the control unit. This is necessary in order to determine if any further actions needs to be taken in order to execute this instruction.

Execute. After finding the type of instruction, any data needed is fetched from the memory. For a lot of CPUs, most instructions can only actually access registers and there are dedicated instructions for accessing main memory. Once the data is in the registers, it can be operated upon. As mentioned earlier, some instructions change the flow of the program and are therefore allowed to change the PC as necessary (when doing this it is often required for the previous value of the PC to be stored so that execution can resume once the branch has completed).

2.2 Harvard Architecture

In the von Neumann architecture, instructions and data are accessed via the same physical and logical pathway (the load/store unit) and there is not formal separation between data and instructions at this level. In this case, the two types of information are stored in the same physical memory but are separated by their locations within the memory. This separation is common practice in all computers as it allows for dynamic repartitioning of the memory according to the needs of the program. The problem is due to the fact that by having a shared interface, instructions and data cannot be accessed simultaneously. This is known as the **von Neumann bottleneck** and it restricts CPU performance to the rate at which is can be supplied with data.

A potential solution is to provide separate memories for small amounts of instructions and data that are likely to be used soon - separate instruction and data caches. It is also common to provide separate interfaces to instruction and data memory which is known as the **Harvard architecture**. In a pure implementation of the Harvard model, instructions and data are stored in physically separate memory but this is not flexible enough for general purpose computational devices that has a single unified memory space. The **modified Harvard architecture** has a single unified memory space (that is partitioned for instructions and data) but with separate buses for instructions and data. In most modern machines, this is the approach that is taken.

2.3 Case Studies

2.3.1 MIPS

The MIPS processor is the canonical example of a modified Harvard architecture and is very similar to the von Neumann model shown earlier. It features an instruction register, a program counter, an ALU etc, but has separate pathways for accessing instructions and for accessing data. In the diagram they are shown as physically separate, but in reality they are part of a physically unified memory.

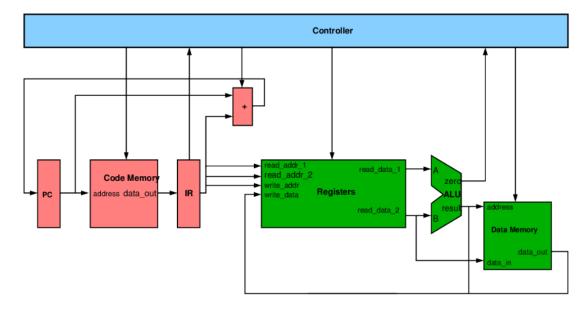


Figure 2: The MIPS processor

2.3.2 Intel x86

Modern versions of x86 are much more complex when compared to the relatively simple MIPS Harvard based architecture. This is partly down to the fact that MIPS has around 60 simple instructions whereas the latest Intel Core machines have hundreds of instructions - many of which are not simple. The benefit of these additional instructions for the x86 architecture is a significant

increase in performance and it can also make it easier for the programmer/compiler writer. This all comes at the cost of a much more complex design - MIPS R4000 contains 1.2 million transistors. A quad-core i7 processor on the other hand has 731 million transistors on a die that is not much larger than that in a MIPS machine. One of the main features of modern Core architecture is that it is highly superscalar and can execute multiple instructions simultaneously. Due to this, the CPU requires several ALUs, several instruction decoders, instruction queues and multiple levels of caching. If you can see through all the complexity, you can see that the Core 2 is essentially a modified Harvard architecture.

3 Instruction Sets and Assembly Language

One of the major defining features between different computers is the choice of **instruction set**. Modern computers are *Turing-complete* - meaning they can perform any computation that can be performed - so in some sense the set of instructions doesn't matter as long as they implement a Turing-complete system. In practice, the choice of instructions can greatly affect the programmers task - especially for very low level programming.

Each type of CPU has a different instruction set which are effectively incompatible with each other. Some machines have instructions that other do not and some may have identical instructions with different binary representations. Due to this, low-level code is extremely machine dependent and written by hand only when totally necessary. Due to the widespread use of von Neumann and Harvard architectures different instruction sets tend to have a general similarity.

3.1 Types of Instruction Set

Instruction sets are frequently classified as **complex** (Complex Instruction Set Computer - **CISC**, Intel designs for example) or **reduced** (Reduced Instruction Set Computer- **RISC**, for example ARM and MIPS).

CISC computers generally have a very extensive range of instructions (this tends to be in the area of several hundered). These can range from simple instructions like addition and subtraction to more complex operations that are often common combinations of simpler instructions to provide specific support for high-level functions. CISC instruction sets have the advantage of possibly making the translation of high-level software into machine language somewhat easier. Microcoding of the complex instructions can also provide a performance benefit over their implementation in software. One drawback is the added complexity of the hardware which may make debugging and optimisation very difficult. The biggest example of CISC processors is the Intel x86 family which even includes instructions such as AESDEC to perform AES decryption.

RISC machines are the opposite of CISC machines - the number of instructions is minimised and each instruction is highly optimised with the ability to make use of performance-enhancing measures such as pipelining and speculatice execution. A popular RISC machine is the MIPS architecture which was very popular in the 1990s but continues to be widely used in embedded systems because of its low power and heat generation. The basic instruction set has around 60 instructions (some MIPS variants did have more), and these are **superpipelined**. Key points on the MIPS architecture:

- 32-bit architecture (instructions, memory addresses and words of data are 32 bits long)
- 32 data registers, \$0 ... \$31. \$0 is a special register and reserved for the value zero. Other registers are reserved for other special purposes by convention only.
- Most instructions can only interact with registers there are special instructions for transferring data to/from memory. Due to this, some stages of the instruction cycle can be omitted.
- Byte-addressed, meaning that an increment of 1 in the program counter points to the next byte, not the next word. This makes the normal PC increment 4, not 1.

3.2 Types of MIPS Instructions

Instructions of the MIPS processor can be divided into eight rough categories:

Load/Store Instructions which fetch/store items from/to memory. Several variants which work on whole(32-bits) or part-words(half-words or bytes).

Arithmetic Instructions which are used to add/subtract etc. two variables being held in the registers, and also to perform comparisons.

Immediate Arithmetic Instructions which are similar to regular arithmetic instructions but used to specifically add/subtract a constant and a variable in a register.

Shift Instructions which are used to perform bit rotations - commonly used in cryptographic protocols

Multiply/Divide Instructions which perform multiplication or division on two variables being held in the registers.

Jump and Branch Instructions which are used to change the normal sequential flow of program instructions. For example, to call subroutines, take branches in the code (at conditionals) and implement loops.

Coprocessor Instructions which are used to send data and pass control to an external coprocessor which might, for example, be a graphics controller or external floating-point processor.

Special Instructions which do no fall into any of the above categories. In this course, we don't consider these at all.

To begin, we'll consider the most common instructions and how they relate to high-level code. Consider the simple code a = a + b;. In the MIPS instruction set, we have an instruction called add which takes three **operands** - analogous to arguments of functions in higher-level languages. The three operands are the **destination** of the result, and the two **sources** of the inputs. Arithmetic instructions can only access the registers so the operands must specify which registers are to be used. Therefore, we may translate our simple example into add \$8, \$8, \$9. This adds the contents of the register 8 to the contents of the register 9, and stores the result into register 8. The order of the operands is important here, the destination comes first!

This code assumes the values of the variables a and b are already in the registers. We will need instructions to load data from memory into the registers and to store the result from registers into memory. An instruction to load a word of data from the memory will need operands which specify which register the data should be loaded into and where in the memory it will come from. For the moment, we will denote memory addresses using C-like notation such as &a. The instruction lw \$8, &a loads the contents of a word of data at memory address &a and puts it in register \$8. It follows that we also have a function sw \$8, &a which takes the contents of a register and stores that word at the specified memory address. We can modify our program to load a and b from memory and to then store a in memory after performing the addition:

```
lw $8, &a;
lw $9, &b;
add $8, $8, $9;
sw $8, &a;
```

This alone does not make up a full MIPS program. We must add various assembler directives and also encapsulate this code into a main function in order to be able to run it.

3.3 MIPS Register Conventions

While all registers apart from \$0 is freely accessible by the programmer there is still a convention for their use. Their conventions are shown in Figure 3. Important registers are \$8 through \$15 which are for 'temporaries' and we can use them to store intermediate values. The assembler also supports the use of the register names rather than numbers, so those will be used from now on making to code a little more readable.

3.4 Machine Code

All MIPS instructions are 32 bits long, with the 32 bits divided into sections. The precise division of bits depends on the type of instruction. If we consider the simple arithmetic operations we

Name	Number	Use
\$zero	\$0	constant 0
\$at	\$1	assembler temporary
\$v0-\$v1	\$2-\$3	function return and expression evaluation
\$a0-\$a3	\$4-\$7	function arguments
\$t0-\$t7	\$8-\$15	temporaries
\$s0 - \$s7	\$16-\$23	saved temporaries
\$t8-\$t9	\$24-\$25	temporaries
\$k0-\$k1	\$26-\$27	reserved for OS kernel
\$gp	\$28	global pointer
\$sp	\$29	stack pointer
\$fp	\$30	frame pointer
\$ra	\$31	return address

Figure 3: The conventions of use for all registers in the MIPS architecture

have used in our examples, the machine level representation must include information about what operation is to be performed and which registers should be used. The precise format for this type of instruction (called a **register operation** is shown below.

0pcode	Source 1	Source 2	Dest	Shift	Func
6	5	5	5	5	6

Figure 4: MIPS instruction format for register operations

The **opcode** field of 6 bits describes which type of instruction is being described (not the specific instruction). This is important as it determines how the rest of the instruction is going to be interpreted. The remaining bits are used to encode the two **source** registers and the **destination** register as 5 bit numbers; a **shift** field which is used by certain bit-shifting operations and denotes how far the shift should be; and finally the **func** field which encodes exactly which instruction is desired. All register type instructions follow the same flow of information so having single opcode for all such instructions simplifies things significantly. We can then use the func code to configure the ALU for specific calculation. For add and **sub** instructions, the opcode is 000000 and the func codes are 100000 and 100010 respectively.

For **load/store** instructions, a slightly different format is required as the instructions require different information.

0pcode	Base Address	Src/Dest	Address Offset
6	5	5	16

Figure 5: MIPS instruction format for load/store operations

These instructions need to be interpreted differently to register operations and therefore have different opcodes - 35 for lw and 43 for sw. A memory address also has to be specified - which is done in two parts - a base and an offset. The 5 bits allocated to the base address do not specify an address themselves, they specify a register where the address is located. The offset field specifies an address relative to the base.

3.5 Further MIPS instructions

4 CPU Microarchitecture

The von Neumann model specifies a quite general and non-specific architecture for a computer, but absolutely none of the details. Here we are concerned with the **microarchitecture** which refers to the detailed structure and organisation of the machine. This can be divided into two broad parts:

- The Datapath is a collection of functional units which implement the instruction set. Each functional unit has a specific purpose: the registers are used for storing data, the program counter bookmarks the code; the instruction register stores the current instruction, and the ALU executes arithmetic and logic operations
- The Control Logic serves to configure the datapath in the right way so that it implements the desired instruction. It ensures that the correct data is going to the correct functional units, that the results are put in the right place and that the ALU is configured to perform the correct operation on the data. Control is the most complex part of the processor. It's somewhat simple in RISC machines due to the few operations they implement but much harder in CISC machines.

4.1 MIPS Microarchitecture

4.1.1 Instruction Fetch

The first step of the instruction cycle is to get the address of the next instruction and then to fetch that instruction. In order to do this, we need to send the contents of the program counter to memory and bring the contents of that address back to the CPU. We will need to use the **program counter**, the **main memory** and the **instruction register**. The main memory is usually not physically part of the datapath but is integral to its operation, so we will include it in our diagrams.

We can put these components together that implements the first stage quite simply. We send the contents of the program counter to the memory via the **address bus**. The returned instruction is then loaded into the instruction register via the **data bus**.

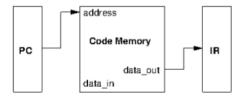


Figure 6: Datapath elements for instruction fetch

Next, we need to increment the program counter. Commonly this is just adding 4 to the PC to move to the next instruction. In the case of a branch or a jump we may need to make a further adjustment to the PC value, but the constant value of 4 is *always* added. As this is done every cycle we add a simple adder component into the datapath as shown in figure 7.

4.1.2 Instruction Execution

Whilst we should look at decoding the instruction after fetching it, it's more instructive to look at the datapath components required to actually execute the instruction set after the decoding has been done. Some components required are obvious (such as the ALU and registers), but others are less so.

Registers are required to provide a fast and local amount of memory for a small amount of data items. Any data that is used by an instruction in the MIPS instruction set must be present in the registers as they cannot access external sources (unless they are special load/store instructions). Because of this, we need to provide mechanisms for writing data into the registers and also for reading data from the registers. For some instructions we may need access up to three registers

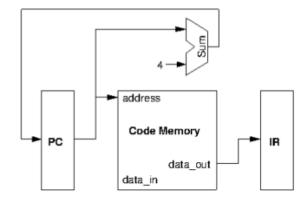


Figure 7: We have integrated an adder to increment the PC

simultaneously. As MIPS has 32 registers, we can use a 5-bit number to uniquely address a single register and these are encoded into the machine code for the instruction. The registers must therefore feature three 5-bit address ports - two of these are for reading and one is for writing - and two 32-bit outputs for retrieving data from the registers. We will also need one final 32-bit input to write data to registers. As not all instructions write to registers we also include an input which signals whether the register specified by the write_addr is to be written or not. All of these are shown in Figure 8.

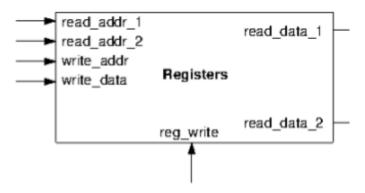


Figure 8: Registers and their input/output ports

Next, we would like to build the ALU. The ALU performs all of the processors arithmetic and logic operations and it's internal circuitry must be able to be configured so that it can perform any of the "ALU instructions". In MIPS, the ALU must also be able to compute address, for example, whena address offsets are used - e.g. lw \$0, 8(\$sp). The ALU handles the addition of the offset. The ALU also handles comparissons required by by branch and set instructions. The ALU must be able to accept up to two 32-bit words from the registers as input, perform the ALU instruction specified in the instruction register as determined by the ALU control unit and to output another 32-bit word (in the case of comparissons, just true/false using the zero output). The ALU is shown in Figure 9.

The registers and parts we have shown so far form the basis for a datapath that can execute R-type instructions (e.g. add). Fifteen bits are dedicated to selecting registers that are involved in ALU computations. We can pass these register addresses directly into the register address ports to select the correct registers. The remaining 17 bits are used by the control to configure the datapath to perform the instruction. This gives us the datapath shown in Figure 10 for these R-type instructions.

Load/store operations are somewhat more complex instructions. They need read and write access to both the registers and main memory, as well as to be able to use the ALU to calculate addresses. A source/destination register address is specified dependent on whether the operation is

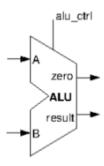


Figure 9: The ALU and it's input/output ports

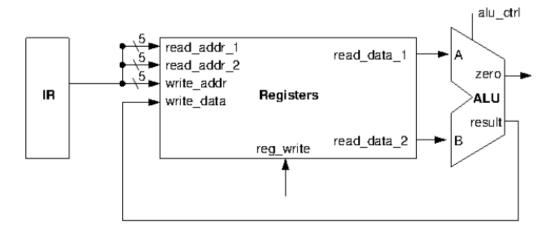


Figure 10: The datapath for R-type instructions

a load/store. The register of the **base address** is also specified and the 16-bit address **offset**. The addition of the base address and the offset is performed by the ALU in the MIPS architecture, but is not as straightforward as it might seem. The base address is 32-bits, whilst the offset is 16-bits. We need to add a new unit to 'extend' the offset so it can be added to the base. This seems trivial (just prepend some zeroes), but branches can go backwards and forwards in the code so we must allow for negative offsets using the two's complement representation. The manipulations that are required to translate the offset are performed by the **sign extension unit**. Once the address is computed by adding the offset to the base address, we must read or write the memory address that we have computed, as well as update the contents of the specified register is we are loading. A signal to the memory that controls whether is is a read or write that will be generated by the CPU control unit. We omit this for now to focus on the flow of data.

4.1.3 Branches and Jumps

Both branches and jumps achieve essentially the same result - they both have to change the contents of the program counter. The datapath mechanisms are essentially the same because of this, but the control differs depending on the specific instruction being executed as branches require a comparisson to be done. We will deal with branches first, and add jumps later.

In order to execute a branch, we have to look at the type of instructions, which is the same as I-type load/store instructions. In the beq instruction for example, we must compare the contents of the two registers specified in the instructions and if they are equal, we add the specified offset to the program counter. Otherwise we increment the PC by 4 as normal. The ALU is capable of performing comparrisons, and the registers are used to store the quantities to be compared. We will also have to extend the offset to 32-bits from 16-bits in order to add it to the PC. As we need the ALU for the comparison, we use a dedicated adder instead of the ALU for the bit extension.

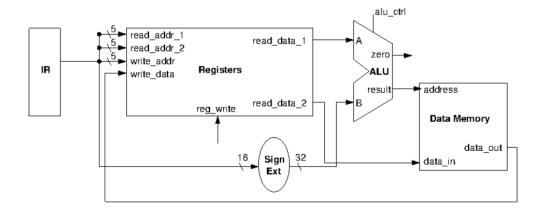


Figure 11: The datapath for load/store instructions

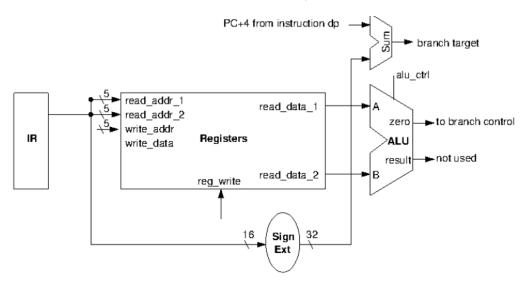


Figure 12: Datapath for branch instructions

This incurs a little overhead but reduces overall complexity. The resulting datapath is shown in Figure 12.

4.2 Integrating the Datapath

We've constructed datapath circuits for a range of the most common MIPS instructions, each of which has their own special requirements. Many of them also have a lot in common, and many of the functional units are used in more than on instruction type. Duplicating functional units is clearly undesirable and because of this we need to be able to control the flow of data within the datapath to ensure that the correct units are employed for each instruction. An example of why this is difficult is that data can be written into the registers from either the ALU or from the data memory, but only one set of wires which these pieces of data can be sent and we cannot simultaneously send two pieces of data down the same wires. The solution is to add switches to the wires, which is known as **multiplexing**.

A 2-1 multiplexer is a switch which allows one of the inputs to pass, based on a control signal C. With N control wires, we can choose 2^N inputs. For example, an 8-1 multiplexer will require 3 control wires. Multiplexing will allow us to combine and integrate different sections of the datapath that we designed for the different types of instruction. Firstly, we will integrate the ALU-based instructions which use the IR, the ALU and the registers; and the load-store instructions which use the IR, the registers, the ALU, data memory and a sign extension unit. The challenge is to integrate these two datapaths using only one ALU, and one set of registers. The two datapaths

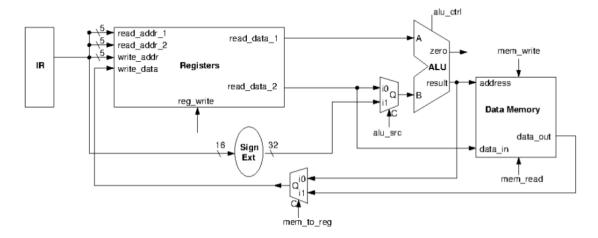


Figure 13: The integrated datapath for ALU and load/store instructions

we are trying to integrate can be seen in Figures 10 and 11.

In an ALU instruction, the ALU inputs both come from the registers, and the result of the operation is written back to the registers. In a load/store instruction, the ALU inputs are from the registers and sign-extension unit, whilst the registers need to be written by the data memory. Therefore we need to multiplex the data that is to be written into the registers; and the ALU inputs, as show in 13.

4.3 CPU Control

5 Digital Logic

5.1 Digital Logic

5.1.1 Transistors

The basic building block of the modern integrated circuit CPU is the **transistor**. For the purpose of this course, we only need to know the transistors are essentially switches which are controlled by their gate voltage. Depending on the type of transistor (NMOS or PMOS), they will either conduct or not conduct when a voltage is applied to the gate. An NMOS transistor will conduct when the gate is at a positive voltage , and a PMOS transistor will conduct when the gate has no voltage - otherwise it doesn't conduct.

Figure 14: (Left) Circuit symbol for an NMOS transistor. (Right) Circuit symbol for a PMOS transistor

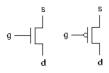
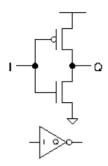


Figure 15: The logical operator **NOT** built from transistors



So far data and instructions have been represented by binary digits that have two possible states - zero or one. These states correspond to zero and positive voltages in physical circuits.

The zero voltage is usually called **gnd**(ground) and the positive voltage **vdd**. A signal that has a connection to gnd is in logic state 'zero', and a signal that has a connection to vdd is in logic state 'one'. Using transistors we can transform signals between 'zero' and 'one' states. To begin with, lets implement an inverter circuit which is identical to the logical operation **not** as shown in Figure 15.

When the input is at low voltage (zero), the lower NMOS transistor is non-conductive but the top PMOS transistor is conductive. The output therefore becomes connected through the PMOS transistor and is therefore logical 'one'. When the input is at high voltage (one), the PMOS transistor is non-conducting while the NMOS transistor is conducting. Q is therefore connected to gnd via the NMOS transistor and is logical 'zero'. We can see dependent on the input, the output of the circuit will be its inverse and it implements the logical **NOT** operator.

Similarly, we can implement other basic logical operators from a few transistors. While we might want to make **AND** and **OR** straight away, it's actually easier to implement their negated versions - **NAND** and **NOR**. We can do both of these in just four transistors, but it requires six to implement their non negated versions.

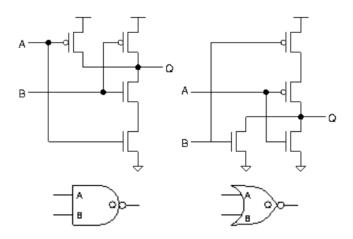


Figure 16: (Left) A **NAND** gate built from four transistors. (Right) A **NOR** gate built from four transistors

5.1.2 Decoders

Combinational logic is the general term for blocks of digital logic which contain no form of memory - the input must be determined solely by it inputs - and is made from networks of the simple logic gates we have just defined. Example of combinational logic we have seen whilst building the MIPS datapath include multiplexers, simple adders, and even an entire (non-pipelined) ALU.

Consider a simple block of combinational logic with two inputs and four outputs - a decoder. This is designed so that each of the four possible input combinations uniquely selects only one of the four outputs. This circuit is common in memories where it is used to select a unique memory location based on the presented address. We start by describing the truth table for this circuit as shown in the figure below. By inspecting the truth table, we can construct its function as a set of logic equations shown next to the truth table.

i0	i1	Q1	Q2	Q3	Q4
0	0		0	0	0
0	1	0	1	0	0
1	0	0	0	1	0
1	1	0	0	0	1

Figure 17: Truth table for a two input decoder

•
$$Q_1 = \neg A \wedge \neg B$$

•
$$Q_2 = \neg A \wedge B$$

•
$$Q_3 = A \wedge \neg B$$

•
$$Q_4 = A \wedge B$$

From here it becomes trivial to translate this into a circuit diagram - you simply have to read off the logical operations (here, the mixtures of **NOT** and **AND**) and insert the appropriate gates.

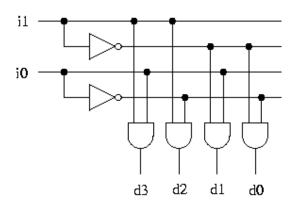


Figure 18: A 2-4 decoder built from logic gates

5.1.3 Multiplexers

We can also build **multiplexers** which are a very common component in the CPU. 2-1 Multiplexers are described by the logic equation $Q = (S \wedge I_1) \vee (\neg S \wedge I_0)$, and the truth table is shown below. Just as we did for a 2-4 decoder, we can easily build a 2-1 multiplexer from a few logic gates as described by its logic equation. In order to build 32-bit multiplexers, all we require are 32 single-bit multiplexers which share the same select signal 'S'.

i0	i1	\mathbf{S}	Q
0	0	0	0
0	0	1	0
0	1	0	0
0	1	1	1
1	0	0	1
1	0	1	0
1	1	0	1
1	1	1	1

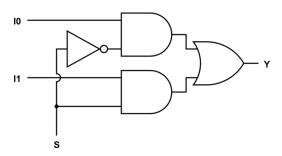


Figure 20: A 2-1 multiplexer built from logic gates

Figure 19: Truth table for a 2-1 multiplexer

5.1.4 A Simple Adder

An **adder** is used both in the ALU and in several other places of the datapath. We will design a circuit which is able to add unsigned integers, two's-complement integers and both unsigned and two's-complement fixed-point numbers. For single digit binary addition, we know that 0+0=0, 1+0=0+1=1 and 1+1=0 with a carry bit of 1. The carry bit of 1+1 means we will need to introduce a second output which will function as our carry bit. It follows that we can easily express the carry bit as $C=A \wedge B$. In order to do the sum itself we need to introduce a new logical operation, **exclusive-OR** (otherwise known as **XOR**). **XOR** is indentical to and **OR** gate but does not output 1 if both input are also 1. Using and **XOR** we can fully implement single-bit addition. The circuit which achieves this is commonly called a **half-adder**.

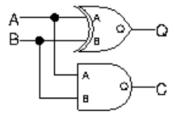


Figure 21: Logic gates which implemet a half-adder

This is clearly very limited as it can only add single bits together and we would like to add multiple-bit binary numbers. We can use a half-adder for the least-significant bit, but for higher bits we need to be able to carry in the the carry bit from the preceding bit. We can easily achieve the sum using a pair of half-adders where the first adds I_1 to I_2 , and the second adds the result of the previous to C_{in} . The carry bit of this addition is then the **XOR** of the carry bits from the two half-adders.

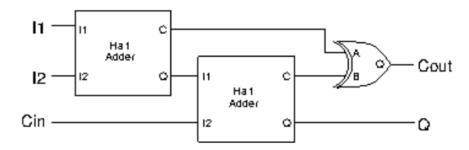


Figure 22: A full-adder built from two half adders and an XOR gate

We can now join lots of full adders (and one half-adder for the least-significant bit) together to create on single 32-bit adder. The carry bit from the most significant bit in the calculation will be used to indicate **overflow** (the result was too large to store in 32-bits).

5.1.5 A Simple ALU

Now we can build an adder, we can put together a simple ALU which can achieve addition of two numbers, bitwise-AND and bitwise-OR. To do this we take our 32-bit adder and modify it by adding one **AND** gate per bit and one **OR** bit per bit (with some multiplexers to select the operation we actually want to perform). We can then program the ALU by using the ALUOp signal to select the chosen function. How to build this ALU is shown in Figure 23. It is not too dificult to see how you could add additional functionality from here.

5.2 Number Representation

5.2.1 Positive Integers

The simplest of numbers we can represent in a binary computer system are positive (or unsigned) integers. This is the simple binary representation that we have already seen before and are common with. In principle, any natural positive number can be represented, but there are practical limitations. Modern computers break data into 32 or 64-bit chunks (referred to as words), which limits the size of numbers that we can store. With n bits, we can store integers from 0 to $2^n - 1$.

Long binary strings are tedious for humans to manipulate and mistakes are easily made. It's common to represent numbers in another base for this reason. Commonly, hexadecimal (base-16) is used. We need 16 symbols to represent values. We us the normal 10 decimal symbols, plus A-F. For example, the value of 0x1C3F in decimal is:

$$(1 \times 16^3) + (12 \times 16^2) + (3 \times 16^1) + (15 \times 16^0) = 4096 + 3072 + 48 + 15 = 7231$$

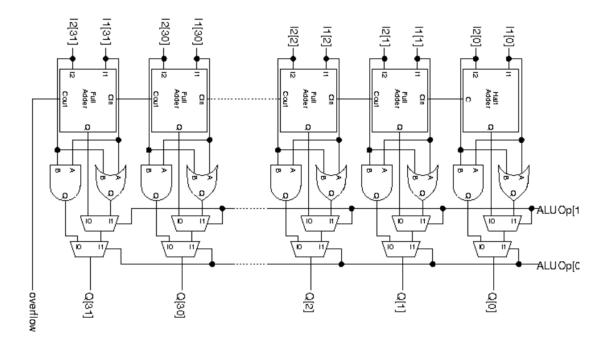


Figure 23: A simple ALU capable of performing addition, bitwise-AND and bitwise-OR operations

As 16 is an integer power of 2, converting between binary and hexadecimal is simply a matter of string substitution. Each group of four bits can be converted directly to their hex equivalent and vice-versa.

5.2.2 Positive Real Numbers

We can extend our representation of unsigned integers to unsigned reals (such as 2.1, 3.14 etc). We do so by allowing a decimal point in our binary number, and extending the powers backwards from 0.

$$101.001_2 = (1 \times 2^2) + (1 \times 2^0) + (1 \times 2^{-3})$$
$$= 4 + 1 + \frac{1}{8}$$
$$= 5\frac{1}{8}$$

We have a choice to make, where do we put the decimal point? Modern computers have 32 or 64-bit words, but many real numbers are still impossible to represent with this. Dependent on the position of the decimal point, we can have an increased accuracy of numbers we can represent or a larger range of values. It's common practice to fix the decimal point at a known location, for example 16-bits for each side. This is called the **fixed point** representation.

5.2.3 Negative Numbers and Two's Complement

One possible way to implement negative numbers in binary is to allocate one of the bits as a sign bit. It makes sense to make this the leftmost bit, as the sign of a calculation cannot be determined until the whole calculation has been done. To do this, we either need an extra bit to represent out numbers or to reduce the range of numbers that we can represent. This method is not optimal for the computer. Take the following example: 01011010 and 10100110 are the negative of each other as they add to give 0 (ignoring any overflow). They are called each other's **two's complement**.

This is a convenient way of representing both positive and negative numbers. To form a numbers two's complement:

- 1. Invert all bits
- 2. Add 1 to the LSB

If we treat positive numbers to always begin with zero, and be otherwise identical to unsigned numbers then it follows that negative numbers must begin with at least one leading one. We'll consider the sum $(7-15)_{10}$ (which is equivalent to $(00111-01111)_2$) and use two's complement to complete this. We convert 01111_2 to it's two's complement and get 10001_2 . Now we can add to get $(00111+10001=11000)_2$. As the lead digit is 1, this must be a negative number. We take its two's complement by inverting all bits, and adding 1 to get 01000_2 .

The range of numbers we can represent in two's complement is different than the number of bits we have. The largest positive value is 0111...11, and the smallest positive value is 000...01. For negative numbers, the smallest ('least negative') value we can represent is 111...11 whilst the largest ('most negative') is 100...00. Therefore, the range of values that can be represented in two's complement is from -2^{n-1} to $2^{n-1}-1$

5.2.4 Floating Point Numbers

We can now represent both integers and real numbers (positive and negative) using **fixed-point** methods. These methods have restrictions. They impose considerable restrictions on the range of numbers that are available to use. Take for example an eight bit representation in two's complement using five bits to store the integer, and three bits to store the non-integer. This is bound by the following limits:

	Two's Complement	Radix 10
Largest Positive Value	01111.111	15.875
Smallest Positive Value	00000.001	0.125
Largest Negative Value	10000.000	-16
Smallest Negative Value	11111.111	-0.125

The range of values supported is rather small, and the precision available is also very limited. Fixed point data has a limited range and limited precision. We can circumvent this partially by allocating more bits, but there are practical limitations on this. To support more accurate numerical calculation, a **floating-point number** type has been devised. Our representation has two parts: the **mantissa** M (sometimes called argument, fraction or significand) and the **exponent** E which are combined to give a value V in the following way:

$$V = M \times 2^E$$

In 8-bits, allowing four bits of mantissa and four of exponent, and placing the radix point after the first bit of the mantissa we have (using two's complement for both exponent and mantissa):

	Floating Point	Radix 10
Largest Positive Value	0.111×2^{0111}	112
Smallest Positive Value	0.001×2^{1000}	4.9×10^{-4}
Largest Negative Value	1.000×2^{0111}	-128
Smallest Negative Value	1.111×2^{1000}	-4.9×10^{-4}

As it can be seen, there is a greatly increased range of values available from the exact same amount of bits. This comes at the cost of resolution. The use on an exponent means that the increments are non-linear, and hence larger numbers have a lower resolution than smaller numbers. As a result, you should be extremely careful when using floating point numbers.

Floating-point arithmetic is rather complex and the details are beyond the scope of this module.

- 5.3 Clocked Logic
- 6 I/O and Peripherals
- 7 Improving Performance
- 7.1 Caches and Virtual Memory
- 7.2 Pipelining and Branch Prediction

List of Figures

The von Neumann Architecture	1
The MIPS processor	4
The conventions of use for all registers in the MIPS architecture	7
MIPS instruction format for register operations	7
MIPS instruction format for load/store operations	7
Datapath elements for instruction fetch	8
We have integrated an adder to increment the PC	9
Registers and their input/output ports	9
	10
	10
	11
Datapath for branch instructions	11
The integrated datapath for ALU and load/store instructions	12
(Left) Circuit symbol for an NMOS transistor. (Right) Circuit symbol for a PMOS	
transistor	12
The logical operator NOT built from transistors	12
(Left) A NAND gate built from four transistors. (Right) A NOR gate built from	
four transistors	13
Truth table for a two input decoder	13
A 2-4 decoder built from logic gates	14
Truth table for a 2-1 multiplexer	14
A 2-1 multiplexer built from logic gates	14
Logic gates which implemet a half-adder	15
A full-adder built from two half adders and an XOR gate	15
A simple ALU capable of performing addition, bitwise-AND and bitwise-OR operations	16
	The MIPS processor The conventions of use for all registers in the MIPS architecture MIPS instruction format for register operations MIPS instruction format for load/store operations Datapath elements for instruction fetch We have integrated an adder to increment the PC Registers and their input/output ports The ALU and it's input/output ports The datapath for R-type instructions The datapath for load/store instructions Datapath for branch instructions The integrated datapath for ALU and load/store instructions (Left) Circuit symbol for an NMOS transistor. (Right) Circuit symbol for a PMOS transistor The logical operator NOT built from transistors (Left) A NAND gate built from four transistors. (Right) A NOR gate built from four transistors Truth table for a two input decoder A 2-4 decoder built from logic gates Truth table for a 2-1 multiplexer A 2-1 multiplexer built from logic gates Logic gates which implemet a half-adder

\mathbf{Index}

adder, 8, 14 address bus, 8 ALU, 2 arithmetic and logic unit, 2	instructions, 1 IO controller, 2 load/store unit, 2
assembly language, 1	long-term memory, 2
bus, 2 bus controller, 2	main memory, 2, 3 mantissa, 17 microarchitecture, 1, 8
cache, 2	microcoding, 5
combinational logic, 13	multiplexer, 14
complex instruction set computer, 5 control logic, 8	multiplexing, 11
control unit, 2	NMOS, 12
data, 1	not gate, 13
data bus, 8	operating system, 1
datapath, 8 digital logic, 1	overflow, 15
entry point, 3	PMOS, 12
exclusive OR, 14	program counter, 2
exponent, 17	reduced instruction set computer, 5
fetch-decode-execute cycle, 3	register operation, 7
fixed point representation, 16	registers, 2
floating-point number, 17	sign extension unit, 10
gnd, 13	
	transistor, 1, 12
half-adder, 14	turing complete, 5 two's complement, 16
harvard architecture, 4	two's complement, 10
high level language, 1	vdd, 13
instruction execution cycle, 2	volatile, 2
instruction register, 2	von Neumann architecture, 1
instruction set, 1, 5	von Neumann bottleneck, 4