

Binary

- k bits can hold 2^k values.

• Big endian : $\begin{array}{cccc|c} 1 & 1 & 0 & 1 \\ \downarrow & \downarrow & \downarrow & \downarrow \\ 2^3 & 2^2 & 2^1 & 2^0 \end{array} = 13_{10}$ "There's a big market for big endian, as this is the 'normal' format"

Little endian : $\begin{array}{cccc|c} 1 & 1 & 0 & 1 \\ \downarrow & \downarrow & \downarrow & \downarrow \\ 2^0 & 2^1 & 2^2 & 2^3 \end{array} = 11_{10}$

- Decimal to other base: division method

Eg. $22_{10} \rightarrow ?_2$

$$\begin{array}{r} 2 \overline{)22} \\ 11 \text{ r } 0 \\ 5 \text{ r } 0 \\ 2 \text{ r } 1 \\ 1 \text{ r } 0 \\ 0 \text{ r } 1 \end{array} \quad 22_{10} = 10100_2$$

- $0.625_{10} = 1 \times 2^{-1} + 0 \times 2^{-2} + 1 \times 2^{-3}$
 $= 0.5 + 0.00 + 0.125$
 $= 0.625_{10}$

- Decimal fractional values to binary: multiplication method
- Eg. $0.625_{10} \rightarrow ?_2$

$$\begin{array}{r} 0.625 \\ \times 2 \\ \hline 1 \leftarrow 1.250 \\ 0.250 \\ \times 2 \\ \hline 0 \leftarrow 0.500 \end{array} \quad \begin{array}{r} 0.500 \\ \times 2 \\ \hline 1 \leftarrow 1.000 \\ 0.000 \end{array} \quad 0.625_{10} = 0.101_2$$

$0.000 \rightarrow \text{Stop.}$

- Hex - Base 16 - $16 = 2^4$ (groups of 4 bits)
- Octal - Base 8 - $8 = 2^3$ (groups of 3 bits)

- Bin \rightarrow Hex

Eg. $0111:0100:1011:1010_2 = 74BA_{16}$

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$
 $7_{10} \quad 4_{10} \quad 11_{10} \quad 10_{10}$

• Unsigned Binary

- No sign, so assumed as +ve

- Addition:

$$\begin{array}{r}
 10110100 \\
 11011001 \\
 \hline
 1100011.01
 \end{array}$$

x x x

* In an 8-bit system, this bit would be overflow.

- Subtraction:

$$\begin{array}{r}
 0\cancel{1}11 \\
 0110 - \\
 \hline
 0101
 \end{array}
 \qquad
 \begin{array}{r}
 0\cancel{1}1 \\
 1011 \\
 \hline
 011
 \end{array}$$

* Problem, as there is no 5th bit to borrow from.
This is underflow.

- Multiplication:

$$\begin{array}{r}
 1011 \\
 \times 101 \\
 \hline
 0000 \\
 1011 \\
 \hline
 110111
 \end{array}$$

- Division :

$$\begin{array}{r} \underline{0\ 1\ 1\ 1} \\ 110 \overline{)1\ 0\ 1\ 0\ 1\ 0} \\ \underline{0\ 0\ 0} \\ \downarrow \\ \begin{array}{l} \textcircled{x}^2 Q\ 1\ 0 \\ \underline{1\ 1\ 0} \\ \downarrow \\ \textcircled{x}^3 Q\ 1 \\ \underline{1\ 1\ 0} \\ \downarrow \\ 0\ 1\ 1\ 0 \\ \underline{1\ 1\ 0} \\ \downarrow \\ 0 \end{array} \end{array}$$

• Signed Magnitude.

- 1 bit for sign, rest for number. For sign bit, 0 = +ve; 1 = -ve.

- A k-bit system has a range of $-(2^{k-1}-1) \rightarrow (2^{k-1}-1)$, which is 2^k-1 unique values.

- Addition:

$$\begin{array}{r} 0\ 1\ 0\ 1 \\ 0\ 0\ 0\ 1 + \\ \hline \underline{0\ 1\ 1\ 0} \end{array}$$

$$\begin{array}{r} 1\ 0\ 1\ 0 \\ 1\ 0\ 1\ 1 + \\ \hline \underline{1\ 1\ 0\ 1} \end{array}$$

$$\begin{array}{r} 0\ 1\ 0\ 0 \\ 1\ 0\ 1\ 0 + \\ \hline \underline{0\ 1\ 1\ 0} \end{array}$$

$$\begin{array}{r} 1\ 0\ 1\ 1 \\ 0\ 0\ 0\ 1 + \\ \hline \underline{1\ 1\ 0\ 0} \end{array}$$

Sign bit is always symmetrical when adding 2 numbers.

* Be careful of answers falling outside of range *

• One's Complement

- One bit for sign, rest for number. For sign bit, 0=+ve, 1=-ve
- A k -bit system has a range of $-(2^{k-1}-1)$ to $(2^{k-1}-1)$, which allows for 2^k-1 unique values.
- Positive numbers: 0 + num. in binary
Eg. $5_{10} = 0101_2$ in 4-bit ones comp.
- Negative numbers: 1 + num in binary, bits flipped
Eg. $-5_{10} = 1010_2$ in 4-bit ones comp.

- Addition:

$$\begin{array}{r} 01101 \\ 11010 + \\ \hline 100111 \\ + \\ \hline 01000 \end{array}$$

Ignore any overflow on sign bit.

Always +1 in all cases.

* Be careful of answers falling outside of range *

• Two's Complement

- One bit for sign, rest for number. For sign bit, 0=+ve, 1=-ve
- A k -bit system has a range of $-(2^k)$ to $(2^{k-1}-1)$, allowing for 2^k unique values.
- Positive numbers: 0 + num. in binary
Eg. $5_{10} = 0101_2$ in 4-bit Two's Comp.

- Negative numbers: 1 + num. in binary, bits flipped, 1 added.

Eg. $+5_{10} = 0101_2$

$1010_2 \rightarrow -5_{10}$, One's C.
+1

$1011_2 \rightarrow -5_{10}$, Two's C.

- Addition:

$$\begin{array}{r} 0100 \\ 0010 \\ \hline \textcolor{pink}{\cancel{00110}} \end{array} \quad \begin{array}{r} 0100 \\ 0110 \\ \hline \textcolor{yellow}{\cancel{01010}} \end{array}$$

$$\begin{array}{r} 1110 \\ 1100 \\ \hline \textcolor{pink}{\cancel{1010}} \end{array} \quad \begin{array}{r} 1100 \\ 1010 \\ \hline \textcolor{yellow}{\cancel{1010}} \end{array}$$

If bits n and $n+1$ are the same, this works.
If different, this doesn't work.

- Subtraction:

$$\begin{array}{r} 0010 \\ 0100 \\ - 1011 \\ \hline 1100 \\ - 0010 \\ \hline 1110 \end{array}$$

Diagram illustrating the subtraction process:

- The first column (0, 0, 1, 1) is labeled "FLIP".
- The second column (0, 1, 0, 1) is labeled "+1".
- The third column (1, 1, 0, 0) is labeled "COPY".
- The fourth column (1, 1, 1, 0) is labeled "ADD".

* Be aware of answers that will be out of range. *

- Multiplication : use Booth's Algorithm

Eg.

$$\begin{array}{r} 1011 \\ 0110 \quad \times \\ \hline 0000 \quad 0000 \\ 0000 \quad 1010 \\ 0000 \quad 0000 \\ \hline 1101 \quad 1000 \\ \hline 1110 \quad 0010 \end{array}$$

$01 \rightarrow A: 1111 \ 1011$
 $10 \rightarrow S: \ 000 \ 0101$

A in top line
S in top line in Two's Comp.
Sign-extend both to 2x bit space.

Step in a bit each time

- Pairs:
- 00 - Do nothing
 - 01 - Use A
 - 10 - Use S
 - 11 - Do nothing

Add values as if they are unsigned, ignoring any overflow.

• IEEE - 754

- Single: 1-bit sign, 8-bit exponent, 23-bit mantissa
- Double: 1-bit sign, 11-bit exponent, 52-bit mantissa
- Single has 127 bias on exponent.
 - To show exp. of 1, show $127+1 = 128_{10} = 1000\ 0000$
 - To show exp of -13, show $127-13 = 114_{10} = 0111\ 0010$
- Double has 1023 bias on exponent.

- Mantissa must be normalised:

$$\begin{aligned} \circ 1011.1101 &= 1.0111101 \times 2^3 \\ \circ 0.00101 &= \underbrace{0.01}_{\text{Mantissa}} \times 2^{-3} \end{aligned}$$

- Eg. $3.75_{10} = 11.01_2 = 1.011_2 \times 2^1$ Exp. = $127+1 = 1000\ 0000_2$

↳ $0100\ 0000\ 0111\ 0000\ 0000\ 0000\ 0000$

↕ ↓ ↓
 Sign Exponent Mantissa
 $0 = +ve$
 $1 = -ve$

- Eg. $12.625_{10} = 1100.101_2$

$$\begin{aligned} &= 1.100101 \times 2^3 \rightarrow 127+3 = 130_{10} \\ &= 1000\ 0010_2 \end{aligned}$$

$0100\ 0001\ 0100\ 1010\ 0000\ 0000\ 0000$

Main Computer Components

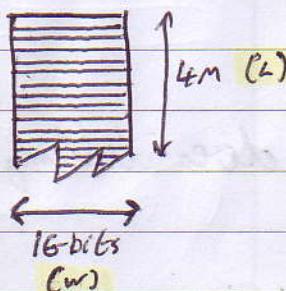
- 1) CPU
- 2) Memory
- 3) Bus
- 4) Clocks / Interrupts
- 5) I/O

1) CPU

- Two main parts: **data path** and **control unit**
- The **data path** is made up of two parts:
 - **Arithmetic Logic Unit (ALU)** - used for arithmetic operations and comparisons.
 - **Registers** - temporary storage areas used to store arrays of bits for the next operation.
- The **control unit** works like a traffic manager to tell the **ALU** what to do and handle events like interrupts.
 - The CU contains and uses registers including the program counter and the status register.

2) Memory

- Array of addressed slots, each with a word of data inside.
- Described as **length × width**, i.e. "4m × 16" means 4m rows, each of 16 bits in length:
- Memory can be word- or byte-addressable.



Log 4m means 2^{2^2}
addresses, from
 $2^2 \times 2^{20} = 2^{22}$
4 m

3) Bus

- Wire connections between components to move data around.
- Types:
 - Point-to-point
 - Fast and dedicated communication channels
 - Multipoint
 - Potential problems with collisions
- Bus arbitration is needed to avoid collisions - master/slave is common.
 - Types:
 - Daisy chain, priority list.
 - #1 master by default, #2 master if #1 is busy, etc. claim the line
 - Problem: lower priority components can be crowded out.
 - Centralised parallel
 - Coordinator - single centralised arbitration circuit decides who can use the bus.
 - Distributed arbitration
 - Self-selection - components decide between themselves who uses the bus.
 - Collision detection - components send data, if collision, re-send.
- Lines on a bus:
 - Data lines
 - Control lines
 - Arbitration does its thing here
 - Address lines
 - Power lines.

4) Clocks

- Synchronises operations within machine incl. moving data and CPU operations
- Measured in MHz (10^6 s^{-1}) and GHz (10^9 s^{-1}).

4) Interrupts

- Alerts CPU to higher-priority events
- Triggered by arithmetic errors, divide, I/O requests, over/underflow, program errors, etc.
- Each interrupt is associated w/ a specific procedure to deal w/ the event: **Interrupt Service Routine (ISR)**

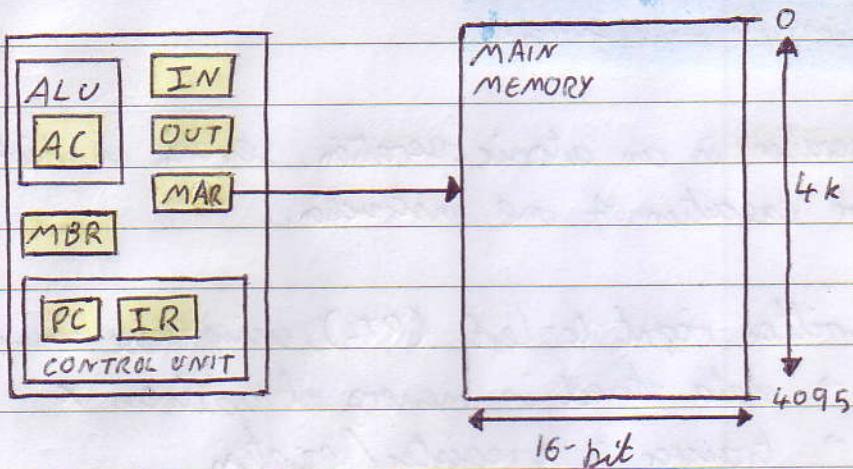
* More on Interrupt Processing later *

5) I/O Subsystem

- Communicate w/ outside world
- Memory-mapped
 - I/O interface appears to CPU like memory
- Instruction-based
 - Dedicated instructions for CPU to handle I/O device
 - Doesn't occupy memory, but requires instruction set

MARIE/Assembly

- Characteristics:
 - Uses 2's comp. binary
 - $4k \times 16$ -bit word-addressable memory.
 - 16-bit ALU
 - 7 registers



- Common data bus connecting registers and memory, but with some P2P buses such as $\text{MAR} \leftrightarrow \text{Memory}$.

• ISA

- 16-bit operations - 4-bit opcode and 12-bit operand/address.

Inst. #	Inst.	Meaning
0001	LOAD X	Load value of X into AC
0010	STORE X	Store value of AC at address X
0011	ADD X	Add value of X to AC and store in AC
0100	SUBTR X	Subtract " " " from " " " " "
0101	INPUT	Read input value into AC
0110	OUTPUT	Send AC value to output
0111	HALT	Terminate
1000	Skip Cond	Skip next line on condition *
1001	Jump X	Takes value at X and places into PC.

* Skip Cond:

1000 1000 0000 0000

↳ 00 → skip if AC < 0

01 → skip if AC = 0

10 → skip if AC > 0

• Register Transfer Language

- A micro-operation is an atomic operation, several of which are required for execution of one instruction.

- RTL is written right-to-left (RTL), using these conventions:

- $M[x]$ → data stored in memory at address x
- \leftarrow → transfer to a register/location

Eg. LOAD X : $MAR \leftarrow x$

$MDR \leftarrow M[MAR]$

$AC \leftarrow MDR$

• Instruction Processing

- Fetch: get instruction from memory, load into IR.

- Decode: check op code, decide what to do, check for an operand, place operand into MDR

- Execute: execute the instruction

Start → Init. PC → Copy PC → $IR \leftarrow M[MAR]$

→ to MAR

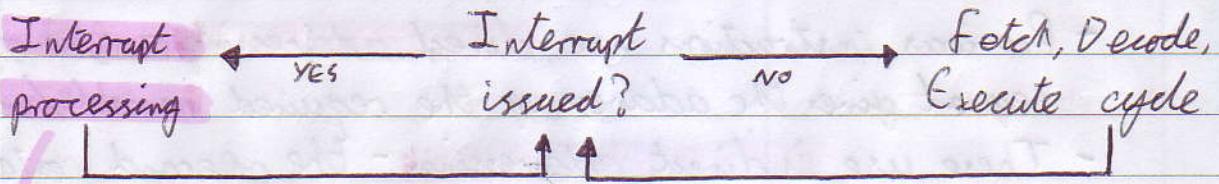
$PC \leftarrow (PC+1)$

Execute ← Decode inst.

$MAR \leftarrow IR[11-0]$

Put the last 12 bits in
MAR to handle any operand

I Interrupt



- Start
- ↓
- Interrupt signal detected
- ↓
- Save variables and registers
- ↓
- Look up ISR address
- PC ← ISR addr
- ↓
- Branch to ISR
- ↓
- Start F/D/E cycle on ISR
- ↓
- Restore variables and registers
- ↓
- Return to F/D/E cycle

* Generally, interrupts are disabled during the ISR, because they are "maskable" interrupts. "Non-maskable" interrupts must be processed to keep the system stable.

• Extended Instruction Set

- Previous instructions use **direct addressing** - the operand segment gives the address of the required variable/value
- These use **indirect addressing** - the operand contains the address of a location in memory, which itself contains the address of the value required.

- ADDI x (Indirect Add)

$MAR \leftarrow x$

$MBR \leftarrow m[MAR]$

$MAR \leftarrow MBR$

$MBR \leftarrow m[MAR]$

$AC \leftarrow AC + MBR$

- JUMPI x (Indirect Jump)

$MAR \leftarrow x$

$MBR \leftarrow m[MAR]$

$PC \leftarrow MBR$

- CLEAR (Clear AC)

$AC \leftarrow 0$ (zero)

- JNS x (Jump and Store)

$MBR \leftarrow PC$

$MAR \leftarrow x$

$m[MAR] \leftarrow MBR$

$MBR \leftarrow x$

$AC \leftarrow 1$

$AC \leftarrow AC + MBR$

$PC \leftarrow AC$

Store PC at x and
jump to $x+1$.

Instruction Set Architectures (ISAs)

• ISA Characteristics:

- Bits per instruction
- Stack- or register-based
- Operands per instruction
- Operand location
- Operand type(s)/size
- Types of operations

• ISA Measurements:

- Main memory space occupied by program
- Instruction complexity
- Instruction length in bits
- Total instructions in instruction set

• ISA Design Considerations

- Instruction length (short/long/variable)
- Number of operands
- Number of addressable registers
- Memory organisation (byte- or word-addressable)
- Addressing modes (direct, indirect, indexed, etc.)
- Byte ordering/Endianness

◦ Big endian:

- More natural
- Sign of the number can be found in byte at offset zero.
- Strings and numbers stored in the same order.

◦ Little endian:

- Easier to place values on non-word boundaries
- Conversion from a 16-bit integer to 32-bit does not require any arithmetic

• CPU Data Storage Structure

- The tradeoffs in these choices are simplicity and cost of hardware design, versus execution speed and ease of use.

- Accumulator Architecture

- One operand of a binary operation is implicitly in the accumulator
- One operand is in memory, creating lots of bus traffic

- Stack Architecture

- Instructions and operands are taken from the stack
- The result of a binary operation is implicitly stored on top of the stack
- A stack cannot be accessed randomly.

- General Purpose Register Architecture

- Registers can be used instead of memory
- Faster than accumulator architecture
- Efficient implementation for compiler
- Results in longer instructions.

* Most systems today use a GPR *

- There are 3 types:

- Memory - memory, where two or three operands may be in memory
- Register - memory, where at least one operand must be in a register
- Load-store, where no operands may be in memory.

• Architectural Differences

- A stack architecture requires different ordering for arithmetic operations.

◦ Normally we use infix : $Z = X + Y$

◦ Stacks require postfix : $Z = XY +$

◦ Postfix makes brackets obsolete:

- Infix : $Z = (A * B) + (C * D)$

- Postfix : $Z = AB * CD * +$

- In a stack ISA, this might look like : PUSH A

PUSH B

MULT

PUSH C

PUSH D

MULT

ADD

POP Z

- In a one-address ISA (like MARIE), the infix $Z = (A * B) + (C * D)$ might look like :

LOAD A

MULT B

STORE TEMP

LOAD C

MULT D

ADD TEMP

STORE Z

- In a two-address ISA this is:

```
LOAD R1, A  
MULT R1, B  
LOAD R2, C  
MULT R2, D  
ADD R1, R2  
STORE Z, R1
```

- In a three-address ISA this is:

```
MULT R1, A, B  
MULT R2, C, D  
ADD Z, R1, R2
```

- Instructions fall into a range of broad categories:

- Data movement
- I/O
- Arithmetic
- Control transfer
- Boolean
- Special purpose
- Bit manipulation

Memory Addressing

- Addressing modes specify where an operand is located
- They can specify a constant, a register or memory location
- The actual location of the operand value is called the effective address.
- Some addressing modes allow us to determine the address of an operand dynamically.
- Modes:
 - Immediate addressing - data is part of the instruction
 - Direct addressing - address of the data is given in the instruction
 - Register addressing - data is located in a register
 - Indirect addressing - the address of a location storing the address of the data is given in the instruction
 - Register Indirect addressing - a register stores the address of a location storing the address of the data
 - Indexed addressing - a register (implicitly or explicitly) is used as an offset, which is added to the address in the operand to determine the effective address.

- **Based addressing** - similar to indexed addressing, but a base register is used instead of an index register.
- **Stack addressing** - the operand(s) is/are assumed to be on top of the stack

* Difference between **indexed** and **based**:

- An **index register** holds an offset relative to the address given in the instruction
- A **base register** holds a base address where the address field represents a displacement from this base.

Example:

Memory:	Register:	LOAD 800
800 → 900	R1 → 800	
...		

900 → 1000

1000 → 500

...

1100 → 600

...

1600 → 700

Accumulator:

Mode	Value
Immediate	800
Direc	900
Indirec	1000
Indexed	700

* **Index:** instruction gives an address, added to the value in a register.

* **Based:** instruction gives a value, added to the address in a register.

Memory

- Two main kinds:

- RAM

- Dynamic RAM (DRAM) - cheap, simple design

- Static RAM (SRAM) - fast, used in caches

- ROM

- Low power consumption

- Hierarchy

- Faster = More expensive

- Faster elements are closer to the CPU:

- Registers

- L1 Cache

- L2 Cache

- Main memory

- Fixed rigid disk

- Optical disk

- Magnetic tape

} System

- Online

- Near line

- Offline

- To access data, CPU checks closest memory first (caches) and works away until it finds a match
 - Once located, the data and nearby elements are fetched into cache.

- "hit" / "miss" → data is found/not found

- "hit rate" / "miss rate" → % of hits/misses for a given level

- "hit time" → time to access data at a given level

- "miss penalty" → time to process a miss, including replacing a block of memory and delivering data to processor

- After a hit, entire blocks are copied, because locality tells us nearby bytes will be needed soon.
- Types of locality:
 - Temporal → recently-accessed data tends to be accessed again
 - Spatial → access tends to cluster
 - Sequential → instructions tend to be accessed sequentially.

• Cache Memory

- Stores recently-used data near CPU to speed up access
- Smaller and faster than main memory
- Content-addressable, so doesn't use normal addresses
 - Smaller is sometimes better, as it takes a long time to search
 - The "content" that is addressed is a subset of bits from the main memory address called a field
 - The fields provide a many-to-one relation between the main memory and the cache
 - A tag field in the cache block distinguishes one memory block from another.

• Direct Mapped Cache

- Simplest type of cache.
- In an N block cache, block x of main memory maps to $x \bmod N$, so in a 10 block cache block 7 maps to 7, 17, 27, etc of main memory.

- Once a block is loaded, a valid bit is set to show that the block contains valid data.

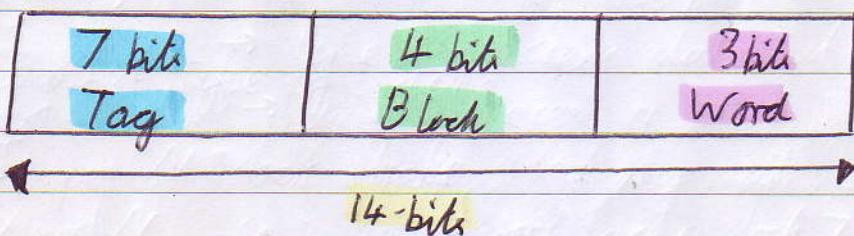
- Example cache:

Block	Tag	Data	Valid
0	0000 0000	Words A, B, C...	1
1	1111 0101	Words L, M, N...	1
2	-	?	0
3	-	?	0

- Block 0 contains multiple words from memory identified with the tag 0000 0000.
- Block 1 contains words identified with the tag 1111 0101.
- Blocks 2/3 are not valid.

- The size of the field into which a memory address is divided depends on the size of the cache.

- o Suppose memory has 2^{14} words, cache has 16 blocks, and a block has 8 words.
- o Memory has $\frac{2^{14}}{2^3} = 2^8$ blocks
- o A main memory address requires 14 bits
- o The three rightmost bits are needed for the word (3 bits are needed to identify one of 8 words in a block)
- o We need 4 bits to select one of the 16 cache blocks.
- o The remaining 7 bits make up the tag field.
- o This gives this format:



• Example: Main memory address: 1AA
In 14-bit binary : 00 0001 1010 1010

This is split up as:

000001	1010	1010
Tag	Block	Word

All 8 words of 0000010101000 - 0000010101111 are loaded into block 0101 of cache with the tag 0000011.

1AA, 1AB, 1AC, etc... will all give cache hits

3AB will cause the whole block to be evicted and replaced.

If the program requests 1AB, 3AB, 1AD, etc., the advantage is lost.

- Disadvantage of direct mapping: repeated calls to different main mem. blocks that would occupy the same cache block destroy any theoretical advantage.

Fully Associative Cache

- A block can go anywhere in the cache
 - The cache can be filled before eviction is necessary.
- Memory address has only two fields: tag and word.
- Suppose again we have 14-bit memory addresses, a 16-block cache and an 8-word block.
The format is : 11 bits 3 bits
 Tag Word
 the rest ↗ ↴ to reference 8 unique words
- All tags are searched in parallel to retrieve data quickly
 - Requires complex, expensive hardware.
- When the cache is full, a block must be evicted.
 - This is the victim block, chosen by methods discussed later.

• Set-Associative Cache

- Combines the ideas of direct mapped / fully associative.
- A main memory block maps to a set of cache blocks, but not "any" block, as with fully-assoc. limited
- The number of blocks per set depends on system design
- Example: a 2-way set-assoc. cache
 - o Each set contains two different blocks:

Set	Tag	Block 0 of Set	Valid	Tag	Block 1 of Set	Valid
0	0000 0000	A,B,C...	1	—	?	0
1	1111 0101	L,M,N...	1	—	?	0
2	—	?	0	P,Q,R...	?	1
3	—	?	0	T,U,V...	?	1

- A memory address takes 3 fields: tag, set and word.
 - o As before word identifies the word in the block and tag uniquely identifies the block.
 - o Set determines the set to which the memory block maps
- Example: Main memory of 2^{14} words, 2-way set-assoc. cache of 16 blocks, and 8 words per block.
 - o 16 cache blocks and 2-way cache = $16/2 = 8$ sets.
 - o Thus 3 bits for word, 3 bits for set and the leftover 8 bits for the tag.

• Replacement Policies

- An optimal policy would look into the future and determine the block that will be needed for the longest time.
 - Impossible, but used as a benchmark.

- Choice of policy used depends on type of locality being optimised, usually temporal

- Least Recently Used (LRU)

- Evicts the block unused for the longest time
- A history must be maintained - this slows the cache

- First in, First Out (FIFO)

- Evicts the block that has been in cache longest, despite recent use rates.

- Random

- Can evict a block needed soon, but avoids thrashing

- These policies must consider dirty blocks that have been updated while in cache. These must be written back to memory.
 - Two write policies:

- Write through : cache/mem updated on every write
 - Reliable, but slows down updates (most accesses are read, so often negligible)

- Write back : update memory only when block selected for replacement

- Less traffic, but memory may not match cache (bad if there are many concurrent users).

• Effective Access Time (EAT)

- EAT is a weighted average taking into account the hit rate and relative access time of different levels of memory.
- The EAT for 2-level memory is:

$$EAT = H \times \text{Access}_c + (1-H) \times (\text{Access}_c + \text{Access}_{mn})$$

When we hit

Access
cache

When we miss

Access cache, get nothing,
so access main mem.

• H = hit rate (%)

• $\text{Access}_{c/mn}$ = Cache/Main Mem. access times.

- The equation can be extended for many levels of memory.

• Types of Cache

- We have discussed a **unified/integrated cache**, with data and instructions together.
- Most systems use **separate caches**
 - Called a **Harvard Cache**
 - Better locality, but more complex
 - Bigger caches can give similar improvements without complexity
- A **victim cache** holds recently evicted blocks to help improve performance.
- A **trace cache** holds decoded instructions for program branches, giving the illusion that noncontiguous instructions are contiguous.
- Many systems use **multiple caches**, creating a mini memory hierarchy.
 - **Inclusive caching**: data can exist on multiple cache levels
 - **Strictly inclusive** : data on one level must exist on all levels below (ie. all larger caches)
 - **Exclusive caching** : only one copy of data allowed.
 - Tradeoffs of the three models include access time, memory size and circuit complexity

Virtual Memory

- Enhances performance by expanding memory capacity without adding memory.
- A portion of the disk serves as an extension of main memory.
- If the system uses **paging**, VM partitions the main mem. into individually managed **page frames** that are written (paged) to disk when not immediately needed.
- **A physical address** is the actual memory address of physical memory
- **Virtual addresses** are mapped to physical addresses by the memory manager.
- **Page faults** occur when a logical address requires a page be brought in from memory.
- **Memory fragmentation** occurs when the paging process results in small, unusable clusters of memory addresses.
- Main and virtual memory are split into equal sized **pages**.
- A process's entire required address space need not be all in main mem. at once - some can be on disk
- Pages allocated to a process do not need to be stored **contiguously** - either on disk or in memory.
- Only the needed pages are in memory - unneeded pages are on slower disk storage.

- Info on the location of each page - disk or mem - is stored in a **page table**:

Virt. Mem. Phys. Mem.

			Page	Page Table	
				Frame #	Valid Bit
0		0	0	2	1
1		1	1	-	0
2		2	2	-	0
3		3	3	0	1
4			4	1	1
5			5	-	0
6			6	-	0
7			7	3	1

- A process generates a **virtual address** and the OS translates it to a **physical address**.
 - o To do this, the virtual address has two fields: the **page** and the **offset**.
 - o Page determines the page; offset determines the location in that page.
 - o Logical page number is translated into a physical page frame through a look-up in the page table.
 - o If the valid bit is 1, the virtual page number is replaced w/ the physical page number.
 - o If the valid bit is 0, the page must be fetched from disk.
 - This is a **page fault**
 - If necessary, a page is evicted from memory, replaced with the page from disk and the valid bit set to 1.
 - o The data is then accessed by adding the offset to the physical frame number.

• EAT

$$EAT = (1-F) \times (2 \times Access_{mn}) + F \times Access_o$$

Where F = page fault rate (%)

$Access_{mn}$ = Main mem access time

$Access_o$ = Time to load a page from disk and access data

2 memory accesses: one for page table, one for the actual data

- This double access affects EAT.
 - Because page tables are read constantly, they are stored in a translation look-aside buffer (TLB) cache that stores the mapping of virtual pages to physical pages in a special associative cache.

• Segmentation

- Instead of equal-size pages, virtual address space is divided into variable-length segments (often under control of the programmer).
- Segments are located through the segment table, containing the segment's mem. location and bounds (thus indicating size).
- After a page fault the segment is loaded from disk and the OS searches for a location in memory large enough to hold the segment.

* Issues

- Paging and segmentation can cause fragmentation
 - Paging is subject to internal fragmentation, where a process may not need all of the space in a page, leaving many pages with unused fragments.
 - Segmentation is subject to external fragmentation, where contiguous chunks of memory become broken up as segments are allocated and deallocated.
- Large page tables are slow, but with their uniform memory mapping, page operations are fast.
- Segmentation allows fast access to the segment table, but segment loading is labour-intensive.
- Each segment has a page table. This means a memory address has three fields: segment, page, offset.

• Summary

- Cache memory gives faster access to main memory.
- Virtual memory uses disk storage to give the illusion of larger main memory.
- Cache maps blocks of main mem. to blocks of cache mem.
- VM maps page frames to virtual pages.
- Three general types of cache
 - Direct mapped
 - Fully-associative
 - Set-associative
- With fully-assoc, set-assoc and VM, replacement policies must be in place
- Replacement policies include LRU, FIFO and random.
 - These must also take into account dirty blocks.
- VM must deal with fragmentation
 - Internal for paging
 - External for segmentation

Instruction-Level Pipelining

- Some CPUs split Fetch/Decode/Execute into smaller steps.
- These steps can be executed in parallel to increase throughput - this is called instruction-level pipelining (ILP).

Example

Suppose the F/D/E cycle had these six stages:

- | | |
|---|-------------------------|
| S1) Fetch instruction | S4) Fetch operand(s) |
| S2) Decode opcode | S5) Execute instruction |
| S3) Determine effective address of operands | S6) Store result |

Suppose we also have a 6-stage pipeline. The steps for each instruction can be overlapped, as so:

Clock cycle:	1	2	3	4	5	6	7	8	9	10
Inst. #1:	S1	S2	S3	S4	S5	S6				
Inst. #2:		S1	S2	S3	S4	S5	S6			
Inst. #3:			S1	S2	S3	S4	S5	S6		

Speed Up

- The theoretical speedup can be determined:
 - Let t_p be the time per stage.
 - Each instruction represents a task, T , in the pipeline.
 - The first task takes $k \times t_p$ to complete in a k -stage pipeline.
 - The remaining tasks complete at a rate of one per cycle, so the remaining time is $(n-1)t_p$, assuming there were n tasks in total.

- Therefore, to complete n tasks in a k -stage pipeline requires:

$$(k \times t_p) + (n-1)t_p = (k+n-1)t_p$$

- Speedup can be calculated as:

$$\text{Speedup } S = \frac{nkt_p}{(k+n-1)t_p}$$

← Normal time
 ← Pipeline time

- As $n \rightarrow \infty$, $(k+n-1) \rightarrow n$, so the theoretical speedup is:

$$\text{Speedup } S = \frac{k t_p}{t_p} = k$$

Assumptions

- These equations make a few assumptions:
 - The architecture supports fetching instructions and data in parallel
 - The pipeline can be kept filled all the time, which is not always the case
 - Pipeline hazards can cause conflicts, flushes or stalls

Pipeline Hazards

- Resource conflicts
- Data dependencies
- Conditional branching/jumping

* More on hazards due to branching later *

Benchmarking and Performance

• Performance Equation

- Describes CPU performance

- CPU time = seconds

$$\text{program} \times \frac{\text{instructions}}{\text{program}} \times \frac{\text{cycles}}{\text{instruction}} \times \frac{\text{seconds}}{\text{cycle}}$$

average

Affected by clock speed

Affected by ISA design

• Amdahl's Law

- Describes the speedup effect of one component in a whole system.

$$\text{speedup} = \frac{\text{old exec. time}}{\text{new exec. time}}$$

$$= \frac{(1-f) + f}{(1-f) + f/k}$$

$$= \frac{1}{(1-f) + f/k}$$

rest of system	old component
$1-f$	f
rest of system	new component
$1-f$	f/k

Where...

f is the fraction of work done by one component

k is the speedup factor of that component

Benchmarking

- ISA, CPU speed, etc. all affect performance in complex ways
- Benchmarking allows a direct, comparable measure of performance independently of these factors

Synthetic Benchmarks:

- Whetstone/Dhrystone

- Linpack

- library calls, math routines

- linear algebra routines

- These were easy to understand, but easy to optimise for and thus not useful.

Improved Benchmarks:

- A group called SPEC designed different forms of benchmark for computer/applications, including:

- CINT2000 - tested CPU integer operations

- CFP2000 - tested CPU floating point operations

- These report peak and base performance:

- Peak - compiler optimisation allowed

- Base - compiler optimisation not allowed

• Factors Affecting Performance

- ISA Design

- Two main notions in ISA design optimization:
 - RISC - Reduced Instruction Set Computing
 - CISC - Complex Instruction Set Computing

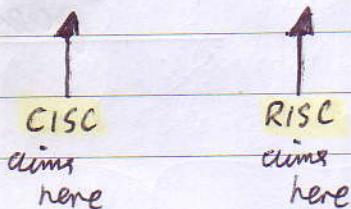
RISC

- Few, fixed-length instructions
- Same number of clock cycles to execute each instruction
- Memory access only through LOAD/STORE operations
- Easier to hardware
- Requires more registers

CISC

- Many, variable length instructions
- Variable number of clock cycles to execute an instruction
- Any instruction can access the memory
- Requires microcode to interpret instructions as they arrive

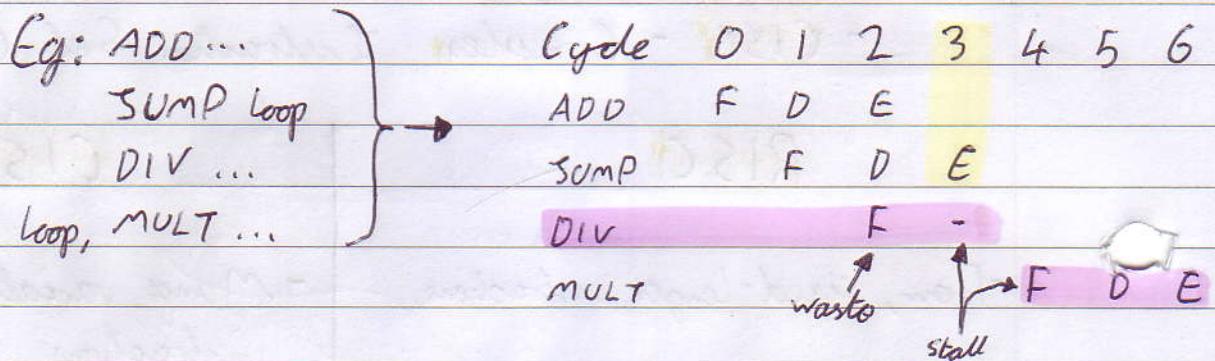
$$\circ \text{Performance} = \frac{\text{instructions}}{\text{program}} \times \frac{\text{cycles}}{\text{instruction}} \times \frac{\text{seconds}}{\text{cycle}}$$



- Pipelining

◦ Can have a +ve impact; see earlier

◦ Conditional branching/jumping causes hazards



In cycle 3 the jump was executed and the program branched to MULT, however in cycle 2 the DIV was fetched creating wasted work and a stall in the stack

◦ There are methods to combat this:

- Delayed branching - order of instructions is adjusted to avoid hazards and improve speedup

- Branch prediction - the outcome of a branch is "guessed" before it is executed

- Hard Disk Usage

- HDDs create bottle necks as they are often the slowest parts
- $$\text{Disk utilisation} = \frac{\text{request arrival rate}}{\text{disk service rate}}$$
 ← Rate of "asks" in requests/second
← Rate of "replies" in I/O ops/second
- This is the probability the disk is busy when an I/O request arrives.
- $$\text{Request wait time} = \frac{\text{service time} \times \text{disk utilisation}}{1 - \text{disk utilisation}}$$

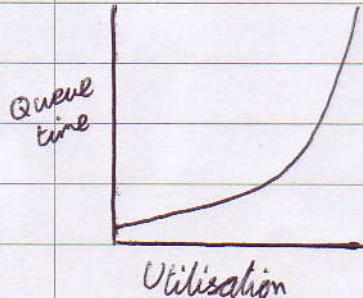
Eg. Service time = 0.015,
Service rate = $\frac{1}{\text{service time}} \approx \frac{66}{\text{sec}}$ } Utilisation = $\frac{33}{66} = 0.5$
Request rate = $33/\text{sec}$ }

$$\begin{aligned}\text{Queue time} &= \frac{0.015 \times 0.5}{0.5} \\ &= 0.015_s \quad (= 15\text{ms})\end{aligned}$$

Change request rate to 60/sec

$$\text{Utilisation} = \frac{60}{66} \approx 0.9$$

$$\begin{aligned}\text{Queue time} &= \frac{0.015 \times 0.9}{0.1} \\ &= 0.135_s \quad (= 135\text{ms})\end{aligned}$$



Rule of thumb: 80% is seen as the max. acceptable utilisation

- Hard Disk Scheduling Approaches

- First Come, First Serve (FCFS)

- + Simple and fair
- Large, random seeks
- I inefficient

- Shortest Seek Time First (SSTF)

- + Minimises arm movement
- Unfair, as remote requests wait a long time

- SCAN

- + Always scans the entire disk
- + Predictable, fair
- Extra arm motion

- C-SCAN

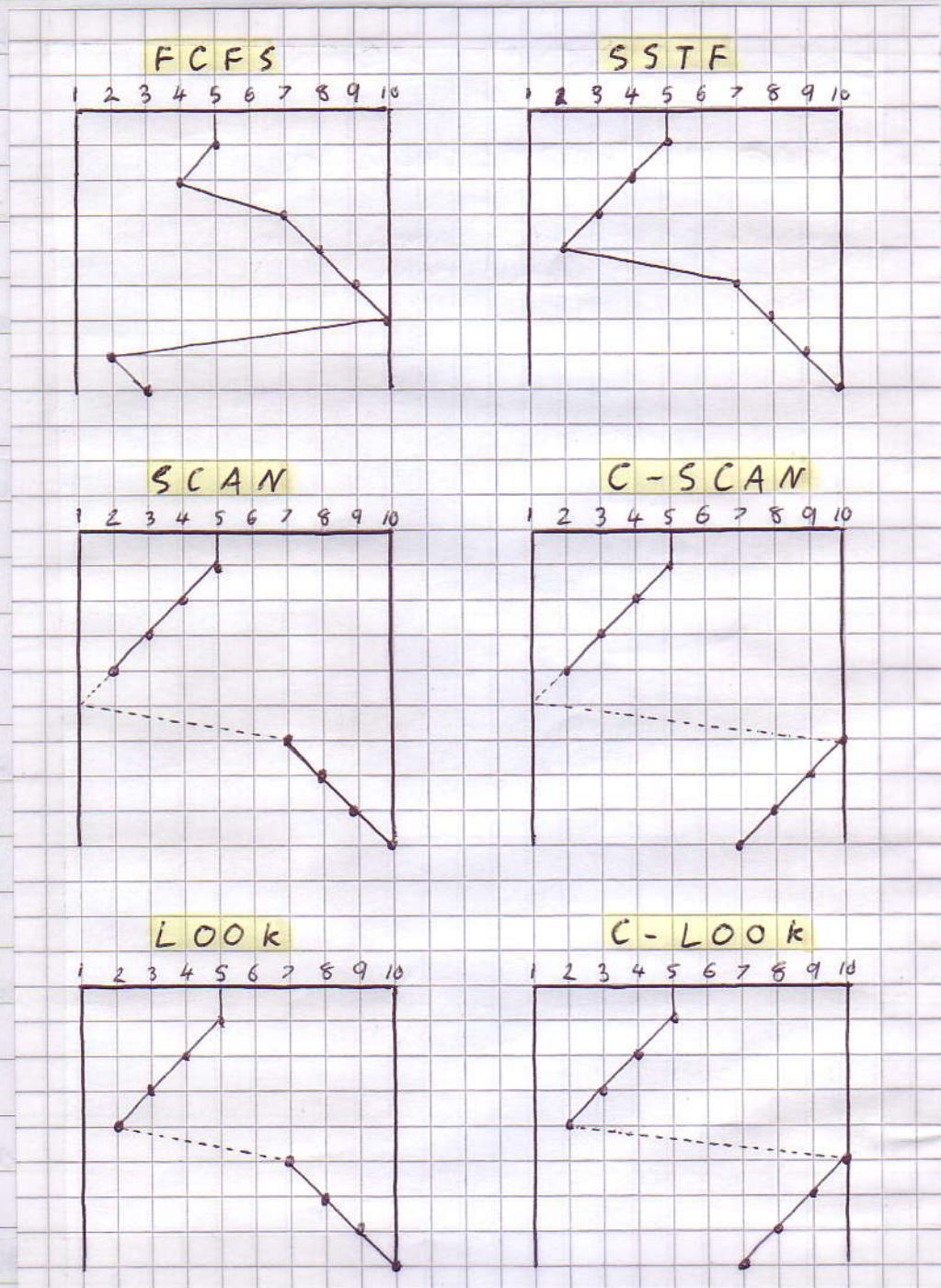
- + Like SCAN, but wraps around
- + Simpler to implement

- LOOK / C-LOOK

- + Like SCAN/C-SCAN, but only scans to outmost track in queue.
- + LOOK has best theoretical performance
- LOOK is most complex to implement

o Examples:

- Assuming the queue is: 5, 4, 7, 8, 9, 10, 2, 3
- Assuming the read head started at 5



- Disk Cache Memory

- Prefetching of data from disk into a localised cache
 - Relies on principle of locality
- Prefetching is subject to cache pollution, when the cache is filled with data that no process needs, leaving less room for useful data.
 - Replacement algorithms like LRU, LFU and random combat this
 - Some systems evict bytes once they are written to the disk (the cache is also a staging area for data going to the disk).
- The disk cache is volatile, but app. believes data is written to disk
 - Power outage = lost data
 - Some have batteries
- Another approach is write-through
 - Data stays in cache but is written to disk immediately
 - OS is signalled that I/O is done only when data is on disk
 - Performance ↓, Reliability ↑
- If throughput is more important than reliability, some use write-back
- Some employ opportunistic writes
 - Dirty blocks remain in cache until arrival of a request for the same cylinder
 - The write operation is then "piggybacked" onto the request
 - Reduces read performance, but increases for write.