

# DDCA

by dcamenisch

## Binary Numbers

An N-bit binary number ranges from 0 to  $2^N - 1$ .  
The rightmost bit is the least significant bit **LSB**.

The **MSB** is defined as the opposite.

**Big Endian:** größtes Byte zuerst

**Little Endian:** kleinstes Byte zuerst

$$2^0 = 1 \quad 2^3 = 8 \quad 2^6 = 64 \quad 2^9 = 512$$

$$2^1 = 2 \quad 2^4 = 16 \quad 2^7 = 128 \quad 2^{10} = 1024$$

$$2^2 = 4 \quad 2^5 = 32 \quad 2^8 = 256 \quad 2^{11} = 1'048'576$$

Hex. Dez. Binary	Hex. Dez. Binary	Hex. Dez. Binary	Hex. Dez. Binary
0 0 0000	4 4 0100	8 8 1000	C 12 1100
1 1 0001	5 5 0101	9 9 1001	D 13 1101
2 2 0010	6 6 0110	A 10 1010	E 14 1110
3 3 0011	7 7 0111	B 11 1011	F 15 1111

**2's complement:** negative numbers go by inverting every bit then adding 1, MSB used as sign flag.

## Boolean Algebra

Rules:  $X + X = X$     $X + \bar{X} = 1$     $X \cdot (Y + Z) = (X \cdot Y) + (X \cdot Z)$   
 $X \cdot X = X$     $X \cdot \bar{X} = 0$     $X \cdot (Y \cdot Z) = (X \cdot Y) \cdot (X \cdot Z)$

De Morgan:  $(X + Y + Z + \dots) = \bar{X} \cdot \bar{Y} \cdot \bar{Z} \dots$ ,  $(X \cdot Y \cdot Z \dots) = \bar{X} + \bar{Y} + \bar{Z} + \dots$

## Product of Sum:

Sum of Product:		
A	B	X
0	0	1
1	0	$\bar{A} + B$
0	1	$\bar{A} \cdot B$
1	1	$\bar{A} + \bar{B}$

Product of Sum:		
A	B	X
0	0	0
1	0	1
0	1	0
1	1	1

maxterm      minterm

$X = (\bar{A} + B) \cdot (\bar{A} \cdot \bar{B})$        $X = (\bar{A} \cdot \bar{B}) + (A \cdot B)$

## NAND Operations

NOT:  $\overline{AA} = \overline{A}$    OR:  $\overline{AA \cdot BB} = \overline{A} + \overline{B}$    AND:  $\overline{AB} \cdot \overline{AB} = \overline{A} \cdot \overline{B}$

**Karnaugh Maps:** used to minimize boolean equations, they work well with up to 4 variables. Some rules:

- use fewest circles possible
- only size  $2^n \times 2^m$
- all must only contain 1's

A	B	C	X
0	0	0	1
0	0	1	1
0	1	0	0
0	1	1	0
1	0	0	1
1	0	1	1
1	1	0	1
1	1	1	0

$\Rightarrow X = \overline{B} + A \cdot \overline{C}$

NR Operation:  $\overline{A + A} = \overline{A}$    OR:  $\overline{(A + B)} + \overline{(B + C)} = \overline{A} + \overline{B} + \overline{C}$    AND:  $\overline{(A + B)} \cdot \overline{(A + C)} = \overline{A} \cdot \overline{B} + \overline{A} \cdot \overline{C}$

## Logic Gates

AND			NAND			OR			NOR		
A	B	X	A	B	X	A	B	X	A	B	X
0	0	0	0	0	1	0	0	0	0	0	1
0	1	0	0	1	1	0	1	1	0	1	0
1	0	0	1	0	1	1	0	1	1	0	0
1	1	1	1	1	0	1	1	1	1	1	0

XNOR			XOR			NOT		
A	B	X	A	B	X	A	X	
0	0	1	0	0	0	0	1	
0	1	0	0	1	1	0	1	
1	0	0	1	0	1	1	0	
1	1	1	1	1	0	1	0	

## Combinational / Sequential Logic

- Comb:** - no memory  
- no cyclic paths  
- combines inputs to get output  
- signals are assigned in every possible condition  
- all addends in sensitivity list
- Seq:** - has memory  
- depends on prior inputs  
- ex. Flip-Flops (register)

## Finite State Machine

Goes through different state, where each state depends on the prev. state and the input.

Moore: Output depends only on current state

Mealy: Output depends on the current state and the input

- State Encodings:** - Binary Encoding (minimizes flip-flops), 00, 01, 10, 11  
- One-Hot (max. flip-flops, min. next state logic), 0000, 0001, 0010, 0100  
- Output (min. output logic)

## Designing a FSM

- identify inputs/outputs
- state transition diagram
- write state transition & output table
- write boolean equations for next state

## Area of FSM

# FF = # bits for state  $\times 2$

# logic gates = count next state and output logic

## Correctness of state diagram

- reset-line
- not multiple transitions for the same input
- no missing transitions
- no unmarked transitions
- initial state (if no reset)
- no mix of Moore/Mealy labeling

## MIPS

J-Type: Jump / Branch Instructions

R-Type: Register for all operands ( $OP=0$ )

I-Type: Instructions with an immediate/constant value

The **caller** calls a function, while the **callee** gets called.

The caller needs to take care of the temporary registers \$t0 - \$t9, while the callee needs to save and restore the preserved registers \$s0 - \$s7.

## ISA and Microarchitecture

The ISA is the interface between software and hardware ("what the programmer sees").

The microarchitecture specifies the underlying implementation that actually executes the instructions.

## Blocking and Non-Blocking Assignment Guidelines

- Use  $\delta$  (posedge clk) and non-blocking (=) to model synchronous sequential logic
- Use continuous assignments to model simple combinational logic
- Use always  $\delta(*)$  and blocking (=) assignments to model more complicated combinational logic where the always statement is helpful
- Do not make assignments to the same signal in more than one always statement or continuous assignment.

- mechanism to cache in a system, call in the QD  
- number of groups - program register  
- index of each bit in programmatic branch predictor configuration register

## ISA

## vs. Microarchitecture

- number of bits required for destination register of a load instruction
- Instructions: opcodes, addressing modes, data types, instruction type and format, registers, condition codes
- Memory: address space, alignment, addressability, virtual memory management, multiprocessor memory location
- Call, interrupt and exception handling
- I/O: memory mapped vs. instructions
- Power & Thermal management
- Multiprocessing / Multithreading support
- Access control, priority and privilege
- Memory-mapped location of exception vectors
- Function of each bit in a programmable branch prediction register
- Order of execution of loads and stores in multi-core CPU
- Program counter width
- Hardware FP-exception support
- Vector instruction support
- CPU endianness
- Virtual page size

## Performance Evaluation

- CPI: cycles per instruction
  - IPC: instruction per cycle
  - MHz: frequency,  $10^6$  cycles/sec.
  - higher MHz  $\Rightarrow$  higher MIPS, IPS could be lower
  - higher MIPS  $\Rightarrow$  less time, could need more instructions
- MIPS: million instructions / sec. = MHz/CP
- Time = # instr.  $\cdot$  CPI  $\cdot$   $\frac{1}{\text{Hz}}$
- Speedup = oldTime/newTime

## Single-Cycle Machines

Each instruction takes a single clock cycle and all state updates are made at the end of the cycle.

- slowest instruction determines cycle time
- + easy to build

## Multi-Cycle Machines

Instruction processing is broken into multiple stages/cycles, state updates happen during execution and architectural updates at the end. Instruction processing consists of two components:

- Data path - relay and transform data
- Control logic - FSM that determines control signals
- + slowest stage determines cycle time

## Dataflow

In a dataflow machine, a program consists of dataflow nodes. A node fires (executes) when all its inputs are ready.

## Write-through

data written to cache block is simultaneously written to main memory

## Write-back

dirty bit (1) is associated with each cache block. write back to main memory when evicted from cache

## Pipelining

The idea is to process multiple instructions at once by keeping each stage occupied. In reality there are a few problems:

- Resource contention, can be fixed by duplication, increased throughput or detection and stalling
- Long latency operations
- Data dependencies, there are flow (read after write), output (write after write) and anti (write after read) dependencies. The last two exist due to a limited amount of registers.

## Handling flow dependencies

- stall
- eliminate at software level
- predict values
- data forwarding
  - $\hookrightarrow W \rightarrow D$ : internal/register file forwarding
  - $\hookrightarrow M \rightarrow E$ : operand forwarding

## Pipeline Stages

- Fetch: CPU reads instructions from instruction memory
- Decode: CPU reads source operands from register file and decodes instruction to control signals
- Execute: CPU performs a computation with the ALU
- Memory: CPU reads/writes data memory
- Writeback: CPU writes result to register file

## Interlocking & Scoreboarding

Detection of data dependencies to ensure correct execution.

## Out-of-Order Execution

Idea to move dependent instructions out of the way of independent ones. Reservation stage as rest are for dependent instructions.

## Reorder Buffer

Complete instructions OoO but reorder them before making results visible to architectural state.

## Tomasulo's Algorithm

Implementation of OoO-Execution. Uses register renaming to eliminate output and anti-dependencies. It further uses reservation stations for individual operations.

1. If reservation station is available:
  - instr. + renamed operands inserted into reservation station
  - rename destination registerElse stall
2. While in reservation station:
  - watch common data bus for tag of sources
  - if tag seen grab value
  - if both operands are valid inst. ready for dispatch
3. Dispatch instr. to functional unit
4. After instr. finishes:
  - put tagged value onto common data bus
  - if register file contains tag, update its value and set valid bit
  - reclaim rename tag  $\rightarrow$  no valid copy of tag in the system

## VLIW

The idea is that the compiler finds independent instructions and statically schedules them into single VLIW instructions.

Lock step execution: if one instruction stalls, the whole VLIW stalls

- + simple hardware
- compiler needs to find  $N$  independent instructions per cycle
- + no dependency checking
- + no instruction distribution
- lock step causes stalls

## Superscalar Execution

Idea is to fetch/decode/... multiple instructions per cycle.

- + higher IPC
- higher complexity for dependency checking  $\Rightarrow$  more hardware

## Systolic Arrays

Instead of a single processing element (PE) we have a array of PE and carefully orchestrate the data-flow between them  $\Rightarrow$  Maximize computation done on a single element.

Difference from pipelining: Array structure is non-linear and multi-dimensional. PE connections can be multi-directional with different speeds. PEs can have local memory and execute kernels.

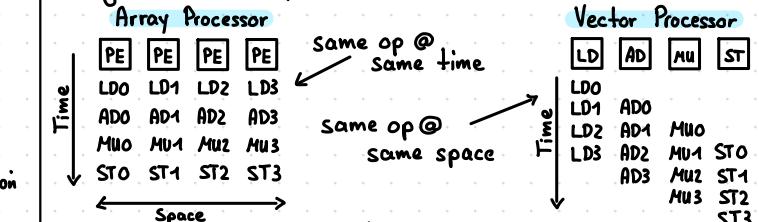
## Fine Grained Multithreading

Hardware has multiple thread contexts (PC+reg) and each cycle the fetch engine fetches from a different thread.

- + no dependencies
- + no branch prediction
- + improved throughput, latency, tolerance, utilization
  - $\hookrightarrow$  dependency checking between threads
- extra hardware
- reduced single-thread performance
- resource contention

## SIMD

Single instruction operates on multiple data.



## Formulas

$$\text{Execution time} = (\# \text{ instructions}) / (\text{instructions}/\text{cycle}) \cdot \frac{\text{cycles}}{\text{seconds}}$$

MIPS: million instructions per sec.  
 $\frac{1}{\text{cycle}} \cdot \frac{\text{cycles}}{\text{seconds}}$

$$\text{Miss Rate} = \frac{\text{Number Of Misses}}{\text{Number Of Total Memory Accesses}} = 1 - \text{Hit Rate}$$

$$\text{Average Memory Access Time} = t_{\text{cache}} + \text{MR cache} (t_{\text{main}} + \text{MR main})$$

$t_{\text{cache}}$  = Main Memory,  $t_{\text{main}}$  = Virtual Memory, MR = Miss Rate

Number of set bits:  $\log_2 S$

Remaining tag bits indicate memory address of the data stored in a given cache set. Two least significant bits  $\rightarrow$  byte offset

## Vector Processing

Performs operation on a whole array. This is only possible if the operations on each element are independent from each other.

The data gets stored in vector registers. Vector chaining describes the vector version of data forwarding. It allows a operation to start as soon as an individual element is ready. The stride is the distance of the vector elements in memory. If a vector is too long it can be split into multiple vectors (strip mining).

- + a lot of work per instruction - works only if parallelism is regular, else it is very inefficient
- + regular memory access pattern
- + no need for loops

## GPU

GPUs are SIMD engines but programmed using threads (SPMD). A set of threads executing the same program are grouped into a warp.

Dynamic Warp merging: Merge threads executing the same instr. after branch divergence. This forms new warps from the warps waiting.

## Delayed Branching

Means that some instructions after a branch are executed regardless of which way the branch goes. A compiler can instructions in such a delay slot if they don't influence the branch itself, else they are filled with NOPs.

## Branch Prediction

A technique used to predict the next address after a branch. If the prediction is wrong, we have to flush the pipeline (misprediction penalty).

## Prediction Direction Schemes

- always (not)-taken (30-40%) 60-70% accuracy
- BTBN: backwards taken forwards not taken (good with loops)
- Last time predictor: single bit stored in BTB (branch target buffer) indicates last direction. Loop accuracy =  $\frac{n-2}{n}$ .
- 2-bit counter based prediction:
  - Local: no interference between different branches
  - Global: single counter for all branches

## Cycles

Total number of cycles can be expressed as  $C = P(I - 1) + B + D$

- C = total number of cycles taken
- P = total number of pipeline stages
- I = total number of instructions
- B = total number of conditional branch instructions executed (instructions executed)
- D = number of cycles stalled for each conditional branch

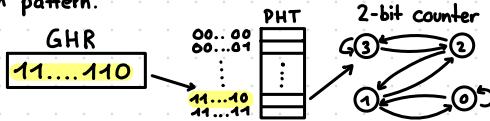
## Delays

Propagation delay: tpd: minimum time from when an input changes until its output (G) reaches their final value.  
Computation delay: tcd: minimum time from when an input changes until any output starts to change its value.

## Global Branch Correlation

The idea is that recently executed branch outcomes are correlated with the outcome of the next branch.

- First level: Global branch history register, keeps track of the last branch outcomes.
- Second level: Pattern history table, keeps a 2-bit counter for each pattern.



## Memory Hierarchy

Memory Array: stores data, address selection logic selects row, readout circuitry reads data.

Memory banking: Multiple memory units with a common data and address bus, helps to resolve long latency. Units can be accessed individually.

Locality: temporal = access to same address in short time  
spatial = access to nearby address

## Blocks & Addressing cache:

- memory is divided into fixed-size blocks
- each block maps to a location in the cache (index bits)
- for a cache hit the tags need to match



Offset = Byte im Block  
Index = Zeile in Tag Store  
Tag = Welches Block im Set

## Associativity:

- multiple blocks have the same index → conflict misses
- n-way associative allows n-blocks with same index
- 1-way → direct mapped      no index → fully associative

## Cache Performance:

- cache size, total data  $c$
- block size  $b$
- associativity  $n$
- #blocks:  $B = c/b$
- #sets:  $S = B/n$

## Replacement Policies:

- FIFO, first-in-first-out
- LRU, least-recently-used
- Random

## Handling Writes

Writeback: write to lower levels when the block is evicted, needs a dirty bit.

Writethrough: write to all levels immediately, simpler but bandwidth intensive.

## Classification of Misses

Compulsory Miss: first reference is always a miss (prefetching)

Capacity Miss: cache is too small

Conflict Miss: all other misses (more associativity)

## Improvement Ideas

- reduce miss rate
- reduce miss latency or cost
- reduce hit latency or cost

$$\text{bits of storage} = \text{associativity} \cdot (\text{tag} + \text{dirty} + \text{valid}) + \text{LRU}$$

## Prefetching

The idea is to improve cache performance by preloading data to avoid misses. There are different techniques to prefetching, some are software based while others hardware dependent.

Stride prefetcher: prefetches cache block in a pattern with a certain stride (if stride = 0, next block prefetching)

Runahead execution: allows the processor to pre-process instr. during cache misses instead of stalling. Therefore it can detect potential cache misses earlier.

- accuracy = #pref. used / #pref. total
- coverage = #acc. predicted / #total acc.

## Virtual Memory

Much larger than physical memory. Virtual address space is divided in pages, while physical address space is divided into frames. Page Table stores mapping: Virtual → Physical together with a valid bit (and more meta data).

$$\# \text{Virtual Pages} = \frac{\text{Virtual Address}}{\text{Page Size}}$$

$$\# \text{Physical Pages} = \frac{\text{Physical Address}}{\text{Page Size}}$$

$$\text{VA: } \frac{\text{Virtual Page Number}}{\text{Page Offset}}$$

$$\text{PA: } \frac{\text{Page Table}}{\text{P. Page Number}} \downarrow \frac{\text{Page Offset}}{\text{Page Offset}}$$

$$\text{Physical Address Space} = \frac{\text{PPN}}{\text{Page Offset}}$$

Multi-level page tables: keeps PT size small

Memory protection: different PT for each program

Translation Lookaside Buffer TLB: cache PT entries to speed up address translation

A memory address  $a$  is valid to one digit when  $a$  is a multiple of  $n$  (where  $n$  is a power of 2). In this case,  $a$  is the smallest unit of memory access, i.e. each memory address specifies a different byte. An n-digit address would have  $n$  log<sub>2</sub>(n) least significant bits that implement address masking.

Cache capacity = (Block Size in Bytes) \* (Blocks per Set) \* (Number of Sets)

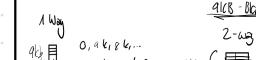
Index Bits = log<sub>2</sub>(Blocks per Set)

Block Offset Bits = log<sub>2</sub>(Block Size in Bytes)

Tag Bits = (Address Bits) - (Index Bits) - (Block Offset Bits)

Software Interlocking: NOPs

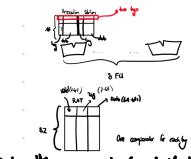
Hardware Interlocking: Detects dependencies and stalls pipeline accordingly



Utilization:  $(\text{Threads} - \text{Instructions each thread is taking}) / (\text{Threads} + \text{Total Instructions taken by all other threads})$

Worst:  $(\text{Total Threads}) / (\text{Threads per core})$

Min Utilization: following the same path (except if there are invalid threads then you can't use many instructions in the full ways as possible and the load in the threads of the failed paths)



Ex. Find the simplest sum-of-products form for this equation:  $F = B + (A + \bar{C}) \cdot (\bar{A} + \bar{B} + \bar{C})$

$$\begin{aligned} F &= B + A\bar{A} + A\bar{B} + A\bar{C} + \bar{C}\bar{A} + \bar{C}\bar{B} + \bar{C} \\ &= B + A + \bar{C} \end{aligned}$$

Ex. Simplify the following min-terms:  $\sum(3, 5, 7, 11, 13, 15)$ .  $\{3, 5, 7, 11, 13, 15\} = \{0011, 0101, 0111, 1011, 1101, 1111\}$

$$\begin{aligned} F &= \bar{A}\bar{B}\bar{C}\bar{D} + \bar{A}\bar{B}\bar{C}\bar{D} + \bar{A}\bar{B}\bar{C}D + \bar{A}\bar{B}CD + \bar{A}BC\bar{D} + ABCD \\ &= CD \cdot (\bar{A}\bar{B} + \bar{A}\bar{B} + \bar{A}\bar{B} + AB) + BD \cdot (\bar{A}\bar{C} + AC) \\ &= CD + BD\bar{C} \\ &= D(B + C) \end{aligned}$$

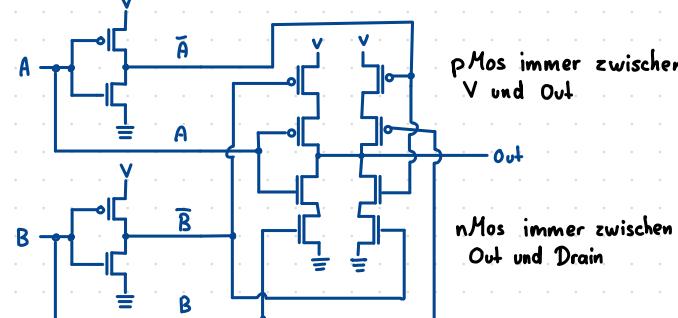
Ex. Convert the following equation to only contain NANDs.

$$\begin{aligned} F &= (\bar{A}\bar{B} + C) + AC = (\bar{A}\bar{B} + C) \cdot \bar{AC} = (\bar{A}\bar{B} \cdot \bar{C}) \cdot \bar{AC} \\ &= AB \cdot \bar{AC} = \overline{\overline{AB} \cdot \overline{AC}} \end{aligned}$$

Ex. Convert the following equation to only contain NANDs.

$$\begin{aligned} F &= (\bar{A} + BC) + \bar{C} = (\bar{A} + BC) \cdot \bar{C} = (\bar{A} + BC) \cdot C \\ &= (\bar{A} \cdot (BC)) \cdot C = (\bar{A} \cdot A) \cdot (B \cdot C) \cdot C \end{aligned}$$

Ex. Draw a XOR-Gate with transistors.



Ex. Sequential or Combinational circuit?

```
module one (input clk, input o, input b, output reg [1:0] q);
  always @(*)
    if (b)
      q <= 2'b01;
    else if (a)
      q <= 2'b10;
endmodule
```

This code results in a sequential circuit because a latch is required to store old values of q if both conditions are not satisfied.

Ex. Is this code a correct multiplexer?

```
module four (input sel, input [1:0] data, output reg z);
  always @(*)
    if (sel)
      z = data[1];
    else
      z = data[0];
endmodule
```

No, the input data is missing in the sensitivity list. A update would not be reflected to the output z.

Ex. Does this result in a D-FlipFlop with a synchronous active-low reset?

```
module mem (input clk, input reset, input [1:0] d, output reg [1:0] q);
  always @(*)
    if (!reset) q <= 0;
    else q <= d;
endmodule
```

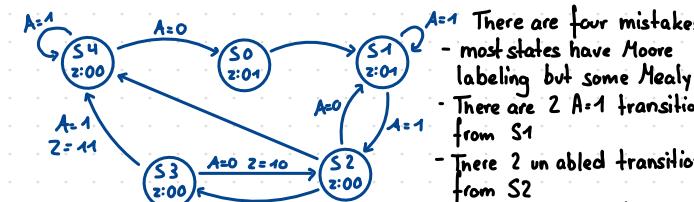
The code implements 2 D-FlipFlops, each works with a asynchronous active low reset.

Ex. Is this code syntactically correct?

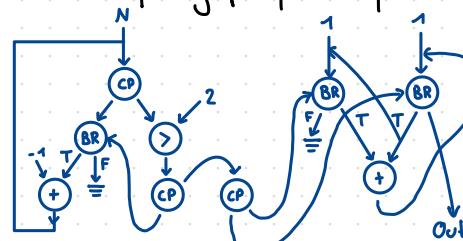
```
module fulladd (input a, b, c, output reg s, c_out);
  assign s = a ^ b;
  assign c_out = (a & b) | (a & c) | (b & c); we use assign there-fore these have to be wires
endmodule
```

```
module top (input wire [5:0] instr, input wire op, output z);
  reg [1:0] r1, r2; should be wires
  wire [3:0] w1, w2;
  fulladd FA1 (.a(instr[0]), .b(instr[1]), .c(instr[2]),
                .c_out(r1[1]), .z(r1[0]));
  fulladd FA2 (.a(instr[3]), .b(instr[4]), .c(instr[5]),
                .c_out(r2[1]), .z(r2[0]));
  assign z = r1 | op;
  assign w1 = r1 + 1;
  assign w2 = r2 << 1;
  assign op = r1 ^ r2; multiple drivers
endmodule
```

Ex. List all the mistakes in this diagram.



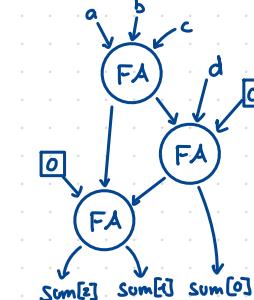
Ex. Draw a data flow graph for the fibonacci function.



Ex. Which designs are compatible with each other?

superscalar - in-order precise exceptions - out-of-order retirement  
superscalar - out-of-order branch prediction - fine-grained multithreading  
single cycle - branch prediction fine-grained multithreading - pipelining  
reservation station - microprogramming Tomasulo's algorithm - in-order  
fine-grained multithreading - single core direct mapped cache - LRU replacement

Ex. Draw the dataflow graph for a four 1-bit addition, you can use Full Adder nodes.



Ex. Any  $n \geq 3$  1-bit addition can be implemented only using Full Adders. Fill out the table.

n	# required FAs	n	# required FAs
3	1	6	4
4	3	7	4
5	3	8	7

Ex. Two programs A, B run on the same machine, both have the same # memory requests, but A needs to stall way more. Why could this be?

A could have a lot of row buffer conflicts, while B has a lot of row buffer hits.

Ex. If a processor executes more IPS, does a program always finish faster?

No, the number of instructions for a program could be different for different processors.

Ex. If a program runs on a processor with a higher frequency, does this imply that it executes more IPS?

No, a processor with a lower frequency can have a much higher number of IPC.

Ex. Write a MIPS 64-bit subtraction (2s-complement) where  $\$4 \$5 - \$6 \$7$ .

```
subu $3, $5, $7
sltu $2, $5, $7
add $2, $6, $2
sub $2, $4, $2
```

Ex. A machine with 5 pipeline stages uses delay slots to handle control dependences. Jump and branch are resolved during execution stage. How many delay slots are needed?

2, since we can fill them during fetch and decode of the jump / branch instruction.

Can we modify the pipeline to reduce the number of delay slots?

Yes, if we move the resolution of the jump/branch target to the decode stage, we only need one delay slot.

Ex. How many delay slots are needed for the following implementations?  
 In-order with branch resolving during 4th stage: 3  
 000 with 64 reservation stages, branch resolving during 2nd cycle of branch execution and 16 stages before the execution stage: Don't know

Ex. Given the following microbenchmark for a pipelined machine.  
 Calculate #dynamic instructions executed, # pipeline stages and #cycles of stall caused by branch instruction.

LOOP1:              Initial R1    #Cycles

SUB R1, R1, #1	4	51
BGT R1, LOOP1	8	63
	16	87

LOOP2:              B LOOP2    all runs execute the same #dynamic instr.

Let:  $C = \# \text{cycles}$

$P = \# \text{stages}$

$I = \# \text{dynamic instr.}$

$B = \# \text{branch instr.}$

$D = \# \text{cycles stall / branch} \Rightarrow P + I = 40, D = 3$

$$C = P + T - 1 + B \cdot D$$

$$51 = P + I - 1 + 4D$$

$$63 = P + I - 1 + 8D$$

$$87 = P + I - 1 + 16D$$

Ex. Given a scalar processor with in-order fetch, out-of-order dispatch and in-order retirement. It has 4 pipeline stages, and 2 reservation stations (one for each type). If the following program gets executed, answer the questions?

```

MOV R0<-1000 FDE1E2E3E4W
LD R1<[R0] FD - - E1E2E3E4E5E6E7E8W
BL R1,100,LB1 FD - - - - E1E2E3E4W
MUL R1<R1,5 FD E1E2E3 // killed
ST [R0]< R1 FDE1E2E3E4W
ADD R1<R1,R0 FD - - - E1W
ST [R0]< R1

```

Cache hit latency? 1 cycle, the last ST instr. is a hit.

Cache miss latency? 8 cycles, the first LD instr. misses.

Cache line size? Unknown

#entries in each reservation station? ALU at least 2, MU unknown

#ALUs? if pipelined at least 1, else at least 2.

Is the ALU pipelined? If there is only 1 ALU yes, else unknown

Does the processor have branch prediction? Yes, because there

are instr. that get killed.

At which stage do branches get resolved? At the end of E4, because in the next cycle the previously fetched instr. get killed.

Ex. Given Tomasulo's Algorithm with: 8 functional units with their own tag/data bus, 32x64 bit registers, 16 reservation stations per functional unit and 2 source register per reservation station, calculate:

$$\# \text{tag comparator / reservation station entry} = 2 \times 8 = 16$$

$$\# \text{tag comparators} = 16 \times 16 \times 8 + 32 \times 8 = 2304$$

$$\min. \text{tag size} = \log(16 \times 8) = 7$$

$$\min. \text{size of register alias table} = 32 \times (7 + 64 + 1) = 2304$$

$$\min. \text{total size of tag store} = 32 \times 7 + 8 \times 16 \times 2 \times 7 = 2048$$

Ex. Comparing a VLIW and an in-order superscalar processor with the same machine width and frequency. For a program A, the VLIW machine is much faster, why could this be?

The superscalar proc. is in-order, requiring bubbles in the pipeline, while the VLIW proc. can re-order instr.

For some other program B, the VLIW is slower, why could this be?

VLIW needs NOPs, while the superscalar proc. doesn't. These NOPs can lead to lower L1-cache hit rate and higher fetch bandwidth.

Ex. Which of the following are goals of VLIW?

- i. Simplify code compilation
- ii. Simplify application development
- iii. Reduce overall hardware complexity
- iv. Simplify hardware dependence checking
- v. Reduce processor fetch width

Ex. Given a vector proc. with these fully interleaved / pipelined instr. VLD/VST 50 cycles, VADD 4 cycles, VMUL 16 cycles, VDIV 32 cycles and VRSHFA 1 cycle. Assume: in-order pipeline, chaining between functional units, first element bank 0, 8KB row buffer / bank, 64 bit vector elements, each memory bank has 2 ports and there are 2 load/store units. What is the minimum (power of 2) # banks so memory accesses never stall? 64 banks, because access latency is 50ms and 64 is the next power of 2.

Executing this program takes 111 cycles, what is the vector length?

VLD	V1, A	50	L-1	111 = 51 + 4 + 16 + 1 + L-1
VLD	V2, B	50	L-1	=> L = 40
VADD	V3, V1, V2	4	L-1	
VMUL	V4, V1, V3	16	L-1	
VRSHFA	V5, V4, 2	16	L-1	

Reducing the banks by a factor of 2, how long does the program take?

VLD	[0]	50		
	[31]	50		
	[32]	50		
	[33]	50		
VLD	[0]	50		
	[31]	50		
	[32]	50		
	[33]	50		

$1 + 50 + 7 + 50 + 4 + 16 + 1 = 129 \Rightarrow 129 \text{ cycles}$

VADD                          4                          tracking the last element

VMUL                          16                          tracking the last element

VRSHFA                          1                          tracking the last element

Now the #banks get reduced further and it takes 279 cycles. How many banks are there?

$$279 = 1 + 16 + 4 + 1 + 7 + \lceil \frac{1}{4} \rceil \cdot 50$$

$$\Rightarrow 5 = \lceil \frac{1}{4} \rceil \cdot 7 \Rightarrow x = 8 \text{ memory banks}$$

In a new version 4 vector proc. share the same memory with 4 times the banks. However the execution is slower than if each program ran on a single proc. with 1/4 banks, why could this be?

Row buffer conflicts as all cores interleave their vectors across all banks.

How can this be fixed?

Partition the memory mappings, or use better memory scheduling.

Ex. Consider the following warps, how can dynamic warp formation be used?

$$X = \{1\ 0\ 0\ 1\ 0\ 1\ 1\ 1\}$$

$$Y = \{1\ 0\ 0\ 0\ 1\ 0\ 0\ 1\}$$

$$Z = \{0\ 1\ 0\ 0\ 0\ 0\ 0\}$$

X' = {1 0 0 1 0 1 1 1}  
 Y' = {1 1 0 0 1 0 0 1}  
 Z' = {0 0 0 0 0 0 0 0}

There are several answers, but notice that X, Y can't be merged.

Ex. How effective is a 16KB, 4-way associative cache with 8B instructions?

Not effective, since the block size is 4B, each instruction needs two accesses. Further it can't exploit spatial locality.

Ex. Given the following access pattern and hit rate for a cache determine its characteristics.

Addresses Accessed	must miss	Hit rate
1. 0 4 8 16 64 128		1/2
2. 31 8192 63 16384 4096 8192 64 16384		5/8
3. 32368 0 123 1024 3072 8192 1	hit	1/3

Cache block size: 8, 16, 32, 64 or 128B

From ① we can see that only 32 or 64 are possible. From ② we can see that 63 must be a hit and therefore it can only be 64B.

Cache Associativity: 1, 2, 4 or 8 way

Combining this with the possible cache sizes of 4 or 8KB we can see that 1 and 2 way would cause too much misses in ② and 8 way would cause another miss in ③, therefore it must be 4 way.

Cache size: 4 or 8 KB

In ③ the access to 0 is a miss and therefore 8192 should be a hit, but with 4KB, 1024 and 3072 would map to the same set and therefore it couldn't be a hit. So cache size must be 8KB.

Replacement Policy: LRU or FIFO

For 8192 to hit in ③ it must be LRU.

Ex. Given a one level cache with 128 B and block size 32 B. Using LRU, the following blocks are accessed:

ABAHBGHHAEHDHGCGCABHDECCBADEF

In a direct mapped cache which blocks are in the same set?

A/B H/D G/C E/F

For a fully associative cache, write down the misses.

ABAHBGHHAEHDHGCGCABHDECCBADEF

Ex. Given a 2-way assoc. write back cache with LRU and a  $2^9 \times 15$  bit tag store. It is virtually indexed, physically tagged. The virtual address space is 1MB, page size 2KB and block size 8B. What is the size of the data store?

$$\text{Tag store} = 2^9 \cdot (2^4 + 5) = 15 \cdot 2^9 \Rightarrow i = 9 \ t = 5$$

$$\text{Data store} = 2^9 \cdot 8 \cdot 2^3 = 8\text{KB}$$

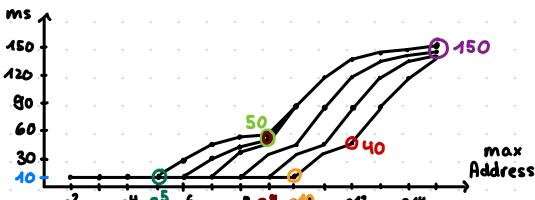
How many bits of the virtual index come from the VPN?

Index	block
9	3

What is the physical address space?

$$\begin{aligned} \text{Page offset} &= 11 \text{ Bits} & \text{Page Tag} &= 5 \text{ Bits} \\ \Rightarrow 2^{11} \cdot 2^5 &= 2^{16} = 64\text{KB} \end{aligned}$$

Ex. Fill in the blanks?



	L1	L2	L3	DRAM
line size	N/A	X	X	N/A
cache assoc.	$2^{10}/1024 = 1$	$2^{10}/1024 = 4$	X	X
cache size	$2^5 = 32$	$2^5 = 512$	X	X
access latency	10ms	40ms	X	$= 100\text{ms}$

"max stride"

Ex. Give a 4-way assoc. write back cache with a  $2^{11} \times 89$  bit tag store, a 9 bit replacement policy, 64B blocks. It is virtually indexed, physically tagged and data from a given adr. can be in up to 8 sets. It uses a 2 level page table with each 1024 entries. How many bits of the virtual address are for the set?

$$\text{Tag Store} = 2^{11} \cdot 89 \rightarrow 11 \text{ Bits}$$

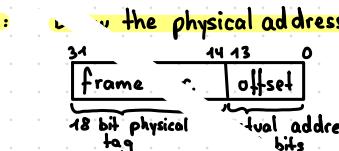
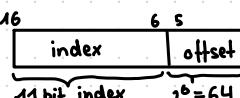
What is the size of the data store?

$$2^{11} \cdot 4 \cdot 64\text{B} = 512\text{KB}$$

How many bits in the PPN overlap with the index bits in the virtual address?

3, since data can be present in up to  $2^{3-8}$  sets.

Draw the virtual address:



$$\text{What is the page size? } 2^{14} = 16\text{KB}$$

What is the virtual address space?

$$2^{\text{VPN}} \cdot \text{Page Offset} = 2^{20} \cdot 2^{14} = 1024 \cdot 2^{14}$$

$$\text{What is the physical address space? } 2^{32} = 4\text{GB}$$

Ex. What is the prefetch accuracy and coverage for A, B using a stride prefetcher?

A: int[100] a;

```
sum = 0;
for (i = 0; i < 1000; i += 4)
    sum += a[i]
```

$$\text{Accuracy: A: } \frac{248}{249} \text{ B: } 0$$

B: int[100] a;

```
sum = 0;
for (i = 0; i < 1000; i += 4)
    sum += a[i]
```

$$\text{Coverage: A: } \frac{248}{250} \text{ B: } 0$$

Ex. Given this code explain which branches correlate locally/globally.

```
for (int i=0; i<N; i++) { //B1
    val = array[i];
    if (val%2==0) //B2
        sum += val
    if (val%3==0) //B3
        sum += val
    if (val%6==0) //B4
        sum += val
```

Locally: only B1, since for B2, B3, B4 the previous value does not matter.

Globally: B4 is correlated with B3 and B2. If one B4 is taken, B2 and B3 are also taken.

Ex. For the same code, calculate the expected value for the PHTE taken-taken after 120 iterations.

W.l.o.g. we take a look at the numbers 1-6. For a single iteration we have 4 chances to increment the PHTE.

B3: Given  $\Pr[B_1.T \wedge B_2.T] = 1/2$  the probability for B3 to be taken is  $1/3$ , resulting in  $1/2 \cdot 1/3 = 1/6$  probability to increase and  $1/2 \cdot (1 - 1/3) = 1/3$  to decrease, therefore B3 contributes  $1/6 - 1/3 = -1/6$ .

B4:  $\Pr[B_2.T \wedge B_3.T] = 1/6 \Rightarrow 1/6 \cdot 1 = 1/6$

B1:  $\Pr[B_3.T \wedge B_4.T] = 1/6 \Rightarrow 1/6 \cdot 1 = 1/6$

B2:  $\Pr[B_4.T \wedge B_1.T] = 1/6 \Rightarrow 1/6 \cdot 1/2 = 1/12 = 0$

Resulting in a total of  $1/6$  per iteration. Therefore after 120 iterations, the expected value is 20.

Fibonacci Code in Assembly

```
int fibfib(n) {
    int a=0;
    int b=1;
    int c=a+b;
    while (n>2) {
        c=a+b;
        a=b;
        b=c;
        n--;
    }
    return c;
}
```

```
branch:
    blt d0, g // use 13 as loop
    add d0, d0, d0 // 2x d0
    add d0, d0, d0 // 4x d0
    add d0, d0, d0 // 8x d0
    add d0, d0, d0 // 16x d0
    add d0, d0, d0 // 32x d0
    add d0, d0, d0 // 64x d0
    add d0, d0, d0 // 128x d0
    add d0, d0, d0 // 256x d0
    add d0, d0, d0 // 512x d0
    add d0, d0, d0 // 1024x d0
    add d0, d0, d0 // 2048x d0
    add d0, d0, d0 // 4096x d0
    add d0, d0, d0 // 8192x d0
    add d0, d0, d0 // 16384x d0
    add d0, d0, d0 // 32768x d0
    add d0, d0, d0 // 65536x d0
    add d0, d0, d0 // 131072x d0
    add d0, d0, d0 // 262144x d0
    add d0, d0, d0 // 524288x d0
    add d0, d0, d0 // 1048576x d0
    add d0, d0, d0 // 2097152x d0
    add d0, d0, d0 // 4194304x d0
    add d0, d0, d0 // 8388608x d0
    add d0, d0, d0 // 16777216x d0
    add d0, d0, d0 // 33554432x d0
    add d0, d0, d0 // 67108864x d0
    add d0, d0, d0 // 134217728x d0
    add d0, d0, d0 // 268435456x d0
    add d0, d0, d0 // 536870912x d0
    add d0, d0, d0 // 1073741824x d0
    add d0, d0, d0 // 2147483648x d0
    add d0, d0, d0 // 4294967296x d0
    add d0, d0, d0 // 8589934592x d0
    add d0, d0, d0 // 17179869184x d0
    add d0, d0, d0 // 34359738368x d0
    add d0, d0, d0 // 68719476736x d0
    add d0, d0, d0 // 137438953472x d0
    add d0, d0, d0 // 274877906944x d0
    add d0, d0, d0 // 549755813888x d0
    add d0, d0, d0 // 1099511627776x d0
    add d0, d0, d0 // 2199023255552x d0
    add d0, d0, d0 // 4398046511104x d0
    add d0, d0, d0 // 8796093022208x d0
    add d0, d0, d0 // 17592186044416x d0
    add d0, d0, d0 // 35184372088832x d0
    add d0, d0, d0 // 70368744177664x d0
    add d0, d0, d0 // 140737488355328x d0
    add d0, d0, d0 // 281474976710656x d0
    add d0, d0, d0 // 562949953421312x d0
    add d0, d0, d0 // 1125899906842624x d0
    add d0, d0, d0 // 2251799813685248x d0
    add d0, d0, d0 // 4503599627370496x d0
    add d0, d0, d0 // 9007199254740992x d0
    add d0, d0, d0 // 18014398509481984x d0
    add d0, d0, d0 // 36028797018963968x d0
    add d0, d0, d0 // 72057594037927936x d0
    add d0, d0, d0 // 144115188075855872x d0
    add d0, d0, d0 // 288230376151711744x d0
    add d0, d0, d0 // 576460752303423488x d0
    add d0, d0, d0 // 1152921504606846976x d0
    add d0, d0, d0 // 2305843009213693952x d0
    add d0, d0, d0 // 4611686018427387904x d0
    add d0, d0, d0 // 9223372036854775808x d0
    add d0, d0, d0 // 18446744073709551616x d0
    add d0, d0, d0 // 36893488147419103232x d0
    add d0, d0, d0 // 73786976294838206464x d0
    add d0, d0, d0 // 147573952589676412928x d0
    add d0, d0, d0 // 295147905179352825856x d0
    add d0, d0, d0 // 590295810358705651712x d0
    add d0, d0, d0 // 1180591620717411303424x d0
    add d0, d0, d0 // 2361183241434822606848x d0
    add d0, d0, d0 // 4722366482869645213696x d0
    add d0, d0, d0 // 9444732965739290427392x d0
    add d0, d0, d0 // 18889465931478580854784x d0
    add d0, d0, d0 // 37778931862957161689568x d0
    add d0, d0, d0 // 75557863725914323379136x d0
    add d0, d0, d0 // 151115727451826646758272x d0
    add d0, d0, d0 // 302231454903653293516544x d0
    add d0, d0, d0 // 604462909807306587033088x d0
    add d0, d0, d0 // 1208925819614613174066176x d0
    add d0, d0, d0 // 2417851639229226348132352x d0
    add d0, d0, d0 // 4835703278458452696264704x d0
    add d0, d0, d0 // 9671406556916905392529408x d0
    add d0, d0, d0 // 19342813113833810785058816x d0
    add d0, d0, d0 // 38685626227667621570117632x d0
    add d0, d0, d0 // 77371252455335243140235264x d0
    add d0, d0, d0 // 15474250491067048628047052x d0
    add d0, d0, d0 // 30948500982134097256094104x d0
    add d0, d0, d0 // 61897001964268194512188208x d0
    add d0, d0, d0 // 123794003928536389024376416x d0
    add d0, d0, d0 // 247588007857072778048752832x d0
    add d0, d0, d0 // 495176015714145556097505664x d0
    add d0, d0, d0 // 990352031428291112195011328x d0
    add d0, d0, d0 // 1980704062856582224387522656x d0
    add d0, d0, d0 // 3961408125713164448775045312x d0
    add d0, d0, d0 // 7922816251426328897550090624x d0
    add d0, d0, d0 // 15845632522852657795100181248x d0
    add d0, d0, d0 // 31691265045705315590200362496x d0
    add d0, d0, d0 // 63382530091410631180400724992x d0
    add d0, d0, d0 // 12676506018282126236080144992x d0
    add d0, d0, d0 // 25353012036564252472160289984x d0
    add d0, d0, d0 // 50706024073128504944320579968x d0
    add d0, d0, d0 // 101412048146257009888641159936x d0
    add d0, d0, d0 // 202824096292514019777282319872x d0
    add d0, d0, d0 // 405648192585028039554564639744x d0
    add d0, d0, d0 // 811296385170056079109129279488x d0
    add d0, d0, d0 // 1622592770340112158218245558976x d0
    add d0, d0, d0 // 3245185540680224316436491117952x d0
    add d0, d0, d0 // 6490371081360448632872982235904x d0
    add d0, d0, d0 // 12980742162720897265745764471808x d0
    add d0, d0, d0 // 25961484325441794531491528943616x d0
    add d0, d0, d0 // 51922968650883589062983057887232x d0
    add d0, d0, d0 // 103845937301767778125966115754464x d0
    add d0, d0, d0 // 207691874603535556251932231508928x d0
    add d0, d0, d0 // 415383749207071112503864463017856x d0
    add d0, d0, d0 // 830767498414142225007728926035712x d0
    add d0, d0, d0 // 166153499682828445001545785207144x d0
    add d0, d0, d0 // 332306999365656890003091570414288x d0
    add d0, d0, d0 // 664613998731313780006183140828576x d0
    add d0, d0, d0 // 132922799746262756001236628165752x d0
    add d0, d0, d0 // 265845599492525512002473256321504x d0
    add d0, d0, d0 // 531691198985051024004946512643008x d0
    add d0, d0, d0 // 1063382397970102048009893025286016x d0
    add d0, d0, d0 // 2126764795940204096019786050572032x d0
    add d0, d0, d0 // 4253529591880408192039572101144064x d0
    add d0, d0, d0 // 8507059183760816384079144202288128x d0
    add d0, d0, d0 // 17014118367521632768158288404576256x d0
    add d0, d0, d0 // 34028236735043265536316576809152512x d0
    add d0, d0, d0 // 68056473470086531072633153618305024x d0
    add d0, d0, d0 // 136112946940173062145266307236610048x d0
    add d0, d0, d0 // 272225893880346124290532614473220096x d0
    add d0, d0, d0 // 544451787760692248581065228946440192x d0
    add d0, d0, d0 // 1088903575521384497162130457892880384x d0
    add d0, d0, d0 // 2177807151042768994324260915785760768x d0
    add d0, d0, d0 // 4355614302085537988648521831571521536x d0
    add d0, d0, d0 // 8711228604171075977297043663143043072x d0
    add d0, d0, d0 // 17422457208342151955594087326285661544x d0
    add d0, d0, d0 // 34844914416684303911188174652571323088x d0
    add d0, d0, d0 // 69689828833368607822376349305142666176x d0
    add d0, d0, d0 // 13937965766673721564475269861029332352x d0
    add d0, d0, d0 // 27875931533347443128950539722058665104x d0
    add d0, d0, d0 // 55751863066694886257901079444117330208x d0
    add d0, d0, d0 // 111503726133389772515802188888234660416x d0
    add d0, d0, d0 // 223007452266779545031604377776469320832x d0
    add d0, d0, d0 // 446014904533559090063208755552938641664x d0
    add d0, d0, d0 // 892029809067118180126417511105877323328x d0
    add d0, d0, d0 // 178405961813423636025235022221175666656x d0
    add d0, d0, d0 // 35681192362684727205047004444235333312x d0
    add d0, d0, d0 // 71362384725369454410094008888470666624x d0
    add d0, d0, d0 // 14272476945073890882018801777694133248x d0
    add d0, d0, d0 // 28544953890147781764037603555388266496x d0
    add d0, d0, d0 // 57089907780295563528075207110776532992x d0
    add d0, d0, d0 // 11417981556059112705615041422155313584x d0
    add d0, d0, d0 // 22835963112118225411230082844310627168x d0
    add d0, d0, d0 // 45671926224236450822460165688621254336x d0
    add d0, d0, d0 // 91343852448472901644920321377242508672x d0
    add d0, d0, d0 // 18268770489694580328984064275448501744x d0
    add d0, d0, d0 // 36537540979389160657968128550897003488x d0
    add d0, d0, d0 // 73075081958778321315936257101794006976x d0
    add d0, d0, d0 // 14615016391756664263187251420358801392x d0
    add d0, d0, d0 // 29230032783513328526374502840717602784x d0
    add d0, d0, d0 // 58460065567026657052749005681435205568x d0
    add d0, d0, d0 // 11692013113405331410549801136287001136x d0
    add d0, d0, d0 // 23384026226810662821099602272574002272x d0
    add d0, d0, d0 // 46768052453621325642199204545148004544x d0
    add d0, d0, d0 // 93536104000000000000000000000000000000x d0
    add d0, d0, d0 // 187072200000000000000000000000000000000x d0
    add d0, d0, d0 // 374144400000000000000000000000000000000x d0
    add d0, d0, d0 // 748288800000000000000000000000000000000x d0
    add d0, d0, d0 // 1496577600000000000000000000000000000000x d0
    add d0, d0, d0 // 2993155200000000000000000000000000000000x d0
    add d0, d0, d0 // 5986310400000000000000000000000000000000x d0
    add d0, d0, d0 // 11972620800000000000000000000000000000000x d0
    add d0, d0, d0 // 23945241600000000000000000000000000000000x d0
    add d0, d0, d0 // 47890483200000000000000000000000000000000x d0
    add d0, d0, d0 // 95780966400000000000000000000000000000000x d0
    add d0, d0, d0 // 191561932800000000000000000000000000000000x d0
    add d0, d0, d0 // 383123865600000000000000000000000000000000x d0
    add d0, d0, d0 // 7662
```