

# 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

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

$$2^1 = 2 \quad 2^4 = 16 \quad 2^7 = 128$$

$$2^2 = 4 \quad 2^5 = 32 \quad 2^8 = 256$$

**Little Endian:** kleinstes Byte zuerst

$$2^3 = 8$$

$$2^4 = 16$$

$$2^5 = 32$$

= Kilo

$$2^6 = 64$$

$$2^7 = 128$$

$$2^8 = 256$$

$$2^9 = 512$$

$$2^{10} = 1024$$

$$2^{11} = 2048$$

$$2^{12} = 4096$$

$$2^{13} = 8192$$

$$2^{14} = 16384$$

$$2^{15} = 32768$$

$$2^{16} = 65536$$

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

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

$$X \cdot X = X \quad X \cdot \bar{X} = 0 \quad X \cdot (Y \cdot Z) = (X \cdot Y) \cdot (X \cdot Z)$$

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

## 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

## NAND Operations

$$\text{NOT: } \overline{AA} \quad \text{OR: } \overline{AA \cdot BB} \quad \text{AND: } \overline{AB} \cdot \overline{AB}$$

**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

Grey Code		
A	B	C
0	0	0
0	0	1
0	1	0
0	1	1
1	0	0
1	0	1
1	1	0
1	1	1

$\Rightarrow$	$\begin{array}{ c c c c } \hline & C & AB & \\ \hline 0 & 00 & 01 & 11 & 10 \\ \hline 1 & 01 & 10 & 11 & 11 \\ \hline \end{array}$	$\Rightarrow X = \bar{B} + A \cdot \bar{C}$
---------------	--	---

$$\text{NOT: } \overline{A+A} \quad \text{OR: } \overline{(A+A)} + \overline{(B+B)}$$

$$\text{AND: } \overline{(A+B)} \cdot \overline{(A+B)}$$

## 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	1	0
0	1	0	0	1	1	0	1
1	0	0	1	0	1	0	1
1	1	1	1	1	0	0	0

## Combinational / Sequential Logic

- Comb:** - no memory  
- no cyclic paths  
- combines inputs to get output  
- signals are assigned in every possible condition  
- all outputs 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, 1000

- 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 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 and 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.

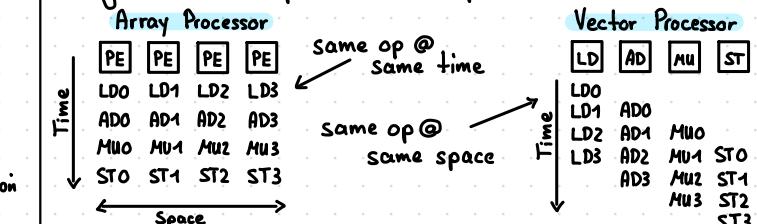
## Fine Grained Multithreading

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

$$\text{Execution time} = (\# \text{ instructions}) / (\text{instructions}/\text{cycle}) \cdot \frac{\text{cycles}}{\text{seconds}} \quad \text{MIPS: million instructions per sec}$$

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

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

L = Main Memory, VM = 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 (including executed and predicted)
- D = number of cycles stalled for each condition branch

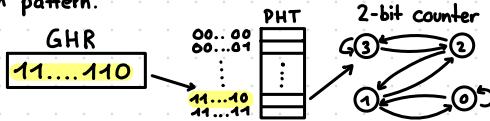
## Delays

Propagation delay: tpd: maximum 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).

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

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

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

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

same a physical tag

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

Cache	Virtual Memory
Block	Page
Block Size	Page Size
Block Offset	Page Offset
Miss	Page Fault
Index/Tag	Virtual Page Number

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 mbyte aligned address would have only  $\log_2(n)$  valid significant bytes when represented in binary.

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

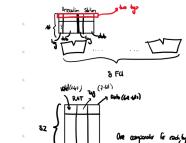
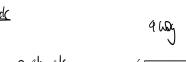
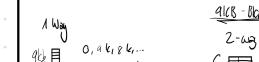
Index Bits =  $\log_2(\text{Blocks per Set})$

Block Offset Bits =  $\log_2(\text{Block Size in Bytes})$

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

Software Interlocking: NOPs

Hardware Interlocking: Detects dependencies and stalls pipeline accordingly



Utilization:  $(\text{Threads} - \text{instructions each thread is taking}) / (\text{Threads} - \text{idle instructions taken at any point in time})$

Worst:  $(\text{Threads} - \text{idle threads}) / (\text{Threads} - \text{idle threads per core})$

Min Utilization: following different paths when in general no many parallel tasks are shared paths, but for each other paths parallel we use 1 thread. If there is a complex task, instead threads there are just not send to further threads + 0 (but needs to be in the memory)

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:  $\sum(3, 5, 7, 11, 13, 15)$ .  $\{3, 5, 7, 11, 13, 15\} = \{0011, 0101, 0111, 1011, 1101, 1111\}$

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

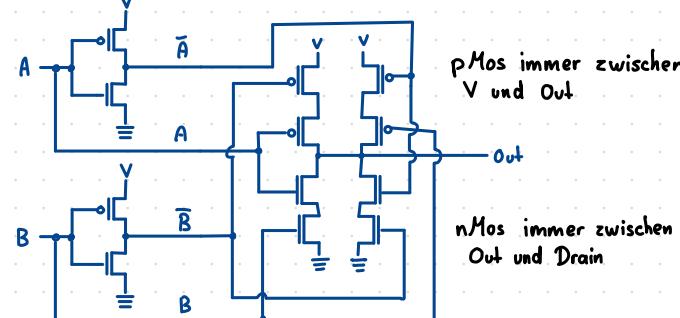
Ex. Convert the following equation to only contain NANDs.

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

Ex. Convert the following equation to only contain NANDs.

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

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
```

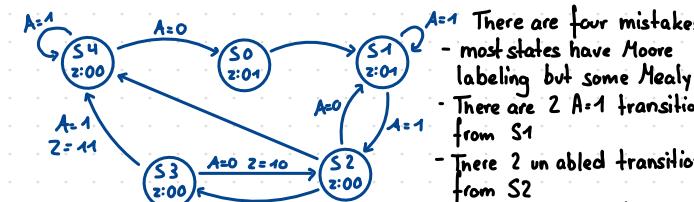
The code implements 2 D-FlipFlops, each works with a asynchronous active low reset.

Ex. Is this code syntactically correct?

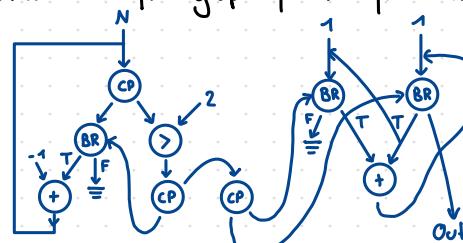
```
module fulladd (input a, b, c, output reg s, c.out);
  assign s = a ^ b;
  assign c.out = (a & b) | (a & c) | (b & c); we use assign therefore these have to be wires
endmodule
```

```
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.



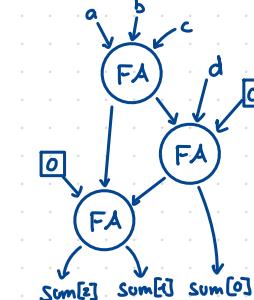
Ex. Draw a data flow graph for the fibonacci function.



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

LOOP2:              16            87

B LOOP2    all runs execute the same #dynamic instr.

Let: C = # cycles

P = #stages

I = # dynamic instr.

B = #branch instr.

D = #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:

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

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

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

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

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

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		
	[30]	7	50	
	[29]			
VLD	[0]	50		
	[31]	50		
	[30]	7	50	
	[29]			
VADD		4		
VMUL		16		
VRSHFA		16		

1 + 50 + 7 + 50 + 4 + 16 + 1 = 129    => 129 cycles

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  $\frac{1}{4}$  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?

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

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

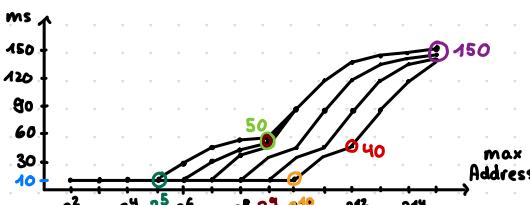
How many bits of the virtual index come from the VPN?

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	$2^5 = 32$	$2^3 = 512$	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?

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

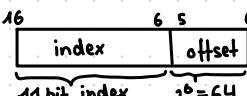
What is the size of the data store?

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

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:



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

What is the virtual address space?

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

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]
```

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

B: int[100] a;

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

Coverage: A:  $\frac{248}{250}$  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  $\Pr[B_1.T \& B_2.T] = 1/2$  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:  $\Pr[B_2.T \& B_3.T] = 1/6 \Rightarrow 1/6 \cdot 1 = 1/6$

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

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

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>0) {
        c=a+b;
        a=b;
        b=c;
        n--;
    }
    return c;
}
```

```
branch:
    blt d1, g, .L1
.L1: add d2, d0, d1
    add d3, d1, d2
    add d4, d2, d3
    add d5, d3, d4
    add d6, d4, d5
    add d7, d5, d6
    add d8, d6, d7
    add d9, d7, d8
    add d10, d8, d9
    add d11, d9, d10
    add d12, d10, d11
    add d13, d11, d12
    add d14, d12, d13
    add d15, d13, d14
    add d16, d14, d15
    add d17, d15, d16
    add d18, d16, d17
    add d19, d17, d18
    add d20, d18, d19
    add d21, d19, d20
    add d22, d20, d21
    add d23, d21, d22
    add d24, d22, d23
    add d25, d23, d24
    add d26, d24, d25
    add d27, d25, d26
    add d28, d26, d27
    add d29, d27, d28
    add d30, d28, d29
    add d31, d29, d30
    add d32, d30, d31
    add d33, d31, d32
    add d34, d32, d33
    add d35, d33, d34
    add d36, d34, d35
    add d37, d35, d36
    add d38, d36, d37
    add d39, d37, d38
    add d40, d38, d39
    add d41, d39, d40
    add d42, d40, d41
    add d43, d41, d42
    add d44, d42, d43
    add d45, d43, d44
    add d46, d44, d45
    add d47, d45, d46
    add d48, d46, d47
    add d49, d47, d48
    add d50, d48, d49
    add d51, d49, d50
    add d52, d50, d51
    add d53, d51, d52
    add d54, d52, d53
    add d55, d53, d54
    add d56, d54, d55
    add d57, d55, d56
    add d58, d56, d57
    add d59, d57, d58
    add d60, d58, d59
    add d61, d59, d60
    add d62, d60, d61
    add d63, d61, d62
    add d64, d62, d63
    add d65, d63, d64
    add d66, d64, d65
    add d67, d65, d66
    add d68, d66, d67
    add d69, d67, d68
    add d70, d68, d69
    add d71, d69, d70
    add d72, d70, d71
    add d73, d71, d72
    add d74, d72, d73
    add d75, d73, d74
    add d76, d74, d75
    add d77, d75, d76
    add d78, d76, d77
    add d79, d77, d78
    add d80, d78, d79
    add d81, d79, d80
    add d82, d80, d81
    add d83, d81, d82
    add d84, d82, d83
    add d85, d83, d84
    add d86, d84, d85
    add d87, d85, d86
    add d88, d86, d87
    add d89, d87, d88
    add d90, d88, d89
    add d91, d89, d90
    add d92, d90, d91
    add d93, d91, d92
    add d94, d92, d93
    add d95, d93, d94
    add d96, d94, d95
    add d97, d95, d96
    add d98, d96, d97
    add d99, d97, d98
    add d100, d98, d99
    add d101, d99, d100
    add d102, d100, d101
    add d103, d101, d102
    add d104, d102, d103
    add d105, d103, d104
    add d106, d104, d105
    add d107, d105, d106
    add d108, d106, d107
    add d109, d107, d108
    add d110, d108, d109
    add d111, d109, d110
    add d112, d110, d111
    add d113, d111, d112
    add d114, d112, d113
    add d115, d113, d114
    add d116, d114, d115
    add d117, d115, d116
    add d118, d116, d117
    add d119, d117, d118
    add d120, d118, d119
    add d121, d119, d120
    add d122, d120, d121
    add d123, d121, d122
    add d124, d122, d123
    add d125, d123, d124
    add d126, d124, d125
    add d127, d125, d126
    add d128, d126, d127
    add d129, d127, d128
    add d130, d128, d129
    add d131, d129, d130
    add d132, d130, d131
    add d133, d131, d132
    add d134, d132, d133
    add d135, d133, d134
    add d136, d134, d135
    add d137, d135, d136
    add d138, d136, d137
    add d139, d137, d138
    add d140, d138, d139
    add d141, d139, d140
    add d142, d140, d141
    add d143, d141, d142
    add d144, d142, d143
    add d145, d143, d144
    add d146, d144, d145
    add d147, d145, d146
    add d148, d146, d147
    add d149, d147, d148
    add d150, d148, d149
    add d151, d149, d150
    add d152, d150, d151
    add d153, d151, d152
    add d154, d152, d153
    add d155, d153, d154
    add d156, d154, d155
    add d157, d155, d156
    add d158, d156, d157
    add d159, d157, d158
    add d160, d158, d159
    add d161, d159, d160
    add d162, d160, d161
    add d163, d161, d162
    add d164, d162, d163
    add d165, d163, d164
    add d166, d164, d165
    add d167, d165, d166
    add d168, d166, d167
    add d169, d167, d168
    add d170, d168, d169
    add d171, d169, d170
    add d172, d170, d171
    add d173, d171, d172
    add d174, d172, d173
    add d175, d173, d174
    add d176, d174, d175
    add d177, d175, d176
    add d178, d176, d177
    add d179, d177, d178
    add d180, d178, d179
    add d181, d179, d180
    add d182, d180, d181
    add d183, d181, d182
    add d184, d182, d183
    add d185, d183, d184
    add d186, d184, d185
    add d187, d185, d186
    add d188, d186, d187
    add d189, d187, d188
    add d190, d188, d189
    add d191, d189, d190
    add d192, d190, d191
    add d193, d191, d192
    add d194, d192, d193
    add d195, d193, d194
    add d196, d194, d195
    add d197, d195, d196
    add d198, d196, d197
    add d199, d197, d198
    add d200, d198, d199
    add d201, d199, d200
    add d202, d200, d201
    add d203, d201, d202
    add d204, d202, d203
    add d205, d203, d204
    add d206, d204, d205
    add d207, d205, d206
    add d208, d206, d207
    add d209, d207, d208
    add d210, d208, d209
    add d211, d209, d210
    add d212, d210, d211
    add d213, d211, d212
    add d214, d212, d213
    add d215, d213, d214
    add d216, d214, d215
    add d217, d215, d216
    add d218, d216, d217
    add d219, d217, d218
    add d220, d218, d219
    add d221, d219, d220
    add d222, d220, d221
    add d223, d221, d222
    add d224, d222, d223
    add d225, d223, d224
    add d226, d224, d225
    add d227, d225, d226
    add d228, d226, d227
    add d229, d227, d228
    add d230, d228, d229
    add d231, d229, d230
    add d232, d230, d231
    add d233, d231, d232
    add d234, d232, d233
    add d235, d233, d234
    add d236, d234, d235
    add d237, d235, d236
    add d238, d236, d237
    add d239, d237, d238
    add d240, d238, d239
    add d241, d239, d240
    add d242, d240, d241
    add d243, d241, d242
    add d244, d242, d243
    add d245, d243, d244
    add d246, d244, d245
    add d247, d245, d246
    add d248, d246, d247
    add d249, d247, d248
    add d250, d248, d249
    add d251, d249, d250
    add d252, d250, d251
    add d253, d251, d252
    add d254, d252, d253
    add d255, d253, d254
    add d256, d254, d255
    add d257, d255, d256
    add d258, d256, d257
    add d259, d257, d258
    add d260, d258, d259
    add d261, d259, d260
    add d262, d260, d261
    add d263, d261, d262
    add d264, d262, d263
    add d265, d263, d264
    add d266, d264, d265
    add d267, d265, d266
    add d268, d266, d267
    add d269, d267, d268
    add d270, d268, d269
    add d271, d269, d270
    add d272, d270, d271
    add d273, d271, d272
    add d274, d272, d273
    add d275, d273, d274
    add d276, d274, d275
    add d277, d275, d276
    add d278, d276, d277
    add d279, d277, d278
    add d280, d278, d279
    add d281, d279, d280
    add d282, d280, d281
    add d283, d281, d282
    add d284, d282, d283
    add d285, d283, d284
    add d286, d284, d285
    add d287, d285, d286
    add d288, d286, d287
    add d289, d287, d288
    add d290, d288, d289
    add d291, d289, d290
    add d292, d290, d291
    add d293, d291, d292
    add d294, d292, d293
    add d295, d293, d294
    add d296, d294, d295
    add d297, d295, d296
    add d298, d296, d297
    add d299, d297, d298
    add d300, d298, d299
    add d301, d299, d300
    add d302, d300, d301
    add d303, d301, d302
    add d304, d302, d303
    add d305, d303, d304
    add d306, d304, d305
    add d307, d305, d306
    add d308, d306, d307
    add d309, d307, d308
    add d310, d308, d309
    add d311, d309, d310
    add d312, d310, d311
    add d313, d311, d312
    add d314, d312, d313
    add d315, d313, d314
    add d316, d314, d315
    add d317, d315, d316
    add d318, d316, d317
    add d319, d317, d318
    add d320, d318, d319
    add d321, d319, d320
    add d322, d320, d321
    add d323, d321, d322
    add d324, d322, d323
    add d325, d323, d324
    add d326, d324, d325
    add d327, d325, d326
    add d328, d326, d327
    add d329, d327, d328
    add d330, d328, d329
    add d331, d329, d330
    add d332, d330, d331
    add d333, d331, d332
    add d334, d332, d333
    add d335, d333, d334
    add d336, d334, d335
    add d337, d335, d336
    add d338, d336, d337
    add d339, d337, d338
    add d340, d338, d339
    add d341, d339, d340
    add d342, d340, d341
    add d343, d341, d342
    add d344, d342, d343
    add d345, d343, d344
    add d346, d344, d345
    add d347, d345, d346
    add d348, d346, d347
    add d349, d347, d348
    add d350, d348, d349
    add d351, d349, d350
    add d352, d350, d351
    add d353, d351, d352
    add d354, d352, d353
    add d355, d353, d354
    add d356, d354, d355
    add d357, d355, d356
    add d358, d356, d357
    add d359, d357, d358
    add d360, d358, d359
    add d361, d359, d360
    add d362, d360, d361
    add d363, d361, d362
    add d364, d362, d363
    add d365, d363, d364
    add d366, d364, d365
    add d367, d365, d366
    add d368, d366, d367
    add d369, d367, d368
    add d370, d368, d369
    add d371, d369, d370
    add d372, d370, d371
    add d373, d371, d372
    add d374, d372, d373
    add d375, d373, d374
    add d376, d374, d375
    add d377, d375, d376
    add d378, d376, d377
    add d379, d377, d378
    add d380, d378, d379
    add d381, d379, d380
    add d382, d380, d381
    add d383, d381, d382
    add d384, d382, d383
    add d385, d383, d384
    add d386, d384, d385
    add d387, d385, d386
    add d388, d386, d387
    add d389, d387, d388
    add d390, d388, d389
    add d391, d389, d390
    add d392, d390, d391
    add d393, d391, d392
    add d394, d392, d393
    add d395, d393, d394
    add d396, d394, d395
    add d397, d395, d396
    add d398, d396, d397
    add d399, d397, d398
    add d400, d398, d399
    add d401, d399, d400
    add d402, d400, d401
    add d403, d401, d402
    add d404, d402, d403
    add d405, d403, d404
    add d406, d404, d405
    add d407, d405, d406
    add d408, d406, d407
    add d409, d407, d408
    add d410, d408, d409
    add d411, d409, d410
    add d412, d410, d411
    add d413, d411, d412
    add d414, d412, d413
    add d415, d413, d414
    add d416, d414, d415
    add d417, d415, d416
    add d418, d416, d417
    add d419, d417, d418
    add d420, d418, d419
    add d421, d419, d420
    add d422, d420, d421
    add d423, d421, d422
    add d424, d422, d423
    add d425, d423, d424
    add d426, d424, d425
    add d427, d425, d426
    add d428, d426, d427
    add d429, d427, d428
    add d430, d428, d429
    add d431, d429, d430
    add d432, d430, d431
    add d433, d431, d432
    add d434, d432, d433
    add d435, d433, d434
    add d436, d434, d435
    add d437, d435, d436
    add d438, d436, d437
    add d439, d437, d438
    add d440, d438, d439
    add d441, d439, d440
    add d442, d440, d441
    add d443, d441, d442
    add d444, d442, d443
    add d445, d443, d444
    add d446, d444, d445
    add d447, d445, d446
    add d448, d446, d447
    add d449, d447, d448
    add d450,
```