

Number systems

$$\begin{array}{lllll}
 2^1 = 2 & 2^2 = 4 & 2^3 = 8 & 2^4 = 16 & 2^5 = 32 \\
 2^6 = 64 & 2^7 = 128 & 2^8 = 256 & 2^9 = 512 & 2^{10} = 1024 \\
 2^{11} = 2048 & 2^{12} = 4096 & 2^{13} = 8192 & 2^{14} = 16384 & \\
 2^{15} = 32768 & 2^{16} = 65536 & 2^{20} = 1048576
 \end{array}$$

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

Kilo	10^3	2^{10}	peta	10^{15}	2^{50}
mega	10^6	2^{20}	exa	10^{18}	2^{60}
giga	10^9	2^{30}	zetta	10^{21}	2^{70}
terra	2^{12}	2^{40}	yotta	10^{24}	2^{80}

Big Endian: Counting Bytes right to left ←
Little Endian: counting Bytes left to right →
unsigned: No sign-flag, all positive $[0, 2^n - 1]$
signed: with sign-flag, pos. & neg. $[-2^{n-1}, 2^{n-1}]$
 ↳ -0 and +0 exists! **No addition!**

2's complement: with sign-flag $[-2^{n-1}, 2^{n-1} - 1]$
 ↳ 0 exists once, Addition work

① Invert each bit ② add 1 **2's compl**

Subtraction: Take 2's compl of 2^0 numb
and than add

Example: Data Prefetching A, A+1, A+2, A+7, A+8, A+9, A+14, ...

- Stride PF observing last 3 requests. Coverage = 0, PF will learn const. Stride of +1, But PF blocks userless, none of the requests will be covered
- Accuracy 0%. None of prefetched cache blocks will be accurate
- New PF, prefetches next N cache blocks $N=2$, f. ex. the PF requests block 22 and 23 after seeing the request to 21. Every 2 out of 3 requests will be covered, and every 3 out of 6 requests will be correct. Cov = $\frac{2}{3}$, ACL = $\frac{1}{2}$
- bandwidth overhead $\frac{5}{3}$. For every group of 3 requests, PF will fetch 2 additional blocks.
- Min. of N to achieve 100% cov. $N=5$, Exmpl. ⇒ bandwidth is $\frac{7}{3}$

Boolean Logic

Product of Sum:

- All FALSE lines,
- Maxterms ANDed,
- Inputs from table inverted

$$\begin{array}{c|ccc}
 A & B & Y \\
 \hline
 0 & 0 & 1 \\
 0 & 1 & 0 \\
 1 & 0 & 1 \\
 1 & 1 & 0
 \end{array}
 \quad
 \begin{array}{c|cc}
 A+B & Y=(A+B)(\bar{A}+\bar{B}) \\
 \bar{A}+\bar{B} &
 \end{array}$$

Sum of Products

- All TRUE lines
- Minterms ORed
- NOT inverted, directly taken

$$\begin{array}{c|ccc}
 A & B & Y \\
 \hline
 0 & 0 & 0 \\
 0 & 1 & 1 \\
 1 & 0 & 1 \\
 1 & 1 & 0
 \end{array}
 \quad
 Y = \bar{A}B + A\bar{B}$$

Identity: $A \cdot 1 = A$ $A + 1 = 1$

Null Element: $A \cdot 0 = 0$ $A + 0 = A$

Idempotency: $A \cdot A = A$ $A + A = A$

Involution: $\bar{\bar{A}} = A$

Complements: $A \cdot \bar{A} = 0$ $A + \bar{A} = 1$

Commutativity: $A \cdot B = B \cdot A$...

Covering: $A(A+B) = A$ $A+(AB) = A$ $AB + A\bar{B} = A$

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

Assoc.: $A(BC) = (AB)C$

$A + (B+C) = (A+B) + C$

Distr.: $A(B+C) = AB + BC$

$A + (BC) = (A+B)(A+C)$

De Morgan: $\bar{ABC} = \bar{A} + \bar{B} + \bar{C}$

$(A+B+C) = \bar{A}\bar{B}\bar{C}$

X: Don't care (value 1 or 0) / Illegal value (value 1 and 0)

Z: Floating value (not 1 nor 0)

Logical completeness: Every formula can be represented by the given set of operations

{AND, OR, NOT} {NAND} {NOR}

Example: $F = (\bar{A}B + C) + A \cdot C$ convert to NAND use only

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

$$= (A \cdot B) \cdot (C \cdot \bar{C}) = (\bar{A} \cdot \bar{B}) \cdot (\bar{A} \cdot B) \cdot (C \cdot \bar{C})$$

Example: $F = \bar{A} + (\bar{B} \cdot C + \bar{A} \cdot \bar{C})$ convert to NOR use only

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

$$\bar{B} \cdot C = \bar{B} + \bar{B} \cdot C + \bar{C} \quad \bar{A} \cdot \bar{C} = \bar{A} + \bar{A} \cdot \bar{C} + \bar{C} + \bar{C} \cdot \bar{A}$$

Combinational Circuits

Transistors: MOS (metaloxide semiconductors)

nMOS: Conducts if $g=1$
pull down well but up poorly

pMOS: Conducts if $g=0$
pull up well but down poorly

Logic Gates

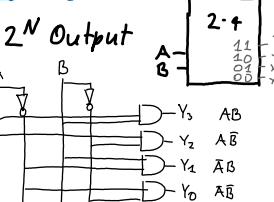
AND	NAND	OR	NOR
$A \cdot B$	$\bar{A} \cdot \bar{B}$	$A + B$	$\bar{A} + \bar{B}$
$A \quad B$	$A \quad B$	$A \quad B$	$A \quad B$
0 0 0	0 0 1	0 0 0	0 0 1
0 1 0	0 1 0	0 1 1	0 1 0
1 0 0	1 0 1	1 0 1	1 0 0
1 1 1	1 1 0	1 1 1	1 1 0

XNOR	XOR	NOT	BUFFER
$\bar{A} \cdot \bar{B}$	$\bar{A} + B$	\bar{A}	A
$A \cdot B$	$A + B$	A	A
0 0 1	0 0 0	0 1 0	0 0 0
0 1 0	0 1 1	1 0 0	0 1 1
1 0 1	1 0 0	1 1 0	1 0 1
1 1 0	1 1 1	1 1 1	1 1 1

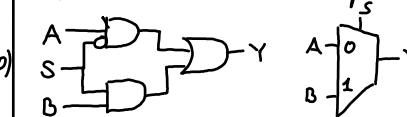
Combinational Building Blocks

Decoder: N Input $\rightarrow 2^N$ Output

A	B	Y_3	Y_2	Y_1	Y_0
0 0	0 0	0 0 0 1			
0 1	0 1	0 0 1 0			
1 0	1 0	0 1 0 0			
1 1	1 1	1 0 0 0			



MUX: selects one of n inputs with $\log_2 n$ select signals



PLA (Prog. Logic Array)

→ directly from Sum of Products form

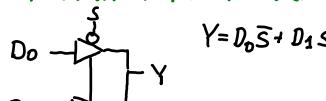
→ For N -Input logic func, we need 2^n n -Input AND Gates

Tristate Buffer

Acts like a switch

E	A	Y
0	0	Z
0	1	2
1	0	0
1	1	1

Mux with Tristate Buf



Combinational VS Sequential

• No Memory

• No cyclic paths

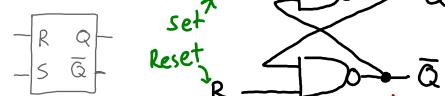
• Every node either input or connects to exactly one output

• Has Memory

• Output depends on prior inputs

Sequential Circuits

SR-Latch



D - Latch

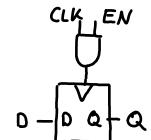


D - Flip-Flop

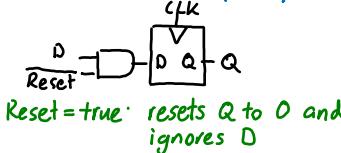


Value from D gets copied to Q at rising clock edge
If clock=0: Q retains old value

Enabled Flip Flop



Resettable Flip Flop



Synchronous Sequential Circuits

- Every element either a register or a comb. circuit
- At least one circuit element is a register
- All registers receive the same clock signal
- Every cyclic path contains at least one reg

Finite State Machines (FSM)

Moore: • output depends only on current state
• output labeled in each state

Mealy: • output depends on current state & input values
• output labeled on the transitions
• A/Y A= value of Input Y= output

A Mealy machines usually requires less states

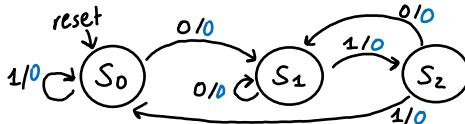
Area of FSM: #bits for state encoding + # logic gates in next-state logic

Example: How to simplify a FSM?

Does a state has only outgoing edges? Delete

Does two states have the same behaviour? Delete one

Example:
Detect 011



R	S	Q
0	0	1
0	1	0
1	0	0
1	1	Q _{prev}

WE	D	Q
0	0	Q _{prev}
0	1	Q _{prev}
1	0	0
1	1	1

Register $\hat{=}$ several D Flip Flops

State Encoding:

Binary Encoding (00, 01, 10, 11)

$\log_2(\# \text{states})$ bits needed \Rightarrow reduces # Flip Flops

One-hot Encoding (001, 010, 100)

states = bits needed \Rightarrow reduces next-state logic

Output Encoding (010, 000, 110, ...)

\Rightarrow reduces output logic

Correctness of a state diagramm

- Must have reset-line or initial state
- no multiple transitions for same input
- no missing transition
- no mix of Moore / Mealy Labeling
- no unmarked transition

Designing a FSM

- Identify inputs & outputs
- sketch the state transition diagramm
- Write a state transition table
- Write an output table
- Select state Encoding
- Write Boolean equations for next-state & output logic
- Sketch circuit

Verilog

\sim	negation	\wedge	bit-wise XOR
\neg	logical negation	$\wedge\wedge$	bit-wise XNOR
$\&\&$	log. AND	{...}	concatenation
$\ \ $	log. OR	$<</>>$	shift left/right
$\&$	bit-wise AND	=	blocking at line
$ $	bit-wise OR	$<=$	non-blocking at end

Wires: • always on the left of assign
• cannot be used at left of = / <=
• are used to connect input & output
• can only be assigned once

Regs: • Cannot be used with assign
• in always-blocks, never outside
• cannot be connected to an output port
• Cannot be used in input ports

CHECK IF:

- All inputs & outputs have to be assigned and names have to match
- not multiple assignments to same signal
- names can't start with numbers / No duplicate-names
- no module recursion
- enough bits for numbers \rightarrow Semicolons
- initialization of variables

Example: D-FlipFlop

```
module flop (input clk, input rst, input [3:0] d, output reg [3:0] d);
  always @ (posedge clk, negedge rst) begin
    if (rst == 1) q <= 0;
    else q <= d;
  end
endmodule
```

Don't include for a sync. reset

Example: Full Adder

```
assign sum = a ^ b ^ c_in;
assign c_out = (a & b) | (a & c_in)
           | (b & c_in);
```

Example: Counter

```
always @ *
```

- begin
- count_next = count;
- count_next += 1;
- end

```
always @ (posedge clk)
  count <= count_next;
```

MIPS Assembly

Name	Number	Use
\$0	0	const. value 0
\$at	1	assembler temporary
\$v0-\$v1	2-3	return values
\$a0-\$a3	4-7	arguments for procedure
\$t0-\$t7	8-15	temporary variables
\$s0-\$s7	16-23	saved variables
\$t8-\$t9	24-25	temporary variables
\$k0-\$k1	26-27	OS temporaries
\$gp	28	global pointer
\$sp	29	stack pointer
\$fp	30	frame pointer
\$ra	31	procedure return address

R-Type

Op 6	rs 5	rt 5	rd 5	shamt 5	function 6
source				destinat.	

I-Type

Op 6	rs 5	rt 5	immediate 16
source		st dest.	

J-Type

Op 6	address 26
------	------------

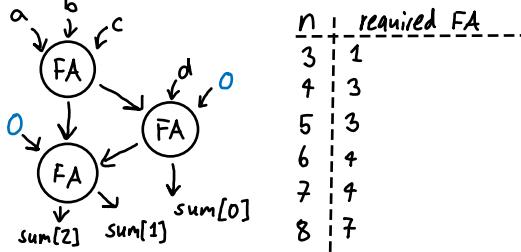
Example: Subtraction 64-bit

```
{$4,$5} - {$6,$7}
subu $3, $5, $7
add $2, $6, $2
slt $2, $5, $7
sub $2, $4, $2
```

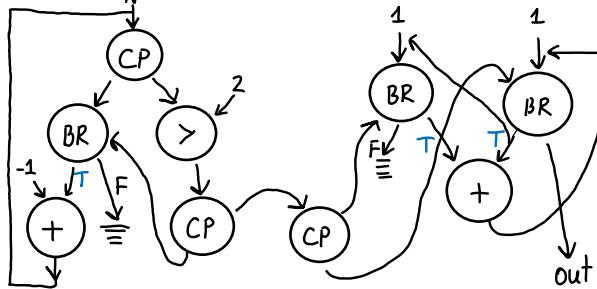
Dataflow

Programm consists of dataflow nodes. A node fires (executes) when all its inputs are ready.

Example Four 1-bit Addition (given Full Adder nodes)



Example Fibonacci Function



Performance Analysis

CPI Cycles per Instruction

$MHz = 10^6 \text{ cycles/sec}$

Clock Freq

IPC Instruction per Cycle = CPI^{-1} doesn't matter!

MIPS Mill. instr. per second = $#MHz / CPI = #MHz / IPC$

Time = #Instructions • CPI / #Hz

Speedup = oldTime / newTime - 1

Higher Frequency \Rightarrow more Instr. per second (Not IPC!)

Higher IPC \Rightarrow faster execution of programm
(but could need more Instr.)

$CPI = \sum_{A \in \text{Instr.}} (\# \text{cycles per } A) \cdot (\text{proportion of } A)$

Execution_time = #Instruction • CPI • clock_time

clock_time = 1 / clock_frequency

Example Fibonacci n in \$4, result in \$2

```

fib: addi $sp, $sp, -16 // allocate stack space
      sw $16, 0($sp) // save r16 in stack
      add $16, $4, $0 // r16 = n
      sw $17, 4($sp) // save r17 in stack
      add $17, $0, $0 // r17 = a = 0
      sw $18, 8($sp) // save r18 in stack
      addi $18, $0, 1 // r18 = b = 1
      sw $31, 12($sp) // save return address
      add $2, $17, $18 // C = a+b
  
```

branch:

```

      slti $3, $16, 2 // use r3 as temp
      bne $3, $0, done // we have finished
      add $2, $17, $18 // C = a+b
      add $17, $18, $0 // a = b
      add $18, $2, $0 // b = C
      addi $16, $16, -1 // n -= 1
      j branch
  
```

done:

```

      lw $31, 12($sp)
      lw $18, 8($sp)
      lw $17, 4($sp)
      lw $16, 0($sp)
      addi $sp, $sp, 16
      jr $31
  
```

} restore

// jump to caller

```

int fib (int n) {
    a=0; b=1;
    C=a+b;
    while (n>1) {
        C=a+b;
        a=b; b=C;
        n--;
    }
    return C;
}
  
```

ISA VS Microarchitecture

ISA Interface between SW & HW. "What the prgr sees"
the machine can't execute subtraction
the machine doesn't have condition codes

A 5-bit imm. can be used in ADD

There are 8 GPR The PC width Memory

The mem-mapped location of exception vectors

Func. of each bit in a progr branch-predictor config. reg

The execution order of load & store in a multicore CPU

The mem.add. nodes available for arith. ops

The ability to flush the TLB from the OS

Exception Handling

SIMD Instr. support

The addr. of mem.-mapped devices of the CPU (e.g. keys)

The #bits required for dest reg of load

Support of integer, division/multiplication

Mechanism to enter in a system call in the OS

size of addr. memory

Microarchitecture: specifies the underlying implementation

ALU doesn't have subtract unit

It takes n cycles to execute an ADD

A 2-to-1 MUX feeds one of the ALU inputs

The reg file has 1 input & 2 output ports

Latency of branch predictor misprediction

The size of physical memory page

The number of pipeline stages in the CPU

The number of cache sets per cache hierarchy

The integer multi. Algo. of the ALU

CPU: size of ROB, fetch width, #non-prog. regs

Branch misprediction penalty

The mem. controller's scheduling algorithm

In order VS OoO Exec

Pipeline

Stages

Fetch: CPU reads Instr. from Instr. memory

Decode: CPU reads source operands from reg. file & decodes the Instr to produce the control signals

Execute: The processor performs the comp. with ALU

Memory: The proc. reads or writes data memory

Writeback: The proc. writes the result to the reg. file

Scoreboarding: Each reg has a valid bit (1 if data is there), stall pipeline if not all sources ready, Problem: stalls for all dependencies (not only flow)

Data Forwarding: Forward the data of a stage directly to where it is needed (solves RAW)

$\hookrightarrow W \rightarrow D$: internal / register file forwarding
 $\hookrightarrow M \rightarrow E$: operand forwarding

Causes of stall: Resource Contention, Data Dependency, Control Dependency, long multi-cycle ops

Data Dependencies: Flow (R/W) "real", Output (W/Aw), Anti (WAR) WAW & WAR only a prob. bc reg. are limited

Handling Flow Dependencies Detect + stall, Detect + forward, Detect + eliminate on Sh level, Detect + move out of way (OoO Execution), Predict, Fine-grained Multithreading

Examples Pipelined Machine

In-order dispatch without forwarding

F D 1 2 3 4 5 6 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

MUL R3 \leftarrow R1, R2
ADD R5 \leftarrow R3, R4
ADD R7 \leftarrow R2, R6
ADD R10 \leftarrow R8, R9
MUL R11 \leftarrow R7, R10
ADD R5 \leftarrow R5, R11

In-order with forwarding

F D 1 2 3 4 5 6 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

F D - - - - D 1 2 3 4 W

Example Prefetching

A: `uint8_t a[1000];`

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

B: `uint8_t a[1000];`

```
sum = 0
for (i=1; i < 1000; i+=4)
    sum += a[i]
```

Accuracy & Coverage

A: `a[0], a[4], a[8], ..., a[996]`

#mem access = $1000/4 = 250$

`a[0], a[4]` are misses, then the prefetcher learns that stride is 4.

248 cache blocks are prefetched, while

248 are used \Rightarrow accuracy = $\frac{248}{250}$ coverage = $\frac{248}{250}$

B: `a[1], a[9], a[16], a[64], a[256] ...` \Rightarrow 5 acc. to Mem. But stride between acc. not constant \Rightarrow prefetcher cannot prefetch \Rightarrow acc. and cov. is 0

b) Better Prefetcher for A: next block prefetch, `a[8]` hit as well \Rightarrow acc/cov = $\frac{249}{250}$
Better Prefetcher for B: Maybe runahead executions / dual-core exec. improve

c) Runahead for A: No. Could provide good perf., but maybe not all blocks gets prefetched, as runahead exec. happens only when stalled on cache miss
Runahead for B: Yes. could potentially prefetch some cache blocks, espec. as the addr. of cache blocks does not depend on the data from the pending cache miss.

Tomasulo's Algorithm

Tag Comp. per res. stat. = #sources of ops • #func. units = $2 \times 8 = 16$

Total tag comp = #units • #st. per unit • #comp. per st. + #reg • #fun. un. = $16 \cdot 16 \cdot 8 + 32 \cdot 8 = 2304$

min tag size = $\log_2(\# \text{units} \cdot \# \text{station per unit}) = \log(16 \cdot 8) = 7$

min size of reg. alias table = #regs • (64/32-bit + tag size + 1) = $32 \cdot (64/32 + 7 + 1)$

min tot. size of tag store = #regs • tag size + #regs • #func. units • #comp. per st. • tag size

Example:

MUL R ₁ , R ₂ \rightarrow R ₃	F D 1 2 3 4 5 6 W = 32 7 2 8 16 7
ADD R ₃ , R ₄ \rightarrow R ₅	F D - - - - 1 2 3 4 W
ADD R ₂ , R ₆ \rightarrow R ₂	F D 1 2 3 + W
ADD R ₈ , R ₉ \rightarrow R ₁₀	F D 1 2 3 4 W
MUL R ₇ , R ₁₀ \rightarrow R ₁₁	F D - - - 1 2 3 4 5 6 W
ADD R ₅ , R ₁₁ \rightarrow R ₅	F D - - - - 1

Reg Valid Tag Value

Reg	Valid	Tag	Value
R ₁	1	-	1
R ₂	1	-	2
R ₃	0	X	-
R ₄	1	-	3
R ₅	0	d(a)	-
R ₆	1	-	4
R ₇	0	b	-
R ₈	1	-	5
R ₉	1	-	6
R ₁₀	0	C	-
R ₁₁	0	Y	-

8 func units, 32-64 bit registers,

16 reservation stations, 2 source/reg/ret. st.

Source 1 + Source 2

V	Tag	Value	V	Tag	Value
a	0	x	1	~	4
b	1	~	2	1	~
c	1	~	8	1	~
d	0	a	0	y	

Source 1 • Source 2

V	Tag	Value	V	Tag	Value
x	1	~	1	1	~
y	0	b	0	c	
z					
t					

Delayed Branching

Delay slot: #Instr. that need to be fetched after a branch. Compiler takes Instr. if they don't influence the branch itself, else delay slot filled with NOP's.

Branch Prediction

Static: Always taken: ~30% - 40% accuracy
Always not taken: ~60% - 70% accuracy
BT FN: Good for loops (fails once)

Dynamic: Last time predictor: Single bit per branch stored in the Branch Target Buffer (BTB). \Rightarrow Misses 2x for loops \Rightarrow 100% accuracy.

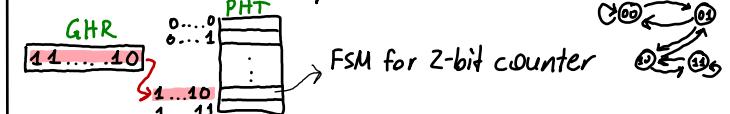
Two-bit counter based prediction: Counter is capped \Rightarrow 00, 1 = 00, 11, 1 = 11. If taken +1, else -1 \Rightarrow Accuracy: ~85% - 90%. (00, 01 \Rightarrow not taken)

Local VS Global: Local: PC based \Rightarrow no interference b/w br.
Global: single counter globally shared

Global Branch Correlation: The idea is that recently executed branch outcomes are correlated with the outcome of the next branch \Rightarrow Prediction based on outcome of last branch with same global branch history

First level: Global Branch history register keeps track of the last branch outcome \Rightarrow GHR Index of PHT

Second level: Pattern history table keeps a 2-bit counter for each pattern



VLIW (Very Long Instruction Word)

The compiler finds independent instructions and statically schedules them into a single VLIW instr.

Lock step Execution: If 1 instr. stalls, whole VLIW stalls + no dependency checking necessary
- compiler needs to find N independent instr. per cycle

Example: VLIW VS in-order Superscalar.

Why VLIW faster: supersc. in-order \Rightarrow bubbles in pipeline. VLIW can reorder

Why Superscalar faster: VLIW needs NOP's \Rightarrow lower I-Cache hit rate & higher fetch bandwidth

Superscalar Execution

Fetch/Decode... multiple instr. per cycle

+ higher IPC - complex dependency checking
- more hardware

Fine-Grained Multithreading

Hardware has multiple thread contexts (PC + Reg for each thread) and each cycle we fetch from a different thread

+ no dependencies - extra hardware
+ no branch pred. - reduced single-cycle perform
+ improved throughput, latency, utilization - resource contention
- dependency checking betw. threa.

Systolic Arrays

Instead of a single Processing element (PE) we have an array of PE and carefully orchestrate the dataflow between them \Rightarrow Maximize comput. on PE's

Difference Pipelining: Array struct. non-linear & multi-dimensional. PE's can have connections multi-directional with diff. speeds and local memory & execute kernels

Example: Matrixmult. $(a \ b)(e \ f) = (ae+bg \ af+bh) \quad (c \ d)(g \ h) = (ce+dg \ cf+dh)$

$$\begin{array}{ccc} a & b & \\ c & d & \end{array} \xrightarrow{\quad} \begin{array}{cc} e & f \\ g & h \end{array} \xrightarrow{\quad} \begin{array}{cc} M & R \\ R & Q \end{array}$$

$$M + R + M \cdot N \quad Q = N \quad P = M$$

SIMD (Single Instruction multiple Data)

"operates" Array Processor: same op @ same time

Vector Processor: same op @ same space

Vector Processor

must be
Performs ops on whole array \Rightarrow ops on each elem. independ.
Data is stored in vector registers

Vector chaining: Data forwarding for vectors (dispatch as soon as individual element is ready).

Stride: distance of element (that are adjacent) in Mem

Strip mining: If vector too long \Rightarrow split in multiple smaller

GPU

GPU's are SIMD engines but using Threads (SIMT)

Warp: Set of threads that execute the same instruction (i.e., at the same PC)

Warps = # threads / (thread per warp)

utilization = (# warps perf. useful work) / (# thread slots)

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

Example

Memory Hierarchy & Caching

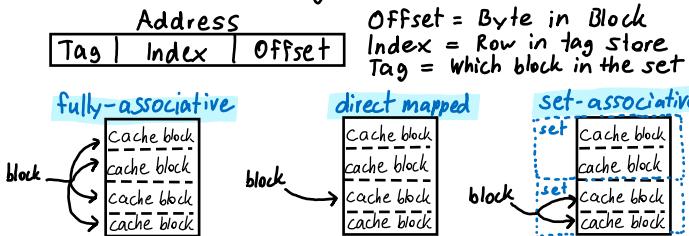
Memory Array: stores data, addr. select. logic selects row, readout circuitry reads data

Memory Banking: Multiple mem. units with common data & addr. Bus, helps to resolve long latency. Units can be accessed individually

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

Blocks & Addressing Cache

- Mem is divided into fix-sized blocks
- Each block maps to a location in the cache (index bits)
- For a **cache hit** the tags need to match



Associativity

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

Cache Performance

Cache size C : total data capacity (excl. tags)

Block size b : data associated with addr. tag

Associativity N : # Blocks represented in same index

$$\# \text{Blocks } B = C/b \quad \# \text{sets } S = B/N$$

Average Memory Access Time:

$$\text{AMAT} = (\text{hit-rate} \cdot \text{hit-latency}) + (\text{miss-rate} \cdot \text{miss-lat.})$$

$$\text{Recursive Latency } T_i = t_i + m_i \cdot T_{i+1} \quad m_i \cdot \text{miss rate}$$

Replacement Policies: • FIFO (first-in-first-out)
• LRU (last-recently-used) • Random

Classification of Misses:

Compulsory Miss: first reference is always a miss

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

Prefetching

Improve cache performance by preloading data before it is used (speculative) in order to reduce mem. latency.

Runahead execution: allows the proc. to pre-process instr. during cache misses instead of stalling

$$\text{Accuracy} = (\text{used prefetches}) / (\text{sent prefetches})$$

$$\text{Coverage} = (\text{prefetched misses}) / (\text{all misses})$$

$$\text{Timeliness} = (\text{on-time prefetches}) / (\text{used prefetches})$$

Virtual Memory (VM)

Virtual address space is divided in **pages**, while phys. addr. space is divided into **frames**. Page Table stores mapping: Virtual → Physical together with a valid bit

$$\# \text{virtual pages} = \frac{\text{Virt. Addr.}}{\text{Page size}} \quad \# \text{physic pages} = \frac{\text{Phys. Addr.}}{\text{Page size}}$$