

## BOOLEAN LOGIC

Theorems:

$$X \cdot Y + X \cdot Y = X$$

$$(X + Y) \cdot (X + Y) = X$$

$$X + X \cdot Y = X$$

$$X \cdot (X + Y) = X$$

$$(X + \bar{Y}) \cdot Y = X \cdot Y$$

$$(X + \bar{Y}) + Y = X + Y$$

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

$$\overline{(X \cdot Y \cdot Z \cdot \dots)} = \bar{X} + \bar{Y} + \bar{Z} + \dots$$

Minterm: product (AND) that includes all input vars

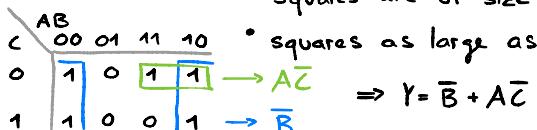
Maxterm: sum (OR) that includes all input variables

SUM-OF-PRODUCTS: Search all TRUE combinations and add (OR) the minterm of each line

$$\begin{array}{|c|c|c|} \hline A & B & Y \\ \hline 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ \hline \end{array} \quad Y = \bar{A}\bar{B} + A\bar{B} = \sum m(0,2)$$

$$Y = (A + \bar{B}) \cdot (\bar{A} + \bar{B}) = \prod M(1,3)$$

PRODUCT-OF-SUMS: Search all FALSE combinations and multiply (AND) the inverted maxterms.

Karnaugh Maps: • use fewest squares of only 1s  
• squares are of size  $2^x$   
• squares as large as possible

Represent 0x cafe1234

Big Endian:

0	1	2	3
ca	fe	01	23

Little Endian:

3	2	1	0
ca	fe	01	23

## Logic Gates

OR	A	B	Y	AND	A	B	Y	NAND	A	B	Y
$\Rightarrow$	0	0	0	$\Rightarrow$	0	0	0	$\Rightarrow$	0	0	1
$\Rightarrow$	0	1	1	$\Rightarrow$	0	1	0	$\Rightarrow$	0	1	1
$\Rightarrow$	1	0	1	$\Rightarrow$	1	0	0	$\Rightarrow$	1	0	1

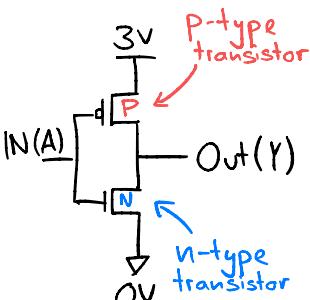
XNOR	A	B	Y	XOR	A	B	Y	NOR	A	B	Y
$\Rightarrow$	0	0	1	$\Rightarrow$	0	0	0	$\Rightarrow$	0	0	1
$\Rightarrow$	0	1	0	$\Rightarrow$	0	1	1	$\Rightarrow$	0	1	0
$\Rightarrow$	1	0	0	$\Rightarrow$	1	0	1	$\Rightarrow$	1	0	0

Example: use only NAND:

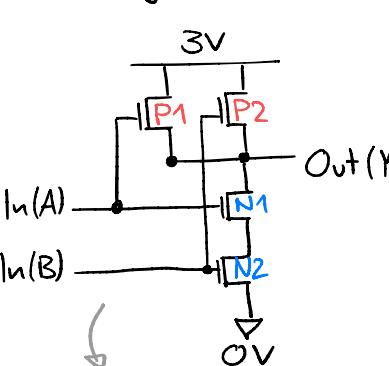
$$\begin{aligned} F &= (A + B \cdot C) + \bar{C} = \overline{(A + B \cdot C) + \bar{C}} \\ &= \overline{(A + B \cdot C) \cdot C} = \overline{\overline{(A + B \cdot C)} \cdot C} \\ &= \overline{(\bar{A} \cdot \bar{B} \cdot \bar{C})} \cdot C = \overline{(\bar{A} \cdot \bar{A}) \cdot (\bar{B} \cdot \bar{C})}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

## CMOS TRANSISTORS

Inverter



NAND gate:



Never pull up and down at the same time!

A	B	P1	P2	N1	N2	Y
0	0	ON	ON	OFF	OFF	1
0	1	ON	OFF	OFF	ON	1
1	0	OFF	ON	ON	OFF	1
1	1	OFF	OFF	ON	ON	0

## VERILOG

Operators:

! logical negation	use this with integers
$\Rightarrow$ any non-zero val $\rightarrow 0$ , only 0 $\rightarrow 1$	
$\&$ bitw./red. AND	$\sim$ bitwise negation
$\sim\&$ reduction NAND	$\sim$ bitw./red. XOR
$\wedge$ bitw./red. XOR	$\sim\wedge$ reduction NOR
$\mid$ logical OR	$\wedge\sim$ bitw./red. XNOR
$\sim\mid$ logical AND	$\wedge\wedge$ logical AND
$\ll$ shift left	$\sim\mid$ 1 iff either input is true or non-zero
$\mid\mid$ logical equality	$\gg$ shift right
$\mid\mid\mid$ case equality	$\neq$ logical inequality
	$\neq\neq$ case inequality

## COMBINATIONAL VS. SEQUENTIAL

↳ All outputs are assigned in every possible way. Also, all inputs are in sensitivity list

## PERFORMANCE ANALYSIS

- CPI: cycles per instruction
- IPC: instructions per cycle
- MHz: frequency,  $10^6$  cycles per second
- MIPS: million instructions / second  $\frac{\text{MHz}}{\text{CPI}}$

## BINARY NUMBERS

$2^1 = 2$	$2^2 = 4$	$2^3 = 8$	$2^4 = 16$
$2^5 = 32$	$2^6 = 64$	$2^7 = 128$	$2^8 = 256$
$2^9 = 512$	$2^{10} = 1024$	$2^{11} = 2048$	$2^{12} = 4096$
$2^{13} = 8192$	$2^{14} = 16384$	$2^{15} = 32768$	$2^{16} = 65536$
			$2^{17} = 131072$
			$2^{18} = 262144$

## SI PREFIXES

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	$10^{21}$	$2^{70}$
tera	$10^{12}$	$2^{40}$	yotta	$10^{24}$	$2^{80}$

## NUMBERS

Unsigned: only positive, 0 to  $2^n - 1$ Signed:  $-2^{n-1}$  to  $2^{n-1} - 1$ , MSB is sign flag

2s Complement: invert each bit and add 1 to the LSB. MSB is sign flag

# FINITE STATE MACHINES

**Moore:** output only dependent on current state, requires more states, synchronous output, output placed on states

**Mealy:** output depends on current state and current inputs, fewer states needed, output placed on transitions

Design a FSM: 1) identify inputs/outputs  
 2) sketch state transition diagram  
 3) write state transition table  
 4) write output table  
 5) select state encoding  
 6) write boolean equations for next state and output logic  
 7) sketch circuit

## State encodings:

- binary:  $\log_2(\# \text{states})$  bits needed  
 reduce # flip-flops to hold state
- one-hot: # states bits needed  
 reduce next-state logic
- output: reduces output logic

A correct state-diagram has:

- reset line
- only one transition per input
- no missing transitions
- no unmarked transitions
- initial state (if no reset)
- no mix of Moore-Mealy labelling  
 output in bubble  $\leftarrow$  output on transition

## Area of FSMs:

- # of FFs = (# bits for state enc.)  $\times$  2
- # of logic gates: count in the next-state logic

Find next-state & output boolean equations

1) Assign state encodings and outputs

2) Create table of the form

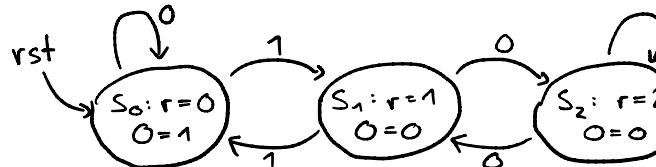
S[0]	S[1]	S[2]	T <sub>A</sub>	T <sub>B</sub>	N[0]	N[1]	N[2]	O[0]	O[1]
0	0	0	x	1	0	0	1	0	0
0	0	0	x	0	0	1	0	0	1
1	1	1	1	0	0	0	1	1	1

3) Find equations for N[x] with sum-of-products:

$$N[2] = (S[0] \cdot S[1] \cdot S[2] \cdot T_A) + (S[0] \cdot S[1] \cdot \overline{S[2]} \cdot T_B)$$

4) Find equations for O[x] with sum-of-products

Example: output 1  $\Leftrightarrow$  input stream %3 = 0



## MIPS ISA (4B instruction size)

r0	\$zero	always 0
r1, r26, r27	\$at, \$k0, \$k1	reserved for assembler
r2-r3	\$v0-\$v1	results
r4-r7	\$a0-\$a3	arguments
r8-r15	\$t0-\$t7	temps, not saved
r16-r23	\$s0-\$s7	saved contents
r24-r25	\$t8-\$t9	temps, not saved
r28	\$gp	global pointer
r29	\$sp	stack pointer
r30	\$fp	frame pointer
r31	\$ra	return address

64-bit integer subtraction:

- A[63:32] in \$4, A[31:0] in \$5
- B[63:32] in \$6, A[31:0] in \$7

$\Rightarrow \text{subu } \$3, \$5, \$7$

$\text{sltu } \$2, \$5, \$7$

$\text{add } \$2, \$6, \$2$

$\text{sub } \$2, \$4, \$2$

fib:

```

addi $sp, $sp, -16
sw $16, 0($sp)
add $16, $4, $0
sw $17, 4($sp)
add $17, $0, $0
sw $18, 8($sp)
addi $18, $0, 1
sw $31, 12($sp)
add $2, $17, $18

```

branch:

```

slti $3, $16, 2
bne $3, $0, done
add $2, $17, $18
add $17, $18, $0
add $18, $2, $0
addi $16, $16, -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

```

allocate stack space  
 save r16  
 r16 for arg n  
 save r17  
 r17 for arg a=0  
 save r18  
 r18 for arg b=1  
 save return address  
 $c = a + b$

use r3 as temp  
 branch if  $n=0$   
 $c = a + b$   
 $a = b$   
 $b = c$   
 $n = n - 1$

} restore registers  
 restore stack pointer  
 return to caller

## REP MOVSB

↳ 3 registers:

ECX, ESI, EDI

↳ repeat ECX-times

copy 1 byte from  
 address ESI to EDI

beq \$1, \$0, AfterLoop

CopyLoop:

lb \$4, 0(\$2)

sb \$4, 0(\$3)

addiu \$2, \$2, 1

addiu \$3, \$3, 1

addiu \$1, \$1, -1

bne \$1, \$0, CopyLoop

Example: ECX  $\rightarrow$  \$1,  
 ESI  $\rightarrow$  \$2, EDI  $\rightarrow$  \$3  
 AfterLoop:

## ISA VS MICROARCHITECTURE

- ISA: interface that is exposed to software changes affect compiler and programmer
  - Instructions: opcodes, addressing modes, data types, instruction types, registers, condition codes
  - Memory: address space, addressability, alignment, virtual memory management
  - Call, interrupt and exception handling
  - Access control, priority and privilege
  - Input-Output: memory-mapped vs. instruction
  - Power and thermal management
  - Multithreading and multiprocessor support
  - Memory-mapped location of exception vectors
  - Function of each bit in a programmable branch-predictor configuration register
  - Order of load/stores in multi-core CPU
  - Memory addressing modes available for arithmetic operations
  - Program counter width
  - General purpose register \$29 is stack pointer
  - Hardware floating point exception support
  - Vector instruction support
  - Memory-mapped I/O port address
  - CPU endianness
  - Virtual page size
  - No condition codes
  - 5-bit immediate can be used in ADD
  - # general-purpose registers
  - size of addressable memory

- Microarchitecture:** specifies how the underlying implementation actually executes the instructions. Changes are invisible to programmer and compiler
- pipelining (e.g. number of stages)
  - in-order vs. out-of-order execution
  - memory access scheduling policy
  - speculative execution
  - superscalar processing
  - clock gating
  - caching: levels, size, associativity, replacement
  - error correction
  - physical structure of registers
  - ALU has no subtract unit
  - # cycles to execute ADD instruction
  - register file has 1 input & 1 output port
  - latency of branch misprediction
  - physical memory page size
  - bit-width of interface CPU  $\leftrightarrow$  L1 cache
  - # cache sets at each level of cache hierarchy
  - max. reservation station capacity
  - ALU critical path
  - Instruction issue width
  - # read/write ports in physical reg. file
  - # ALUs
  - # instructions fetched per clock cycle
  - # prefetches of hardware prefetcher
  - Reorder Buffer size

## PIPELINING

- Fetch: CPU reads instr. from instr. memory
- Decode: CPU reads source operands from register file and decodes the instruction to produce the control signal
- Execute: computation with ALU
- Memory: reads or writes to/from memory
- Writeback: write result to register file

Pipeline stalls happen:

- resource contention
- data / control flow dependencies
- multi-cycle operations

Handle flow dependencies:

- detect and wait
- detect and forward data to dependent instr.
- detect and eliminate in software ( $\rightarrow$  no hardware checks)
- detect and move out of way for indep. instr.
- predict, execute speculatively, verify
- do something else (fine-grained multithreading)

**Reorder Buffer:** complete instructions out-of-order, but reorder them before making results visible. Access register file first, and point to ROB entry which will contain result if reg not valid.

Valid	Dest. Reg. ID	Dest. Reg. Value	Store-Addr.	Store-Data	PC	Val. Bits	Exception?
-------	---------------	------------------	-------------	------------	----	-----------	------------

Handle WAR/WAW dependences:

register renaming (register ID  $\rightarrow$  ROB entry ID)

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
1	MUL	RS, R6, R7											F	D, E, E <sub>2</sub> , E <sub>3</sub> , E <sub>4</sub> , M, W		
2	ADD	R4, R6, R7											F	D, E, E <sub>2</sub> , E <sub>3</sub> , M, W		
3	ADD	RS, R5, R6											F	D, E, E <sub>2</sub> , E <sub>3</sub> , M, W		
4	MUL	R4, R7, R7											F	- D, E, E <sub>2</sub> , E <sub>3</sub> , M, W		
5	ADD	R6, R7, RS											F	D, E, E <sub>2</sub> , E <sub>3</sub> , M, W		
6	ADD	R3, R0, R6											F	- D, E, E <sub>2</sub> , E <sub>3</sub> , M, W		
7	ADD	R7, R1, R4											F	- D, E, E <sub>2</sub> , E <sub>3</sub> , M, W		

ADD: 3 cycles, MUL: 4 cycles

min read/write ports reg file: 2xread, 1xwrite

data forwarding: cycles 7 and 10  
internal forwarding in reg file: cycle 13 same time  
optimize: remove (2) (WAW), move (4) up (stalls for D),  
switch (6) and (7) (no RAW dependency)

# TOMASULO'S ALGORITHM

Implement dynamic scheduling with register renaming to eliminate output (WAW) and anti-dependences (WAR).

Reservation stations: holds the operands of an instruction. Instructions are executed out-of-order, when the operands are ready

OP	RS 1 (0 if ready)	RS 2	Val 1	Val 2	Imm. / Addr.	Busy
----	----------------------	------	-------	-------	-----------------	------

reservation station ID only one of these has a value  
When a value is computed, it is streamed on the common data bus to all RS and reg file.

1. insert instr. & data (source tag or value) into a free reservation station
2. while in reservation station, watch common data bus. If tag seen, grab value. If all operands available, mark as ready for dispatch
3. dispatch instr. to functional unit when ready
4. after FU completed: put tagged value on CDB, reg file has tag for RS that writes to it. If see tag, grab value.

Example: draw RAT and RS after cycle 7:

MUL R1, R2 → R3	F D E <sub>1</sub> E <sub>2</sub> E <sub>3</sub> E <sub>4</sub> E <sub>5</sub>
ADD R3, R4 → R5	F D — — —
ADD R2, R6 → R7	F D E <sub>1</sub> E <sub>2</sub> E <sub>3</sub>
ADD R8, R9 → R10	F D E <sub>1</sub> E <sub>2</sub>
MUL R7, R10 → R11	F D —
ADD R5, R11 → R5	F D

## Register Alias Table

Reg	Valid	Tag	Value
R1	1	-	1
R2	1	-	2
R3	0	x	-
R4	1	-	4
R5	0	d(a)	-
R6	1	-	6
R7	0	b	-
R8	1	-	8
R9	1	-	9
R10	0	c	-
R11	0	y	-

## Source 1 + Source 2

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

## Source 1 \* Source 2

V	Tag	Val	V	Tag	Val
x	1	~	1	1	~
y	0	b	0	0	c
z					
+					

# tag comparators per RS = #FU \* # registers  
↳ if every FU has separate data bus  
total # tag comparators =  $\uparrow$   
usually 2  
#FU \* #RS per FU \* # tag comp. per RS  
+ #FU \* # architectural registers  
min size of tag = #FU \* # RS per FU  
min size of RAT = (data+tag+valid)\* # registers  
min size of tag storage = tag size \* # registers  
+ tag size \* #RS per FU \* #FU \* # registers

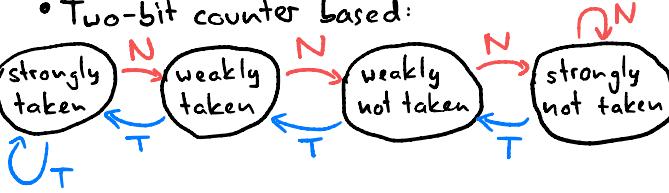
## BRANCH PREDICTION

Static prediction direction schemes:

- Always not-taken
- Always taken
- Backward taken, forward not-taken

Dynamic prediction direction schemes:

- Last-time: single bit per branch, stored in branch target buffer (BTB)
- Two-bit counter based:



- Local: PC-based, for each branch
- Global: counter shared among all branches
- Two-level global:



- ⇒ associate branch outcomes with "global T/NT history" of all branches
- Two-level local: have a per-branch history, choose based on last time with the same local branch history

Example: pipeline w/ multiple stages, each stage completes in 1 cycle. stalls only when conditional branch instruction. Consider: 4

LOOP1: SUB R1, R1, #1 // R1=R1-1  
BGT R1, LOOP1 // Branch if R1 > 0

LOOP2: B LOOP2 // Branch until killed

The microbenchmark has following results:

Initial R1 value	Number of cycles taken
4	51
8	63
16	87

⇒  $C = P + I - 1 + B \cdot D$  # of cycles stalled per conditional branch  
total # of cycles # pipeline stages # of conditional branch instructions  
total # of dynamic instrs. executed  $\Rightarrow 51 = P + I - 1 + 4 \cdot D$   
 $63 = P + I - 1 + 8 \cdot D$   
 $87 = P + I - 1 + 16 \cdot D$

stall per cond. branch: D=3 cycles  
# pipeline stages + # dynamic instructions: P+I=40  
Example: for(i=0; i<10<sup>6</sup>; i++) // B1  
if(i%3=0) // B2

⇒ 2-bit GHR, PHT with 2-bit saturating counter. Initial values for 100% mispredict in first 5 iterations:

PHT Entry	Value
TT	01
TN	00
NT	01
NN	NN

Example: for(i=0; i<N; i++) // B1 (has local correlation)  
if(a[i] % 2=0) // B2 } globally correlated  
if(a[i] % 3=0) // B3 } B4 taken  
if(a[i] % 6=0) // B4 } ⇐ B2 & B3 +.

2-bit GHR, PHT with 4 entries: TT, TN, NT, NN

PHTE inc. when taken. Calc. PHTE of TT after 120 it.

B3: B1-T & B2-T =  $\frac{1}{2}$  ⇒ INC =  $\frac{1}{2} \cdot \frac{1}{3} = \frac{1}{6}$ , DEC =  $\frac{1}{2} \cdot \frac{2}{3} = \frac{1}{3}$

B4: B2-T & B3-T =  $\frac{1}{6}$  ⇒ INC =  $\frac{1}{6} \cdot 1 = \frac{1}{6}$ , DEC = 0

B1: B3-T & B4-T =  $\frac{1}{6}$  ⇒ INC =  $\frac{1}{6} \cdot 1 = \frac{1}{6}$ , DEC = 0

B2: B4-T & B1-T =  $\frac{1}{6}$  ⇒ INC =  $\frac{1}{6} \cdot \frac{1}{2} = \frac{1}{12}$ , DEC =  $\frac{1}{6} \cdot \frac{1}{2} = \frac{1}{12}$

chance of B2-T if B4-T and B1-T

Total:  $(\frac{1}{6} - \frac{1}{3}) + \frac{1}{6} + \frac{1}{6} + (\frac{1}{12} - \frac{1}{12}) = \frac{1}{6} \Rightarrow 120 \cdot \frac{1}{6} = \underline{\underline{20}}$

## SYSTOLIC ARRAYS

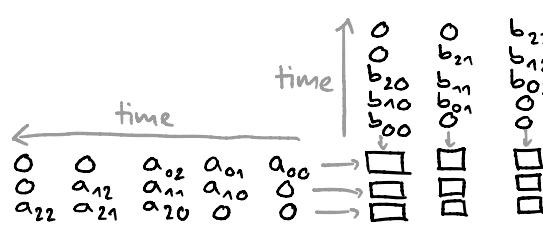
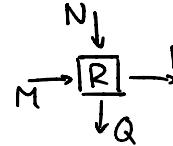
⇒ Array of PEs instead of single PE

Example matrix multiplication:

$$\begin{bmatrix} C_{00} & C_{01} & C_{02} \\ C_{10} & C_{11} & C_{12} \\ C_{20} & C_{21} & C_{22} \end{bmatrix} = \begin{bmatrix} A_{00} & A_{01} & A_{02} \\ A_{10} & A_{11} & A_{12} \\ A_{20} & A_{21} & A_{22} \end{bmatrix} \times \begin{bmatrix} B_{00} & B_{01} & B_{02} \\ B_{10} & B_{11} & B_{12} \\ B_{20} & B_{21} & B_{22} \end{bmatrix}$$

Processing Element:

$$P = M, Q = N, R = R + (M \cdot N)$$



## VECTOR PROCESSING

**Memory banks:** split up memory, so many data elements can be fetched concurrently.

**Stride:** constant distance that separates elements that are adjacent in a vector register.

The data for vector instructions is stored in vector registers, which holds multiple "vectors".

Vector chaining is data forwarding with vectors.

It allows processing of individual elements as soon as they are ready.

Example:  $i=0 \dots 49$   $C[i] = (A[i] + B[i]) / 2$

MOVI VLEN = 50

1

MOVI VSTR = 1

1

VLD VO = A

11 + VLEN - 1

VLD V1 = B

11 + VLEN - 1

VADD V2 = VO + V1

4 + VLEN - 1

VSHFR V3 = V2 >> 1

1 + VLEN - 1

VST C = V3

11 + VLEN - 1

Without chaining: (285 cycles)

1111111 49 1111 49 111 49 111 49 111 49 1  
 $VO = A[0 \dots 49]$   $VI = B[0 \dots 49]$  ADD 1 SHIFT 1 STORE 1

With chaining: (182 cycles)

1111111 49 1111 49 111 49 111 49 111 49 1  
 $141 \quad 49 \quad 1$   
 $111 \quad 49 \quad 1$   
 $111 \quad 49 \quad 1$   
 $111 \quad 49 \quad 1$

Assuming each memory bank has a single port

If # data elements > vector length: break loops and execute multiple times with (different) VLEN

# banks should be at least equal to the number of cycles of VLD and VST if stride=1. Stride and # banks should be relatively prime if stride > 1.

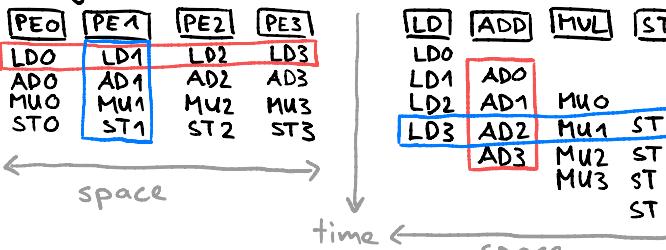
## SIMD PROCESSING

**Warp:** a set of threads that execute the same instr. at same PC on different data

$$\Rightarrow \# \text{ warps} = (\# \text{ threads}) / (\# \text{ threads per warp})$$

**SIMD utilization:** fraction of SIMD lanes executing a useful operation

Array Processor vs. Vector Processor



■ Same OP at same time  
■ different OP at same space

■ Same OP at same space  
■ different OP at same time

## Examples

Tracing the cache: (fully-assoc. / 4 blocks)

Req A B A H B G H H A E H D ...

M7 x x x x x x x

MRU A B A H B G H H A E H D

: A B A H B G H H A E H

LRU A A A B G G A

$$\Rightarrow \text{Cache miss rate} = \frac{\# \text{M}}{\# \text{Req}}$$

## Cache Reverse Engineering

Addresses Accessed Hit Rate

0 16 24 25 1k 255 1100 305 | 2/8

31 65k 65k 131k 262k 8 305 1060 | 3/8

262k 65k 4 | 2/3

• Block size: seq1: 8B: {0}, {16}, {24}, {25}, {255}, ...  
16B: {0}, {16}, {24}, {25}, {255}, {305}, ... ⇒ 2 hits!

• Associativity: seq2: 3 hits w/ block size 16B  
⇒ 31, 65k, 305 hit ⇒ 8 must miss, replaced by blocks mapping to set 0 (65k, 131k, 262k) ⇒ 2-way

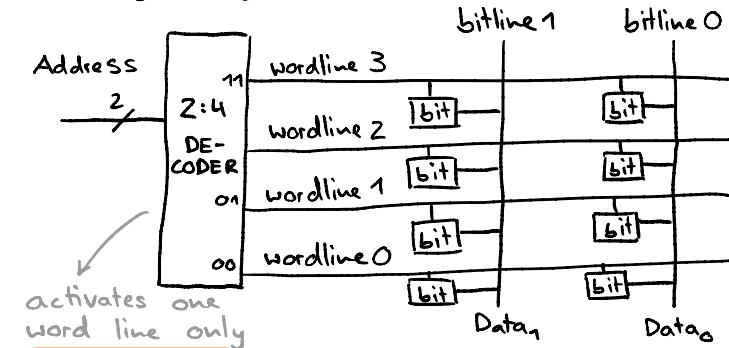
• Replacement: seq3: with LRU only 262k would hit ⇒ with FIFO, 4 also hits

• Cache size: we know block size 16B, 2-way.  
⇒ 4 bits for indexing byte in block

⇒ 7 bits if size=4KB, 8 bits if size=8KB  
⇒ 2048 maps to set 0 ⇒ size is 4KB

# CACHE & MEMORY

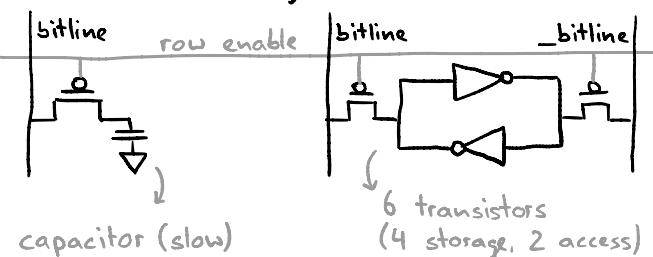
Memory Array:



**Memory Bank:** divides a large array into multiple banks which can be accessed indep. Has rows (many) and columns (fewer) that have distinct addresses. A row buffer stores the most recently accessed row (simple cache)

## DRAM vs. SRAM

- slow, but cheap
  - requires refresh
  - used in main memory
- fast, but expensive
  - no refresh required
  - used in caches



Memory locality:

- temporal: program accesses same memory location many times in small window of time
- spatial: program accesses nearby memory location in a window of time

Memory is divided into blocks: **Data (64B)** **Tag**

The block address maps to a potential address in cache (determined by index bits):

8-bit-address: **2b 3bits 3bits**  
tag index byte in block

If the stored tag is valid & matches the tag of the requested block: cache hit.



Associativity:

- fully-associative: every memory chunk can be mapped to any cache block
  - ↳ **5 bits | 3 bits** → no index, only 5bit tag
    - ↳ location needs to be calculated w/ logic
- direct-mapped: a memory chunk can only be mapped to a single cache block
- set-associative: allows same index multiple times (n-way associative)
  - ↳ 2-way-associative: **3 bits | 2b | 3 bits**
    - ↳ increased tag index byte in block
    - ↳ points to a set instead of single block

Replacement Policy:

- FIFO: first in, first out
  - LRU: least recently used
  - Not MRU: not most recently used
- ↳ for every program a different replacement strategy might be optimal

Cache performance:

- cache size (total data capacity excl. tags)
- block size (data associated w/ address tag)
  - ↳ **b = 2<sup>offset</sup>**
- associativity (how many blocks can be represented in same index/set?)
- number of blocks: **B = C / b**
- number of sets: **S = B / N**
- tag store size:  $2^{\text{index}} \cdot (\text{associativity}) \cdot (\text{tag+valid+dirty}) + [\text{LRU}]$
- data store size:  $2^{\text{index}} \cdot [\text{associativity}] \cdot [\text{block size}]$

# PREFETCHING

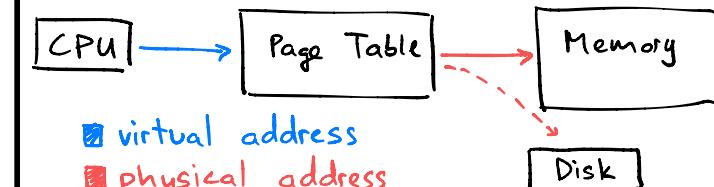
Performance: • Accuracy: **used sent** prefetched

- Coverage: prefetched misses / all misses
- Timeliness: **on-time used** prefetched

Runahead execution: with small instruction window, pre-execute instructions if data has to be fetched from memory → prefetch

Stride prefetches: if there is a constant stride between memory accesses → use to predict

# VIRTUAL MEMORY



**Virtual Address:** **Page table number | Page table offset | Page offset**  
↳ uses 2 levels of page tables

**Page table entry:** **V Physical Page Number**  
↳ page in virt. memory is block in cache

**Physical Address:** **Physical page number | Page offset**  
↳ same

Translation Lookaside Buffer (TLB):

Small cache of recently used PTEs (i.e. recently used Virtual-to-Physical translations)

Page faults: if a page is not in physical memory but on disk: which to replace?  
→ clock page replacement: pointer to last accessed frame, set R bit in PTE if accessed. when replace, look for first PTE with R bit not set (and unset the traversed PTEs)

Physical address space:  $2^{[\text{page offset}] + [\text{physical tag}]}$