



```

module SmilingSmile(input clk, reset, output y);
reg [1:0] state, nextstate;
parameter SO = '1'b0;
parameter S1 = '1'b1;
always @ (posedge clk, posedge reset)
begin
    if (reset) state <= SO;
    else state <= nextstate;
end
always @ (*) begin
    case(state)
        0: if ((a & b) == 1'b0) nextstate = S1;
        1: if ((a & b) == 1'b1) nextstate = SO;
        2: if ((a & b) == 1'b1) nextstate = S0;
        3: if ((a & b) == 1'b0) nextstate = S1;
    end
end
endmodule
//output logic (mealy dep.)
assign y = a ? state == S3 : assign y = en ? a : 4'bzz;
endmodule // Mealy dep. end module

```

```

always @ (posedge clk) begin
    alb1c8d = al(b^c & d);
    $int2, $int1, $int03 = $1'000; NAND: ~& NOR: ~^ XNOR: ~^
    (Can be vectorised) 5
    4'b1122: int2 = 1'b1;
    4'b0112: int1 = 1'b1;
    Behavioral models: Define logic behavior
    default: $int2, $int1, $int03 = 3'10; Structural models: Define interactions
    endcase // $c1 dododo co 101 of instanced submodules
    assign y = {c[2:1], f3'd[0]3, c[0], 3'b101}; (Synthesizable)
    Reduction operators are bitwise: &a |a ~a ^a
    decoder hi #(f(c)(g)); module decoder #(parameter N=3) (in-

```

```

Testbenches (DUT = Device under test)
module testbench();
example.tv
1-0
1-0
:
endmodule

```

```

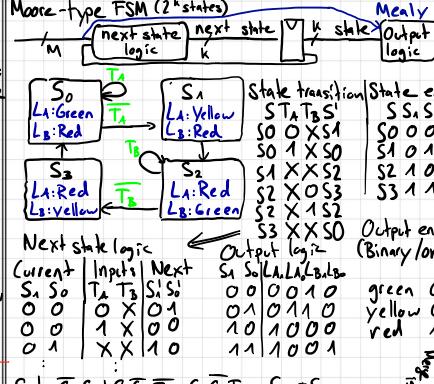
reg [a:b:c:wire y];
$function dut(a,b,c);
reg [2:0] test vectors [100000:0];
initial begin
    a=0; b=c=0; #10; // indicates a delay
    c=1; #10; // old values will stay
    if(y != 0) $display("001 failed");
    $readmem("example.tv", testvectors);
end
always @ (posedge clk) begin
    #1; ea, bc3 = test vectors [ea0];
end
endmodule

```

Variables modified in initial blocks must be registers (2D register arrays)

1 timescale 1ns/1ps: (Units are 1ns and precision of sim. is 1ps)

Simplicity favors regularity = Large funct.  $\Rightarrow$  Multiple instr. Make the common case fast = Small instr. set instr. Good design demands good compromises Smaller is faster



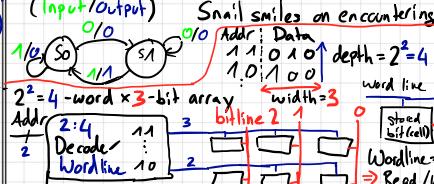
Next state logic

Divide-by-N counter

$$S_1 = \overline{S}_1 S_o + S_1 \overline{S}_o \overline{T}_B + S_o T_B = S_1 \oplus S_o$$

$$S_0 \rightarrow S_1 \rightarrow \dots \rightarrow S_{n-1}$$

Moore: Output of node is shown when moving away.  
Mealy: Output of edge as seen as it moves.



Bitline = 2  $\Rightarrow$  Read Bitline = (0/1)  $\Rightarrow$  Write 0/1

Three-portled memory: [A1]  $\Rightarrow$  RD1, [A2]  $\Rightarrow$  RD2,

RAM (Random access memory): Volatile = On, when power

ROM (Read-only memory): Nonvolatile = Keep data

DRAM (Dynamic RAM): Store data using capacitor charge

SRAM (static RAM): Store data using cross-coupled inverters

refresh

DRAM: SRAM:

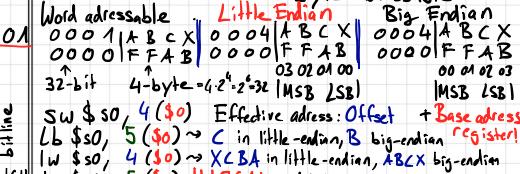
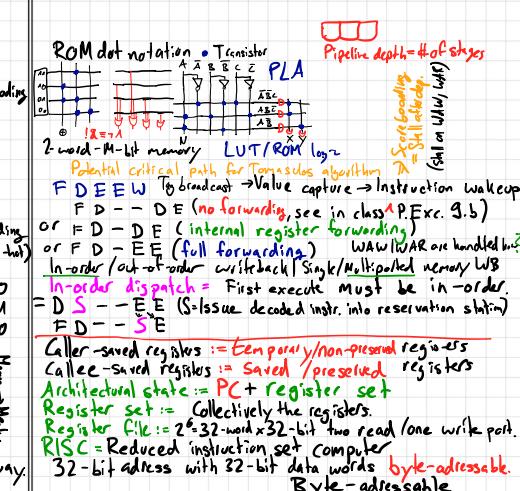
$\xrightarrow{+t_{on}}$   $\xrightarrow{+t_{off}}$  can leak current (I<sub>leak</sub>)

FPGA (Field programmable gate array)

Consists of: (LB/configurational logic block) = LUT + registers

FF (super expensive)

Multiplexers



OxCAFEBABE Little-Endian: 3a 2b fe ca  
Big-Endian: ca fe 2b 3a  
Little Endian: Kleinstre Stelle an Wichtigste Memory Adresse!

[C-3b] Memory model 2<sup>16</sup>; 16-bit (2 byte) addressability  
2<sup>16</sup> word alignment

0 0 0 0   C D   LDW DR, Base, offset 6	0 0 0 2   A B   DR=Mem[BaseR+LSHF(SEXT(offset6), 1)]
0 0 0 4   F	0 1 0 0   LEA ... DR=PC + LSHF(SEXT(PC, offset1), 1)]

Only loads effective address, not memory there!

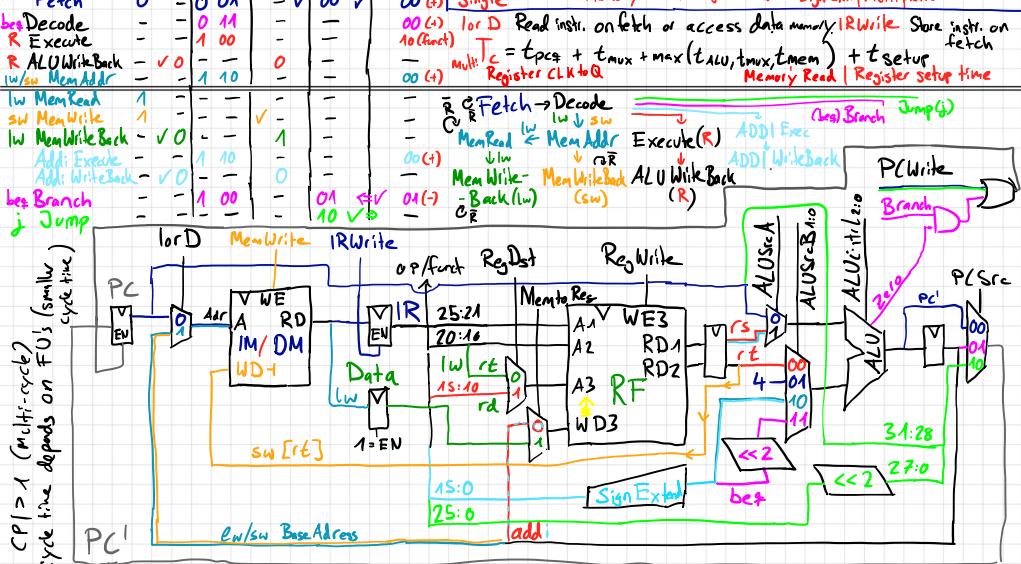
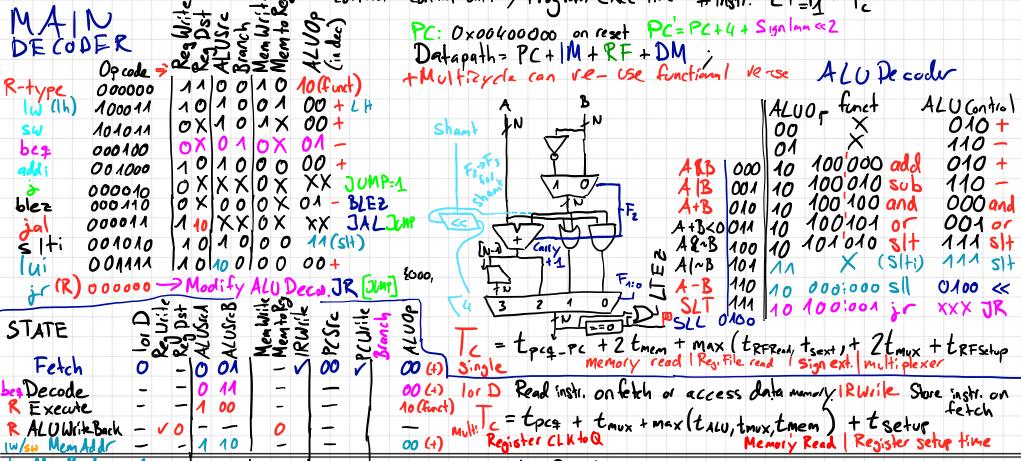
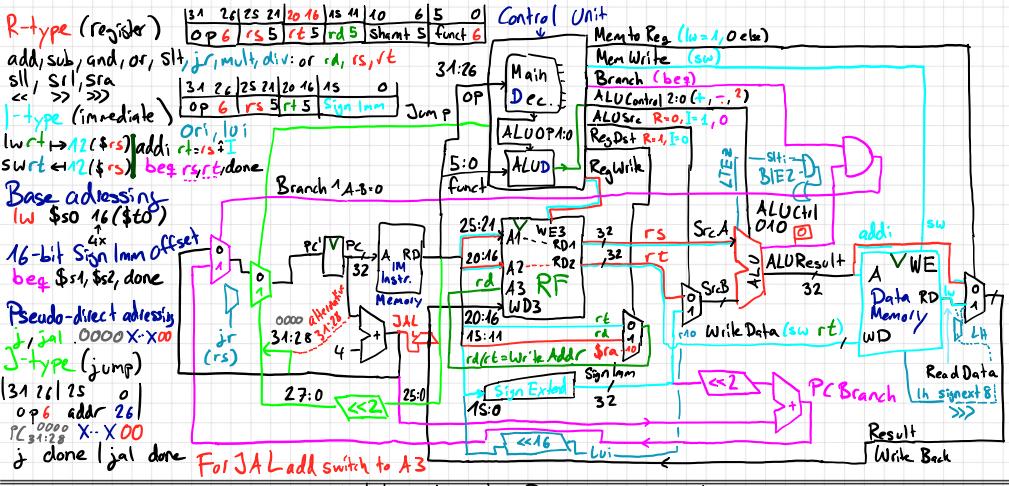
Loading: Enable / disable input to be written to a register

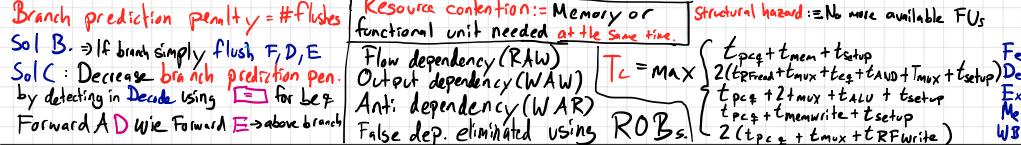
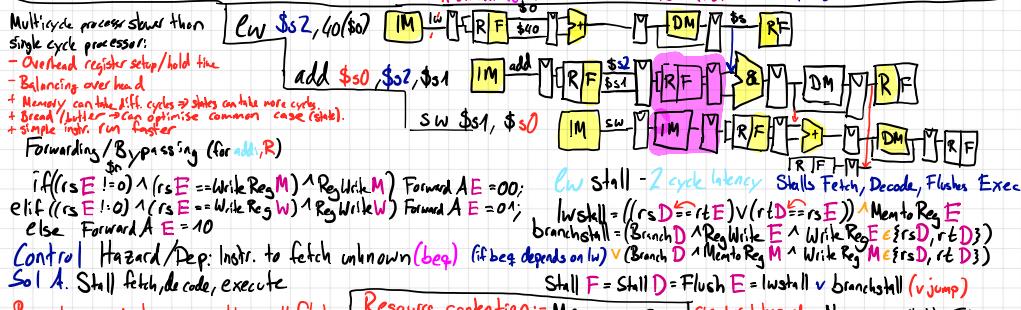
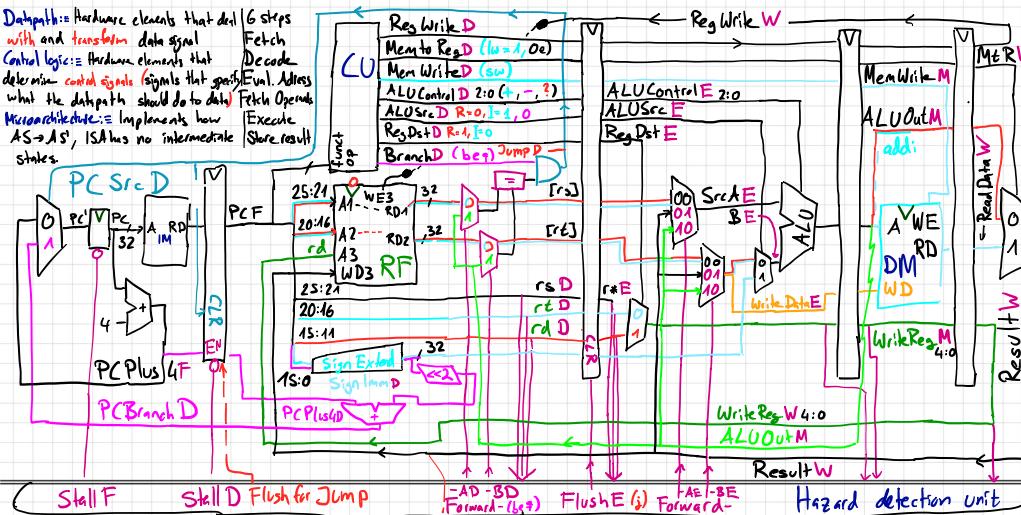
Gating: Enable / disable input to bus using tristates

General purpose registers: 8x16 bit registers

Condition codes: P (written reg is positive), N (negative), Z (zero)

Memory mapped I/O XFFEOO-XFFF for I/O





**Handle flow dependencies:** Detect & Wait (stalling), Detect & Bypass (Forwarding), Detect & Eliminate dep. on software level (inserting NOPs)

- Predict value & verify (flush) - Do something else (threading)
- Inserting NOPs = #instructions, Inserting bubbles = #CPUs
- Out-of-order execution (Execute later ind. code in same thread)

**Control dependency handling**

A. Solve by stalling/flushing

B. Early branch resolution

MIPS uses [ ] for beq in Decode

- Introduces another data dep.

C. Dynamic branch resolution

-> 2-bit branch predictor

D. Cons Pipelining:

- Resource contention/Structural hazard

- External fragmentation (Some stages not required for all instr.)

- Internal fragmentation (Not all stages have the same speed)

- Skills add HW, PC, CP code dep.

Dependence detection mechanism

Scoreboarding (Valid bit on each reg)

+ Simple - Still not only RAW

Combinational dep. check logic

+ Fast (only RAW)

- More complicated esp. on Superscalar machines.

**ROB** ≥ # registers + max stage count = max size of instr. window.

**Indirection (Alternative impl.)**

1. Access RF, if valid bit = 1 → OK, else load/overwrite pointer of where it will land in ROB.

2. Check if ROB is valid, else stall and wait.

- Increases latency (for ROB-to-RF)

3. Finished computation, broadcast ROB 106.

+ Simple - Adds another stage (incr. latency)

**Memory dependencies are dynamic**

**MOB** is equivalent concept to ROB.

- Huge, cannot solve all deps.

- Additionally needs to know size of access.

(lw \$t0, \$0(\$s0) Data dep / hazards)

add \$t1, \$t0, \$s1 RAW (flow dep)

sub \$t2, \$s2, \$s3 WAR (antidep.)

add \$t3, \$s4, \$s5 WAW (output dep.)

The last instr. can be squashed.

**Instruction window** = All decoded but not yet retired instr.

## Out-of-order execution (Data-flow order)

Dep. instr.  $\rightarrow$  moved away into reservation stations



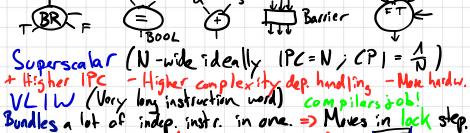
- Link consumer op to value of producer
- Buffer instructions until they are ready in reservation station
- Keep track of readiness (V-bit) = Broadcasting (compare)
- Dispatch to functional unit (FU)  $\rightarrow$  Wake up and select instr
- #resv. statn entries = (usually  $K = \#FU$ )

$\text{tag size} = \lceil \log_2(\#\text{resv. statn entries}) + 2 \rceil$   
 $\#\text{tag comparators} = \#\text{resv. statn entries} \cdot \#\text{source reg.} \cdot (\#\text{FU} + \#\text{RO} + \#\text{reg.})$

**Re-Order** (makes it Von-Neumann, used for except.)  
After execution update the RAT (frontend register file)  
After retirement update the architectural reg. file  
= offset instr. in machine  $\rightarrow$  executed  $\rightarrow$  updated in program order  
On exception: Flush pipeline  $\wedge$  copy architectural RF  $\rightarrow$  frontend RF  
+ OoOE tolerates latency of single ops.  
Only so much can be tolerated, size of instr. window depends on scheduling window and RAT | ROB size.

Data flow at ISA level = unsuccessful

- + Good for irregular parallelism (hard to spot dependencies).
- + Removed PC (instr. counter) dependency
- Bad for debugging  $\rightarrow$  Not just sequential state.



Superscalar ( $N$ -wide ideally  $\text{IPC} = N$ ;  $\text{CPI} = \frac{1}{N}$ )

+ Higher IPC - Higher complexity dep. handling - More hardware

VLIW (Very long instruction word) compiler job!  
Bundles a lot of indep. instr. in one  $\Rightarrow$  Moves in lock step.

+ Simple (no sched / no dep. checking - no renaming)  $\rightarrow$  compiler job!

- Compiler needs to find  $N$  indep. instr. Problem:  $\xrightarrow{\text{lock step problem}}$

- Code size increases (NOPs)  $\rightarrow$  Lock step problem

- Lower L1-cache hit-rate / Higher fetch memory bandwidth

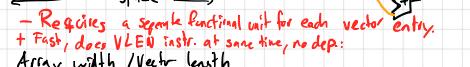
- Change width = Recompile (change of ISA)

+ Compilers for VLIW can be used for superscalar

SIMD processing (Single instruction multiple data)

For e.g. dot product on 2 vectors. Modern = Array -  $\&$  Vector proc.

Array processor : SIMD at the same time using different spaces

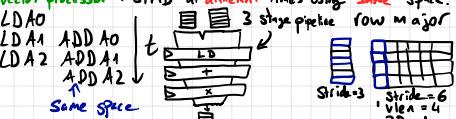


- Requires a separate functional unit for each vector entry

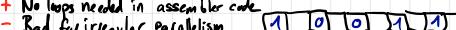
+ Fast, does VLEN instr. at same time, no dep.

Array width / vector length

Vector processor : SIMD at different times using same space.



- + A single vector instr. has no intra dependencies
- + No loop needed in assembly code
- Bad for irregular parallelism
- Slow for single-part memory



Memory banks VLEN register VMASK = 100111 VSTRIDE = 2

Word-interleaving = Row buffering N parallel accesses for

stride = 1 relatively prime to banks  $\rightarrow$  MDR  $\xrightarrow{\text{32}} \text{Bank 1}$   
MDR  $\xrightarrow{\text{32}} \text{Bank 2}$   
VLEN  $\rightarrow$  should use 64 banks.

If load latency is 50 cycles should use 64 banks.

- Row buffer conflicts at single bank  $\rightarrow$  Single value

$\rightarrow$  Resolve with more banks, random addressing, better mem. scheduling

Dynamic instructions = # Instr. executed by machine.

Static instructions = # Instr. written down in program

Vector chaining = Cannot chain loads/stores to the same memory bank!

VLEN = 32+8 = 40 Chaining & Banking B = 32 = Bank size

VLD A Load latency + 31 + 2 posted

VLD B Hit load latency + 31 + Hit latency + 7 +

VADD V2 S0,V1 F0 + 31 + 7 + 7 +

VMUL V3 S1,V0 F0 + 16 + 24 + 16 + 7 +

Banking 1/0 chaining

VLD A  $\xrightarrow{\text{LL}} + 12 -$

VLD B  $\xrightarrow{\text{LL}} + 12 -$

Chain VADD  $\xrightarrow{\text{LL}} + 4 + 21 -$

No chain VADD  $\xrightarrow{\text{LL}} + 4 + 21 -$

30 cycles incl. fetch

General banking B = bank size

VLD VAC4 VLD VAC4

$\xrightarrow{\text{LL}} + 3 - + (\text{K}+1) - \text{LL} + C - 1$

$\xrightarrow{\text{LL}} + 3 - + \text{LL} + C - 1$

$\xrightarrow{\text{LL}} + 3 - + \text{LL} + C - 1$

30 cycles incl. fetch

Density-time: Only load masked. Tag[1] Value Tag[1] Value

RAT R0 Y T0 13

R0 X T0 13

R1 0 A 8

B. mul A 0 - - - 7

## Flynn's taxonomy

SISD (Scalar sequential): Pipeline / Superscalar / OoO / VLIW

SIMD (Data parallel): Array / Vector processor

MISD: Systemic arrays / streaming processors SIMD: Multicore Superscalar

GPU: SIMD - Single instr., Multiple threads

Warp: Set of threads that exec. same instr. but with diff. data

SIMD vs. SIMD

[LD, ADD, ST], VLEN  $\rightarrow$  [LD, ADD, ST], # of threads

$\downarrow$  Vector instrs.  $\uparrow$  Scalar instrs.

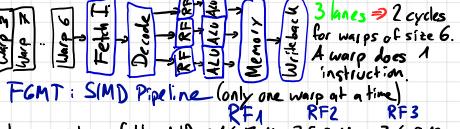
Sequential instr. on ISA level Threads are grouped dynamically

+ Can treat each thread separately

$\rightarrow$  Can be executed on a sequential machine if desired.

$\rightarrow$  leads to SIMD processing

+ Granularity on the level of warps



Contain registers of thread IDs: 1,4,7,10,2,5,8,11,3,6,9,12

**Warp Instruction level parallelism:** Can issue 1 warp / cycle. Different functional units: (load, multiply, add)

Example: 8 threads per warp, 4 lanes, 3 functional units  $\downarrow$  load



Assume answers must be

Branch divergence  $\rightarrow$  w/o branch divergence!

Every thread can take another control flow path, warps are hypothetically regrouped on every instr. Spacers are NOP d. to act only after branch divergence.

Works only on same SIMD lang if separate RF.

## Systolic arrays

+ Specialised (good efficiency, high concurrency, simple design)

- Less memory bandwidth / balanced computation

- Specialised (not generally applicable, more hardware)

Convolution (SFC)(g-a-T) = S(fa-T)g(r-T)

PE = Processing element for e.g. WARP

Memory  $\xrightarrow{\text{Yres}, \text{Xres}}$   $\xrightarrow{\text{Yout}, \text{Xout}}$   $\text{Y}_i = \text{W}_i \cdot \text{X}_i + \text{W}_{i+1} \cdot \text{X}_{i+1} + \dots + \text{W}_k \cdot \text{X}_{k+i-1}$

Cycles 1 2 3 4 5 6 7  $\rightarrow$  Can optimise loop iterations

Issue X:  $\xrightarrow{\text{X}_1, \text{X}_2, \text{X}_3}$

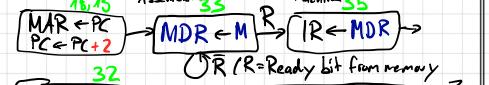
Issue Y:  $\xrightarrow{\text{Y}_1, \text{Y}_2, \text{Y}_3}$

Result Y:  $\xrightarrow{\text{Y}_1, \text{Y}_2, \text{Y}_3}$  loop { A  $\rightarrow$  B  $\rightarrow$  C }

## Decoupled Access/Execute Processors (DAE)

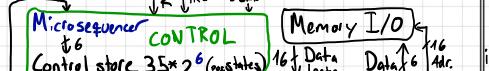
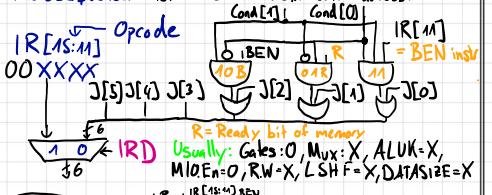
- Fetch: Split into A-instr and E-instr part
- Write/Read memory done over AEQ | EAQ
- Also 2 RFs which can be interchanged with Queues
- + Queues lower needed number of registers
- + Access/Execute can run ahead of each other.
- Branches require synchronization of A/E.
- 2 instr. streams (could be done in 1 though)
- Loop Unrolling** + Good if you have bad loop prediction.
- + Enlarges basic block  $\Rightarrow$  Enables code optimization/scheduling
- Needs more code, also check if not multiple of unroll factor

**LC-3b** (critical path: 5 stages + 2 mem. stages + mem. latency  $\downarrow$  control path  
 Pseudoinstructions  $\rightarrow$  Instructions  $\rightarrow$  Microinstructions  
 Assembler 33 Machine 35



$[R[15:11]$  Opcode,  $IR[11] = 1$  if JSR(label), 0 if JSR(ic)  
 $B R np$  Branch if CC neg. or pos.  $P_{rel}$  on branch

**Microinstructions:** Control signals of data path  $\uparrow$  to  
 $\uparrow$  are stored in the control store determining the next state.  
 Microsequencer determines the next state address:  
 $Cond[11] \quad Cond[0]$



**Single-bus design:** + Low hardware - Reduced concurrency  
 SIMD Utilization =  $\frac{\text{Total # of possible ops at a time}}{\text{# useful ops of a warp}}$

Loop hot 32 threads, code:  
 $i: A[i] > 0 \quad 4 \text{ ops} \quad 3$   
 $\text{Entweder } \frac{32}{32} / \frac{160}{160} = \frac{32+4}{160} = 4 + (8+a) = 8+a$   
 $\cdot 8+a = \frac{4}{40} \Rightarrow 8+a = 8 \Rightarrow a=0 \quad \& [1:3] \Rightarrow \text{impossible}$   
 $a \text{ in every } 32 \text{ of } 8 \text{ elements is less than } 0$

## LC-3b

- + Simple design for powerful computation (hidden p-ISA)
- + Easy extensibility of the ISA
- + Portable Microcode
- Cache, TLB SRAM DRAM  $\sim 10\text{Gb}$   $\sim 1\text{ns}$
- MM = Physical memory, Page T. DRAM  $\sim 100\text{Gb}$   $\sim 10\text{ns}$
- Rest of virtual memory Hard Disk  $\sim 1\text{TB}$   $\sim 100\text{ms}$
- Miss rate = # Misses / # memory accesses =  $1 - \text{HitRate}$
- Hit rate = # Hits / # memory accesses =  $1 - \text{MissRate}$

**AMAT (Average Mem. Access time):**  
 $t_{cache} + t_{MR \text{ cache}} (t_{mem} + t_{VM})$

**Spatial locality:** When accessing data, bring data nearby.  
**Temporal locality:** Keep recently accessed data in higher levels of cache.

$$\begin{aligned} C &= \text{Capacity} (\# \text{ of bytes}) \\ b &= \text{Block size} (\text{bytes of data brought in at once}) \\ B &= C/b = \# \text{Blocks in cache} \approx \# \text{Cache lines} \\ N &= \# \text{of blocks in a set} (\text{Degree of associativity}) \\ S &= B/N = \# \text{of sets} \approx \# \text{of cache rows} \end{aligned}$$

$$\begin{aligned} \text{Tag store size N-way cache: } S &= N(C + 1 + 1) + (\log_2(N)(N-1)) \\ \text{Data store size N-way cache: } S &= N \cdot b \text{ (bytes)} \end{aligned}$$

**Compulsory miss:** First time accessing this particular memory.

**Conflict miss:** Won't happen in a fully assoc. cache with LRU rep.

**Capacity miss (victim):** Even happen in fully assoc. cache with LRU rep.

**LRU replacement:** Oldest / MRU replacement = Newest

**Direct Mapped Cache:**  $(N=1=L, S=B=2^3)$

Memory address 2<sup>3</sup> cache:  $6 \text{ Tag: Set: } 000 + 8 \text{ byte offset}$

**N-way set associative cache** ( $1 \leq N \leq B$ )  $N=2, S=B/N=2^3=8$



$\{[\text{TagWay1} == \text{TagWay2} \text{ Valid of Way1}] \mid [\text{TagWay2} == \text{TagWay1} \text{ Valid of Way2}] \}$

**Additioally needs a multiplexer to choose data!**  
**Fully associative cache** ( $N=B, S=1$ )

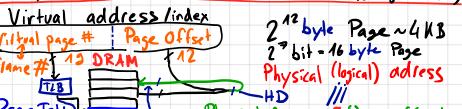
Usually small due to the many comparisons!

Example: Direct mapped cache Block size =  $2^2 = 4$  (offset)

$S=2^4=16 = B$



**0 ② ④ 8 16 32** has 2 cache hits with Block size 8  
 $0 \quad 8 \quad 16 \quad 24 \quad 16 \quad 8 \quad 0 \quad B=16, S=2$   
 temporal locality  $\uparrow$  = Set  $\uparrow$   
 spatial locality  $\uparrow$  = block size  
 Set diff. index  $\downarrow$  = N (degree of associativity)



**Page Table:** Page Table links chunks (size of offset) 0-4096-4096-8192

**DRAM:** Frame # Frame content - Check T-LB  $\Rightarrow$  Cache hit  
**LRU - Policy:**  $\begin{cases} 0 & \text{if } 0 \\ 1 & \text{if } 1 \\ 2 & \text{else} \end{cases}$  Check page table for valid  
 $\begin{cases} 0 & \text{if } 0 \\ 1 & \text{if } 1 \\ 2 & \text{Not valid} \Rightarrow \text{page fault} \end{cases}$

**Translation Lookaside Buffer (TLB):** fully associative

Physical frame # & Page offset  
 Physical tag = size of tag of cache

Page table