

Calculating performance

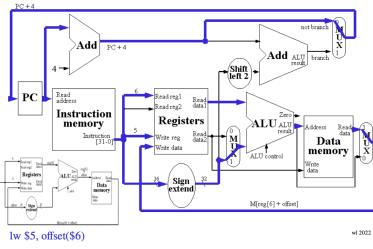
- CPI: average clock cycles per instruction
- number of cycles for program P = number of instr. for P \times CPI
- execution time for P = clock cycle time \times number of cycles for P
- $= \frac{1}{\text{clock speed}} \times \text{number of cycles for P}$
- average exe. time for P₁, P₂..P_n = $\frac{1}{n} (\text{exe. time for P}_1 + \dots + \text{exe. time for P}_n)$ (assume equal workload)

execution time equation:

$$\text{exe. time} = \text{instr. count} \times \text{CPI} \times \text{cycle time}$$

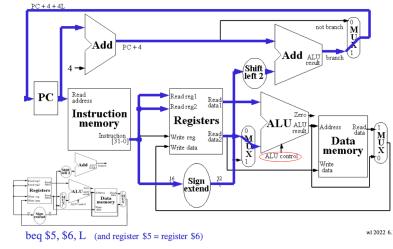
wl 2023.2.3

Combined datapath for load



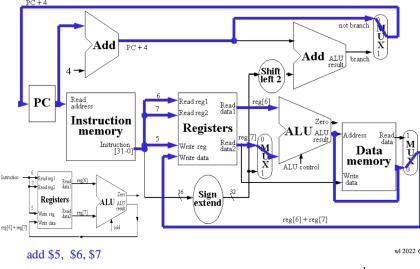
wl 2023.6.13

Combined datapath for branch



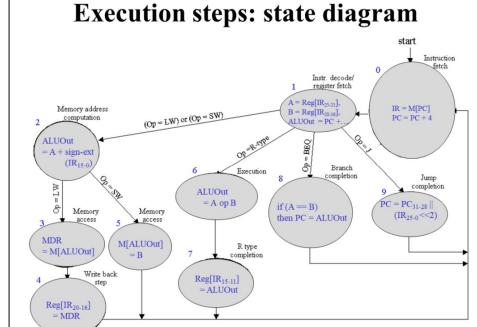
wl 2023.6.14

Combined datapath for R-type

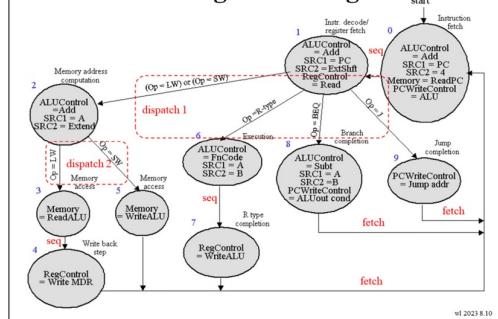


wl 2023.6.12

Execution steps: state diagram

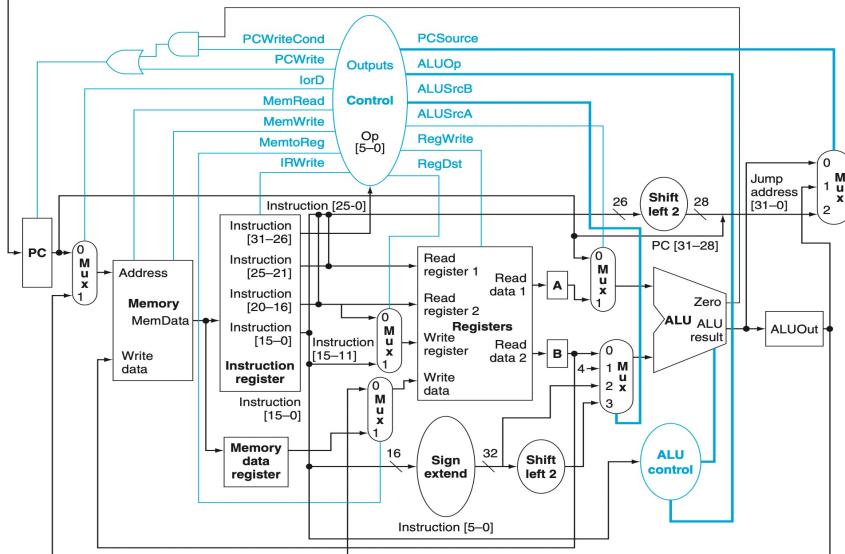


Field assignment diagram



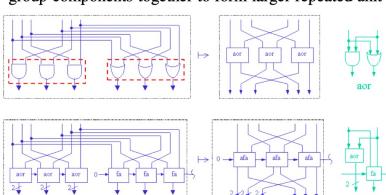
wl 2023.8.10

Instruction	RegDst	ALUSrc	Mem-to-Reg	Reg-Write	Mem-Read	Mem-Write	Branch	ALUOp1	ALUOp0
R-format	1	0	0	1	0	0	0	1	0
lw	0	1	1	1	1	0	0	0	0
sw	X	1	X	0	0	1	0	0	0
beg	X	0	X	0	0	0	1	0	1



Deriving ALU cell by interleaving components

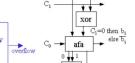
- group components together to form larger repeated unit



- the dotted boxes have the same function and interface

Selecting ALU operation

- programmable inverter for b_i (using xor)
- connecting mux in series
- d_{0,1}: 00 and, 01 or, 10 add, 11 subtract
- detecting overflow: exercise



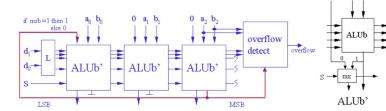
wl 2023.4.9

Comparison operations

- slt: set on less than, if a < b then 1 else 0
- if a < b, a-b < 0, so MSB of (a-b) is 1 (32 bits)

Implementation

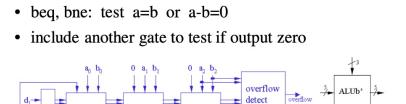
- provide additional input to each cell, left of a_b, b_i
- LSB input from MSB ALUB output, other inputs set to 0
- include additional mux in cell for selection
- to select s, 0, d, 1, d, 1



wl 2023.4.10

Zero detection

- beq, bne: test a=b or a-b=0
- include another gate to test if output zero
- summary: s d_{0,1} d_{2,3} function: set on less than



wl 2023.4.11

Instructions

Data movement

`movq Src, Dest` Dest = Src
`movsbq Src, Dest` Dest (quad) = Src (byte), sign-extend
`movzbq Src, Dest` Dest (quad) = Src (byte), zero-extend

Conditional move

`cmove Src, Dest` Equal / zero
`cmoveq Src, Dest` Not equal / not zero
`cmoveq Src, Dest` Negative
`cmoveq Src, Dest` Nonnegative
`cmoveq Src, Dest` Greater (signed >)
`cmoveq Src, Dest` Greater or equal (signed \geq)
`cmoveq Src, Dest` Less (signed <)
`cmoveq Src, Dest` Less or equal (signed \leq)
`cmoveq Src, Dest` Above (unsigned >)
`cmoveq Src, Dest` Above or equal (unsigned \geq)
`cmoveq Src, Dest` Below (unsigned <)
`cmoveq Src, Dest` Below or equal (unsigned \leq)

Control transfer

`cmpq Src2, Src1` Sets CCs Src1 Src2
`testq Src2, Src1` Sets CCs Src1 & Src2
`jmp label` jump
`je label` jump equal
`jne label` jump not equal
`js label` jump negative
`jns label` jump non-negative
`jg label` jump greater (signed >)
`jge label` jump greater or equal (signed \geq)
`jl label` jump less (signed <)
`jle label` jump less or equal (signed \leq)
`ja label` jump above (unsigned >)
`jb label` jump below (unsigned <)
`pushq Src` %rsp = %rsp, %rsp[8] = Src
`popq Dest` Dest = Mem[%rsp], %rsp = %rsp + 8
`call label` push address of next instruction, jmp label
`ret` %rip = Mem[%rsp], %rsp = %rsp + 8

Arithmetic operations

<code>leaq Src, Dest</code>	Dest = address of Src
<code>incq Dest</code>	Dest = Dest + 1
<code>decq Dest</code>	Dest = Dest - 1
<code>addq Src, Dest</code>	Dest = Dest + Src
<code>subq Src, Dest</code>	Dest = Dest - Src
<code>imulq Src, Dest</code>	Dest = Dest * Src
<code>xorq Src, Dest</code>	Dest = Dest ^ Src
<code>orq Src, Dest</code>	Dest = Dest Src
<code>andq Src, Dest</code>	Dest = Dest & Src
<code>negq Dest</code>	Dest = - Dest
<code>notq Dest</code>	Dest = ~ Dest
<code>salq k, Dest</code>	Dest = Dest $\ll k$
<code>sarq k, Dest</code>	Dest = Dest $\gg k$ (arithmetic)
<code>shrq k, Dest</code>	Dest = Dest $\gg k$ (logical)

Addressing modes

- Immediate
`$val Val`
`val:` constant integer value
`movq $7, %rax`
- Normal
`(R) Mem[Reg[R]]`
`R:` register R specifies memory address
`movq (%rcx), %rax`
- Displacement
`D(R) Mem[Reg[R]+D]`
`D:` constant displacement 1, 2, or 4 bytes
`R:` register specifies start of memory region
`D:` constant displacement D specifies offset
`movq 8(%rdi), %rdx`
- Indexed
`D(R,Ri,S) Mem[Reg[Rb]+S*Reg[Ri]+D]`
`D:` constant displacement 1, 2, or 4 bytes
`Rb:` base register: any of 8 integer registers
`Ri:` index register: any, except %esp
`S:` scale: 1, 2, 4, or 8
`movq 0x100(%rcx,%rcx,4), %rdx`

Choosing instructions for conditionals

compare or test	cmp b,a	test b,a	Examples:
<code>je "Equal"</code>	<code>a == b</code>	<code>a b</code>	<code>cmp 5, (%rax)</code>
<code>jne "Not equal"</code>	<code>a != b</code>	<code>a b</code>	<code>je: (%rax) == 5</code>
<code>js "Sign"</code> (negative)	<code>a < b</code>	<code>a b</code>	<code>jne: (%rax) != 5</code>
<code>jns "non-negative"</code>	<code>a b</code>	<code>a b</code>	<code>js: (%rax) < 0</code>
<code>jg "Greater"</code>	<code>a > b</code>	<code>a b</code>	<code>jns: (%rax) >= 0</code>
<code>jge "Greater or equal"</code>	<code>a >= b</code>	<code>a b</code>	<code>jg: (%rax) > 0</code>
<code>jl "Less"</code>	<code>a < b</code>	<code>a b</code>	<code>jge: (%rax) >= 0</code>
<code>jle "Less or equal"</code>	<code>a <= b</code>	<code>a b</code>	<code>jl: (%rax) < 0</code>
<code>ja "Above" (unsigned >)</code>	<code>a > b</code>	<code>a b</code>	<code>jle: (%rax) < 0</code>
<code>jb "Below" (unsigned <)</code>	<code>a < b</code>	<code>a b</code>	<code>ja: (%rax) > 0</code>
			<code>test %rdi, %rdi</code>
			<code>je: %rdi == 0</code>
			<code>jne: %rdi != 0</code>
			<code>jg: %rdi > 0</code>
			<code>jl: %rdi < 0</code>
			<code>test %rax, 0x1</code>
			<code>je: %rax == 0</code>
			<code>jne: %rax != 0</code>

What to do on a write-miss?

Write-allocate (load into cache, update line in cache)

- Good if more writes to the location follow
- More complex to implement
- May evict an existing value
- Common with write-back caches

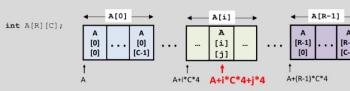
No-write-allocate (writes immediately to memory)

- Simpler to implement
- Slower code (bad if value consistently re-read)
- Seen with write-through caches

Nested Array Element Access

Array Elements:

- `A[i][j]` is element of type `T`, which requires `K` bytes
- Address = $A + i \cdot (C \cdot K) + j \cdot K$
 $= A + (i \cdot C + j) \cdot K$



Condition codes

CF Carry Flag
ZF Zero Flag
SF Sign Flag
OF Overflow Flag

Integer registers

%rax	Return value
%rbx	Callee saved
%rcx	4th argument
%rdx	3rd argument
%rsi	2nd argument
%rdi	1st argument
%rbp	Callee saved
%rsp	Stack pointer
%r8	5th argument
%r9	6th argument
%r10	Scratch register
%r11	Scratch register
%r12	Callee saved
%r13	Callee saved
%r14	Callee saved
%r15	Callee saved

Nested Array Element Access Code

```
int get_pgh_digit(int index, int digit){  
    return pgh[index][digit];  
}
```

```
leaq (%rdi,%rdi,4), %rax      # 5 * index  
addq %rax,%rsi                # digit + 5 * index  
movl pgh,(%rsi,4),%eax        # *(pgh + 4*(digit + 5*index))
```

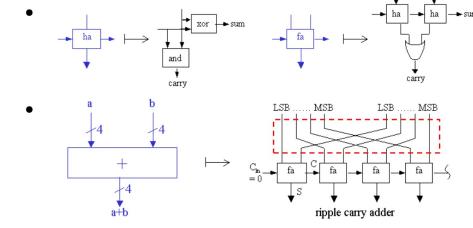
Array Elements:

- `pgh[index][digit]` is int and `sizeof(int)=4`
- Address: $pgh + 5 \cdot 4 \cdot index + 4 \cdot digit$

Assembly Code:

- Computes address as: $pgh + ((index+4*index) + digit) \cdot 4$
- `movl` performs memory reference

Add / subtract



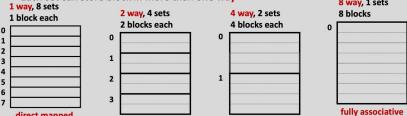
subtractor: a-b

- $-(\sum 2^i x_b) + (\sum 2^i x_b) = -1$ e.g. $0101_2 \oplus 1010_2 = 1111_2 = -1_{10}$
- $-(\sum 2^i x_a) - (\sum 2^i x_b) = (\sum 2^i x_a) + (\sum 2^i x_b) + 1$
i.e. invert b bitwise and set $C_{in} = 1$ for adder

wl 2023.4.7

So we combine two ideas:

- Each address maps to exactly one set
- Each set can store block in more than one way



Microprogram: field assignment

- Fetch: ALUControl = Add, SRC1 = PC, SRC2 = 4, Memory = ReadPC, PCWriteControl = ALU, Sequencing = Seq
- 1: ALUControl = Add, SRC1 = PC, SRC2 = ExtShift, ReadControl = Read, Sequencing = If (opcode=lv) or (opcode=sr) then goto Mem1 else if (opcode=R-type) then goto Rformat1 else if (opcode=B/EQ) then goto BEQ1 else if (opcode=D) then goto JUMP1
- 2: Mem1: ALUControl = Add, SRC1 = A, SRC2 = Extend, Sequencing = If (opcode=lw) then goto LW2 else goto SW2
- 3: LW2: Memory = ReadALU, Sequencing = Seq
- 4: RegControl = WriteMDR, Sequencing = Fetch
- 5: SW2: Memory = WriteALU, Sequencing = Fetch
- 6: Rformat1: ...



Microsequencer details

- Dispatch ROM 1 (DR1): R-type 000000 0110 state 6
jump 000010 1001 state 9
beq 000100 1000 state 8
lw 000101 0010 state 2
sw 101011 0010 state 2
- if R-type then state 6
else jump then state 9

C declaration	Intel data type	Assembly code suffix	Size (bytes)
char	Byte	b	1
short	Word	w	2
int	Double word	l	4
long int	Quad word	q	8
char*	Quad word	q	8
float	Single precision	s	4
double	Double precision	t	8
long double	Extended precision	tt	10/12

What to do on a write hit?

- Multiple copies of data exist. What is the problem with that?
- Write-through
 - Write immediately to memory and all caches in between
 - Memory is always consistent with the cache copy
 - Slow: what if the same value (or line!) is written several times
- Write-back
 - Defers write to memory until line is evicted (replaced)
 - Need a dirty bit
 - Indicates line is different from memory
 - Higher performance (but more complex)
- Capacity miss:
 - Occurs when the set of active cache blocks (the working set) is larger than the cache
 - Note: Fully-associative only has Compulsory and Capacity misses

Satisfying Alignment with Structures

- Within structure:
 - Must satisfy each element's alignment requirement
- Overall structure placement:
 - Each structure has alignment requirement K_{max}
 - K_{max} = largest alignment of any element
 - Counts array elements individually as elements
 - Address of structure and structure length must be multiples of K_{max}
- Example: Internal fragmentation
 - Internal fragmentation
 - Multiple of 8
 - Multiple of 4
 - Multiple of 2
 - Multiple of 1

K_{max} is 8, due to double element

P10 P4 P8 P16 P24

Multiple of 8 Multiple of 4 Multiple of 2 Multiple of 1

Multiple of 8