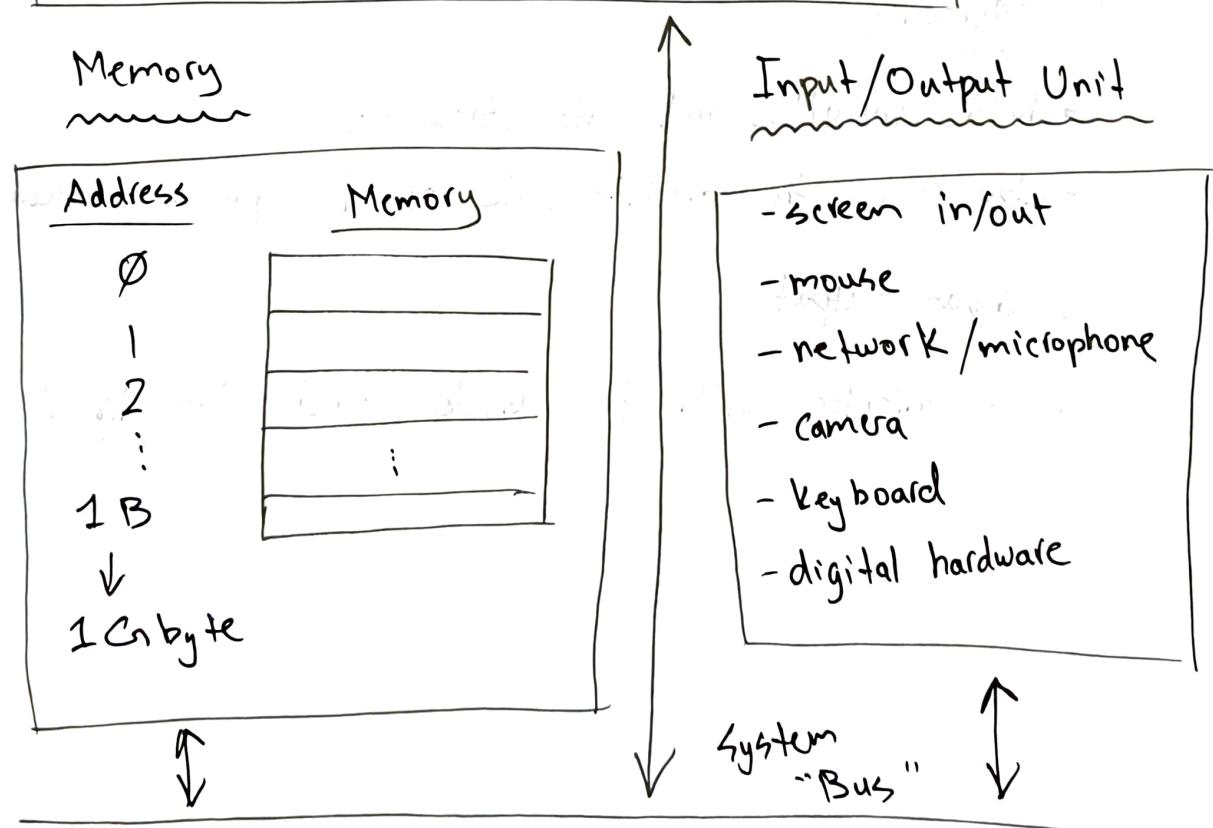
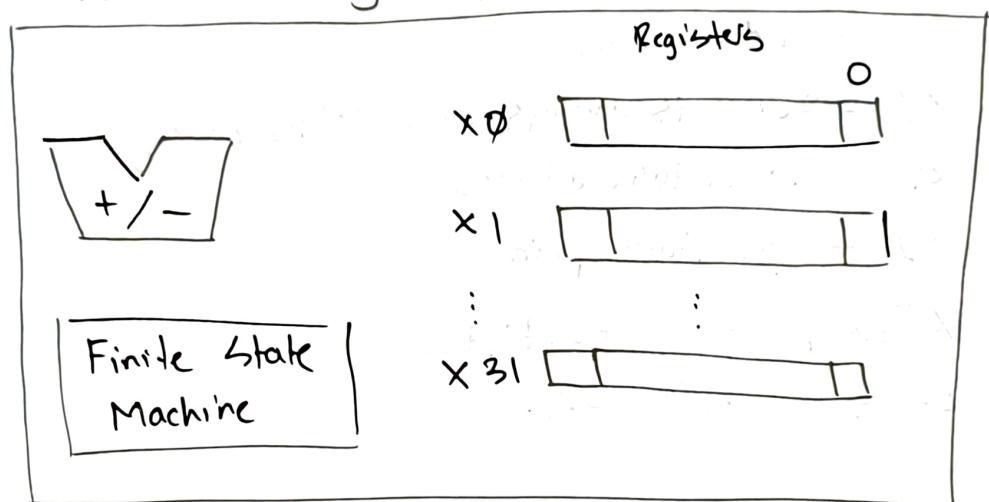


Today: Basic Computer Organization and Assembly Language Programming

Computer Organization

- * A computer has 3 main parts

Central Processing Unit (NIOS V/RISE V)



① CPU:

- controls performs computations
- makes decisions and controls flow of program
- using NIOS V in this course
- main computation is done in the 32×32 bit registers
- each register has a second, more useful name
 - ↑ $x_0 \rightarrow x_{31}$

e.g. x_0 — zero : is always zero (0 is handy)

x_1 — ra : return address

x_2 — sp : stack pointer

x_5, x_6, x_7 — t $_0, t_1, t_2$

Others — s $_0 \rightarrow s_8$

② Memory:

- holds variables and data structures
- also holds the machine code of the program executing

③ Input / Output

- communicates with the outside world thru digital logic

Programming (C program magic)

int A, B, C, D; ①

+ 4 × 32-bit numbers

+ Let's put A, B, C, D into registers t0, t1, t2, t3

t0	1	A
t1	8	B
t2	9	C
t3	7	D

A = 1; ②
B = 8; ③
C = A + B; ④
D = B - A; ⑤

Assembly

- ① no assembly language
- ② li t0, 1
- ③ li t1, 8
- ④ add t2, t0, t1
- ⑤ sub t3, ~~t0~~, t1
t1, t0

Function

(nothing to do)

t0 ← 1

t1 ← 8

t2 ← t0 + t1

t3 ← ~~t1 - t0~~ ? might have it backwards?
t1 - t0

loop: j loop

• li means
"load immediate"

• label called loop
says jump to
loop ... an infin.
~~end~~ loop

Today: Labels, Loops and Conditional Branches
in NIOS V Assembly Language.

Lab 1 Part 2: learn CPU later, single stepping, breakpoints

Part 3: write a program

Part 4: learn about the monitor program

Part 5: run part 3 program

Conditional Branch

- * to make a loop, you need a "conditional branch" assembly statement
- * have an unconditional "branch" = "jump"

0x20
↑
address

loop: j loop
label

- * a conditional branch → transfers execution to an address/label if a specific condition is true (e.g. $t2 = t3$) otherwise → execution goes to the next instruction in sequence
- * an example: a loop that executes 4 times

(and does nothing)

.global _start ← label called _start
_start: li t0, 1 | * t0 ← 1
 li t1, 4 | * t1 ← 4
 represents address where -- goes

t0 is one of 32 registers

myloop: addi t0, t0, 1 * $t0 \leftarrow t0 + 1$
 ble t0, t1, myloop * "branch" less than or
 equal to ... is $t0 \leq t1$?

fin: j fin

- general form of conditional branch statement:

bXX rA, rB, DEST-LABEL

- rA and rB are the two registers being compared
- XX is the condition

~~eg.~~ if rA \square rB then branch to DEST-LABEL,
else go to next statement

ne	" != "
ge	" \geq "
le	" \leq "
lt	" < "
gt	" > "

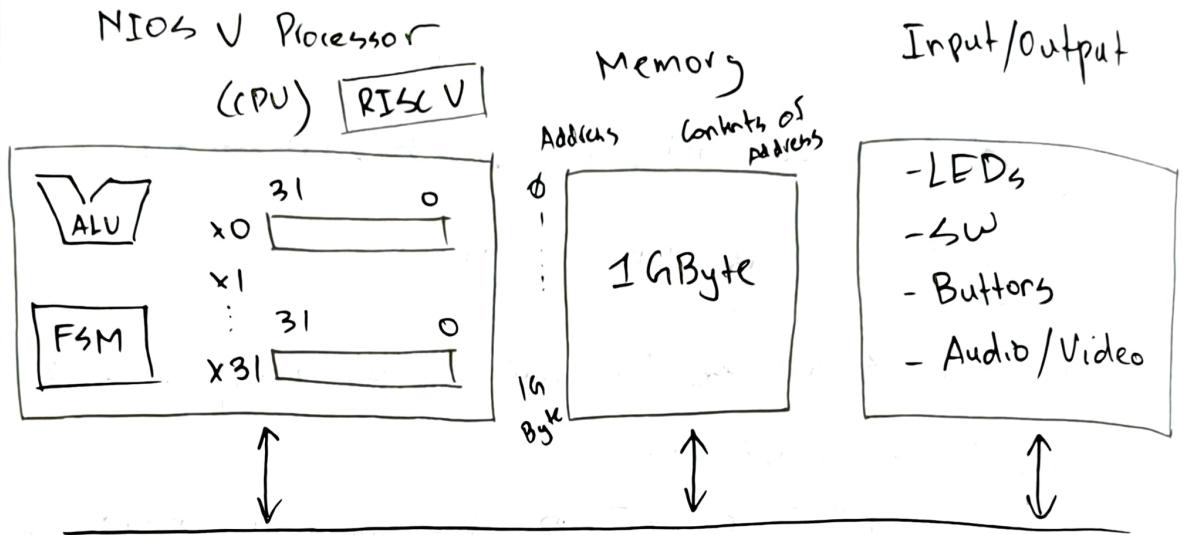
} you add "u" at the end
(e.g. leu, ltu, gtu...) for
unsigned

- program gets turned into 5 words of numbers in memory
- "_start" has the value that the assembler knows about
 ↗ address where the instruction "li t0, 1" is stored
 ↗ equals "00000000"
- "myloop" becomes "00000008", which is address that stores
the "addi t0, t0, 1" instruction.

Today: Memory organization, in NIOS V processor

Recall Computer Organization

32-bit



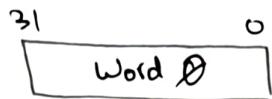
- max # of addresses $\leq -2^{32}$
- NIOS V is a 32-bit processor
 - registers are 32 bits wide
 - number of addressable locations in memory = 2^{32}
 - max amount of bits retrieved/sent to memory is 32
 - "word size" of processor is 32 bits
 - however, the memory is byte (8 bits) addressable
 - every byte gets a separate address
 - even though the processor can access 32 bits at a time = 4 bytes at 4 different addresses

Memory

Address

0x0

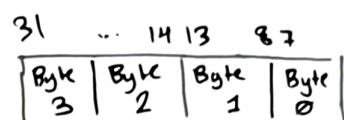
Memory
Contents



this 32-bit word
contains 4 bytes

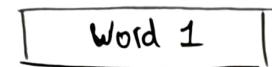
MSB

LSB



a "x" or
actually
"0x"
is zero
and it's
in hex

0x4



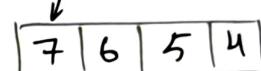
most sig.
byte

least
sig. byte

0x8

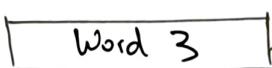


byte at
address 7



".word"
is a
command
to assemble
to put
"10" into
mem addrs.
0x2c

0xc



→



→



→



→



Program to add 3 nums
(10, 20, 30) in memory

→ need to initialize memory like this in C

```
int list[3] = {10, 20, 30};
```

→ Skeleton of the program to start: and how it
is laid out in memory

Address

Assembly Code

0x0

la t0, list

takes address of list (0x2c)
and puts it into t0

0x8

lw t1,(t0)

pointer dereferencing ... load ~~word~~
into t1 what's stored at the
mem address fitting in t0

0x28

done: j done

→ 0xa

0x2c

list: .word 10

→ 0x14

0x30

.word 20

→ 0x1e

0x34

.word 30

→ ~~0x1a~~

base 16

0x1e

Today: Add 3 numbers in full

1. Program to add 3 numbers in memory

- ①
 - global — start
 - start la t \emptyset , list
 - "la" means load an address (32 bits) into a register
 - in this case, list = 0x2c
 - now, $t\emptyset \leftarrow 0x2c$

done: j done
 list: .word 10
 .word 20
 .word 30
 .word 40
 answer: .word 0

⑥

0x2c	10
0x30	20
0x34	30
0x38	0

②

lw t1, (t \emptyset)

- * "lw" means "load the word" from the memory address given in register t \emptyset
- * $t1 \leftarrow (t\emptyset)$ — just like dereferencing a pointer
- * you should read the brackets () as "contents of memory address given by what's inside ()"

$$\therefore t1 \leftarrow (0x2c) = 10$$

③

lw t2, 4(t \emptyset)

* similar, except we add 4 to t \emptyset and then load from memory that address

④

lw t3, 8(t \emptyset)

$$\therefore t2 \leftarrow (t\emptyset + 4) = (0x2c + 4) = (0x30) = 20$$

$$\therefore t3 \leftarrow (t\emptyset + 8) = (0x34) = 30$$

⑤

add t2, t2, t1

$t2 \leftarrow t2 + t1$

add t2, t2, t3

$t2 \leftarrow t2 + t3$

la t0, answer

* puts answer = 0x38 \rightarrow t0

sw t2, (t0)

* "sw" means "store word" in
register t2 into memory address

+ last 2 lines can be given by la, t0

be replaced by

"sw t2, 12(t0)"

- Today:
- Accessing smaller items in memory: bytes/half words
 - a more complex program from lab 2

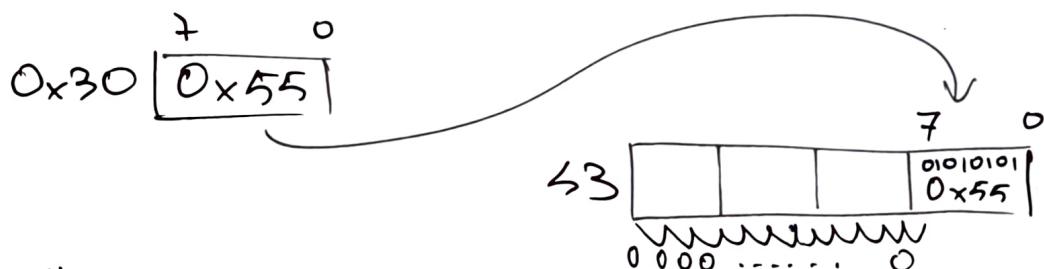


- an instruction like `lw t2, (t3)`
loading words so far = 32 bits / 4 bytes

→ NIOS V has byte-addressable memory

• we can also access a single byte from memory using the load byte (lb) instruction

eg: `lb $3, ($2)` # load the byte @ address given in reg \$2 into reg \$3
0x30



- lb use sign extension
- the most significant bit, bit 7, is copied into bits 8 → 31
- if don't want sign extension → use lbu
↳ unsigned

- + similarly, `lh` loads half a word = 2 bytes from memory and sign extends from bit 15 to bits 16 → 31
- `lh` is unsigned

Initializing Memory With Values

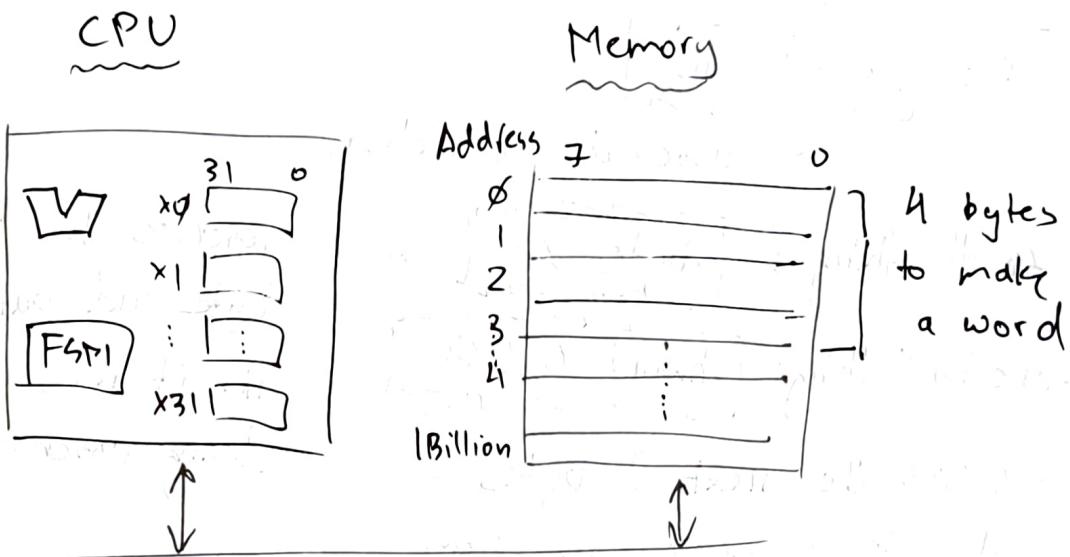
`list .word 10`
 assembler directive

- + small thing: `.byte 50` ← reserves the "next" byte and put 50 in it in memory during assembly
- + one half thing: `.hword 0x6600`
 reserves the next 2 bytes and puts 0x6600 into them
- + bytes can be @ any address
- + half-words must have an address on 2 byte boundaries ← address must be even
- + words must have addresses that are divisible by 4

Today:

- + more types and details of NIOS V instructions
- + move, add, load
- + logic and shift

Context of all this



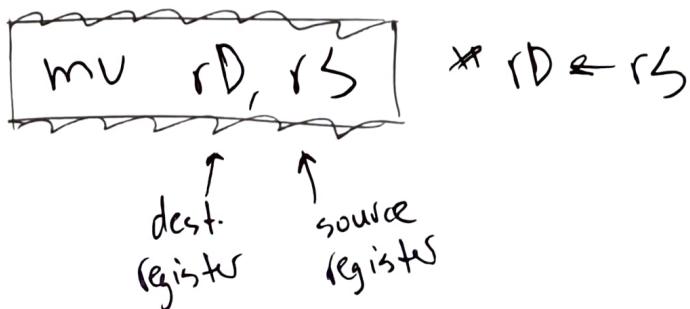
General Forms of Instructions in NIOS V

→ eg: move instruction

eg: $mv\ t_3, t_2 \quad *t_3 \leftarrow t_2$

copies values in t_2 into t_3

general form:



- * the "la" and "li" instructions are really more like move immediate instructions
 - eg: la a \emptyset , $\emptyset \times 2\emptyset$ $\Rightarrow a\emptyset \leftarrow \emptyset \times 2\emptyset$ (full 32 bits of hex 0x20)
 - immediate constant
 - encoded inside the instruction
 - + may be implemented with an "li"
- the add and subtract instructions
 - $\boxed{\text{add } rC, rA, rB} \Rightarrow rC \leftarrow rA + rB}$ registers only
 - $\boxed{\text{addi } rC, rA, imm} \Rightarrow rC \leftarrow rA + imm$
 - $\boxed{\text{sub } rA, rB, rC} \Rightarrow rC \leftarrow rA - rB$
 - a 12-bit number that is sign extended
- * sub does not have "subi"
- + use addi but use a negative number

→ Recall: load word

$lw t4, (a\emptyset)$ * load into register $t4$
the word beginning at
memory address given
in $a\emptyset$: $t4 \leftarrow (a\emptyset)$

General form

$lw rB, D(rA)$ |
 $\underbrace{r}_{\curvearrowright} B \leftarrow (rA + D)$

* load into rB the word beginning at $rA + D$

→ Similarly, the store word

$sw rB, D(rA)$ | * copies all of rB into
the word in memory beginning
at address $rA + D$

→ recall lb, lh, sb, sh

Logic and Shift Instructions

* need to isolate and access individual bits

* use bit-wise AND, OR, XOR → in C: $\$\$$ ||

* these have two general forms

and
or
xor } } $XXX rC, rA, rB$ | * bit-wise operation
of rA w/n rB result
into rC

$XXXi rC, rA, imm$ |

* immediate form

eg: and t3, t1, t2 $t3 \leftarrow t1$ BIT wise AND

t1: ... 0 0 1 0 1

t2: ... 0 1 1 0 1

t3: ... 0 0 1 0 1

BIT wise AND

eg: ori t2, t1, 0x007

t1: ... 0 1 0 1 0 1

imm~~t2~~: ... 0 0 0 1 1 1

t2: ... 0 1 0 1 1 1

sign extended
(max 12 bits)

→ Shift Instruction

also need to move bits within a register
(shifting them left or right)

eg: srl i t2, t1, 3

* shift t1 3 bits
to the "right", bring
in 0's into higher
order bits

T1~~reg~~: 1 1 0 1 ... 1 1 0

T2 : 0 0 0 1 1 0 ...

"shift right
logical immediate"

can replace
this with a
register, and

* shifting
right by 1
divide
by 2

"srl" will shift
right by x and
x is stored in
that register

* shifting
left by 1
multiply
by 2

- ∴ srl i ← immediate (right)
- srl ← not // (right)
- sll i ← (left) immediate
- sll ← (left) not //

Today: Subroutines, linkage and parameter passing

Subroutines

- + a subroutine is separate code that is re-used in different parts of a program
- + API (application programmer interface)
- Key to well-structured code
- makes your code/software useable by someone else and vice-versa.

Ex: Sample C program w/h a function (subroutine)

```
int main(void){  
    int x, z;  
    x = my_sub(3);  
    z = my_sub(4);  
}  
  
int my_sub(int p)  
    return (p+p);
```

Q1: How does my-sub "know" where in the code* to go back to? *that called it

Q2: How do parameters (arguments) like 3 and 4 get sent to my-sub and the answer returned?

Q1:

one of
the "argument"
registers

+ recall that each assembly instruction (eg: li a0, 3) is encoded as a number (eg: 00x00300513)

by the assembler:

+ the PC points to the next instruction to be executed

program
counter

= 0x0 x00300513

not one of
the $x0 \rightarrow x31$

→ 0x4 x05621ABC

↳

+ there is another register, called the program counter (PC) and it contains the address of the next instruction to be executed

+ we need to save the program counter when calling a subroutine

• will save it in register x1, also known "ra"
which means "return address"

Ex: Assembly Version

main: li a0, 3 *a0 ← 3

call my-sub *ra ← pc, pc ← my-sub

next: li a0, 4 *a0 ← 4

call my-sub *ra ← pc, pc ← my-sub

done: j done

my-sub: add a0, a0, a0 *a0 ← a0 + a0

ret

*pc ← ra

we are
not showing
the saving
of answers
into x, z

Today: How the stack works and using it for subroutine "linkage"

Recall:

- the program counter always knows address of next instruction to be executed
 - the "call" copies pc into ra : $ra \leftarrow pc$
 - the "return" copies ra into pc : $pc \leftarrow ra$
- we need a way to save and restore ra appropriately as a subroutine calls another subroutine calls another...
- a stack data structure → a ~~set~~ chunk of memory
- can do 2 operations: "push" a new item onto the stack and ② "pop" an item off of the stack into a register
 - in H10sV, register X2, also called "SP" (top of ~~point~~ stack pointer) → stack pointer for short.

- stacks grow downwards in memory addresses

ECE468

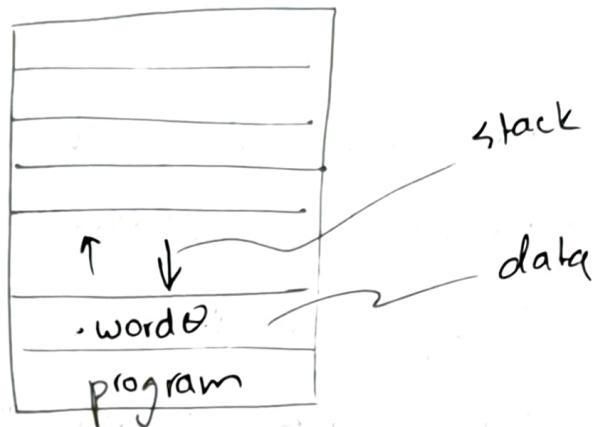
20000

1F FFF

1F FDx

:

0x0



- we will initialize our stacks in the course/lab to begin @ address 0x20000 → grows downwards when pushing items more

- let's push 0x1234f678 onto stack

first initialize stack pointer

la \$p, 0x20000

* \$p ← 0x20000

li t0, 0x1234f678

* want to push onto
the stack

Address Data



- we want to push a 4 byte (word) onto stack

Step 1. subtract 4 from sp

$\boxed{\text{addi } \text{sp}, \text{sp}, -4} \quad * \text{sp} \leftarrow \text{sp} - 4$

Step 2. store the item at the address

$\boxed{\text{sw } t0, (\text{sp})} \quad * \text{memory @ sp address} \leftarrow t0$

* that's a push

* to pop an item off of the stack

$\boxed{\text{lw } t1, (\text{sp})} \quad * t1 \leftarrow (\text{sp})$

$\boxed{\text{addi } \text{sp}, \text{sp}, 4} \quad * \text{sp} \leftarrow \text{sp} + 4$

* now removed from stack, (but is in $t1$)

Let's use the stack to correctly allow subroutines to call subroutines and save and restore correctly with the "last in first out" behaviour

la \$p, 0x20000

li \$a0, 1

call sub-1 * \$ra ← PC, PC ← sub-1

done: j done

*1 sub-1: addi \$p,\$p,-4 } push ra into
 sw \$ra,(\$p) } stack

*2 → call sub-2
 :
 lw \$ra,(\$p) } pop ra off
 addi \$p,\$p,4 } stack
 ret

sub-2: addi \$p,\$p,-4 } push ra
 sw \$ra,(\$p)

*3 → call sub-3
 :
 lw \$ra,(\$p) } pop ra
 addi \$p,\$p,4
 ret

- Today:
- Introduction to Memory-mapped Input/Output
 - Connection between the virtual (inside computer) and the outside, real, physical world
 - Nios V/Risc V calling convention for subroutines

Memory-Mapped I/O

From ECE241 and this course, we use DE1-SoC Board

• it has an FPGA on it → ?

• it also has

10 LEDs (LED_R 9:0)

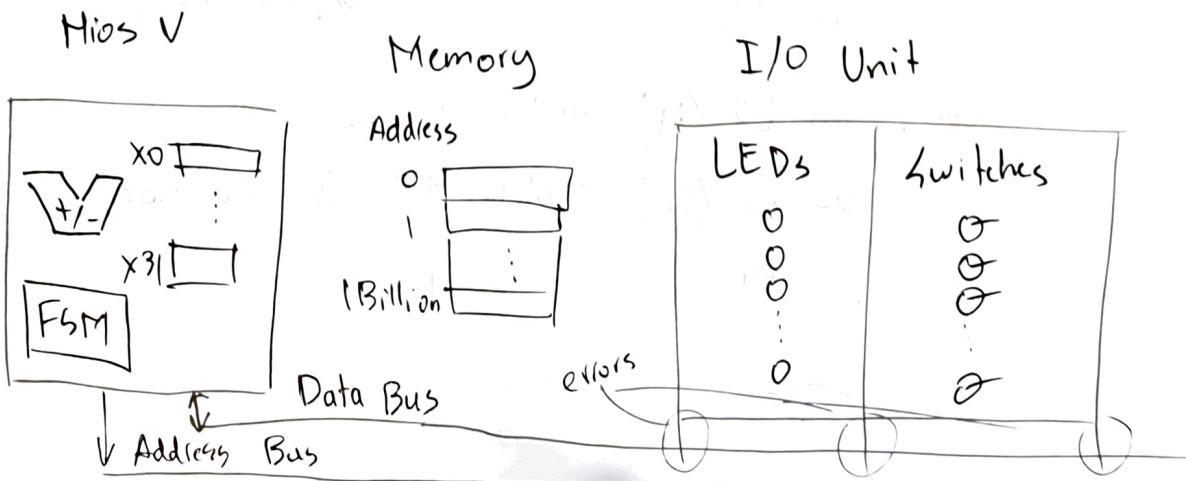
10 Slider switches (SW 9:0)

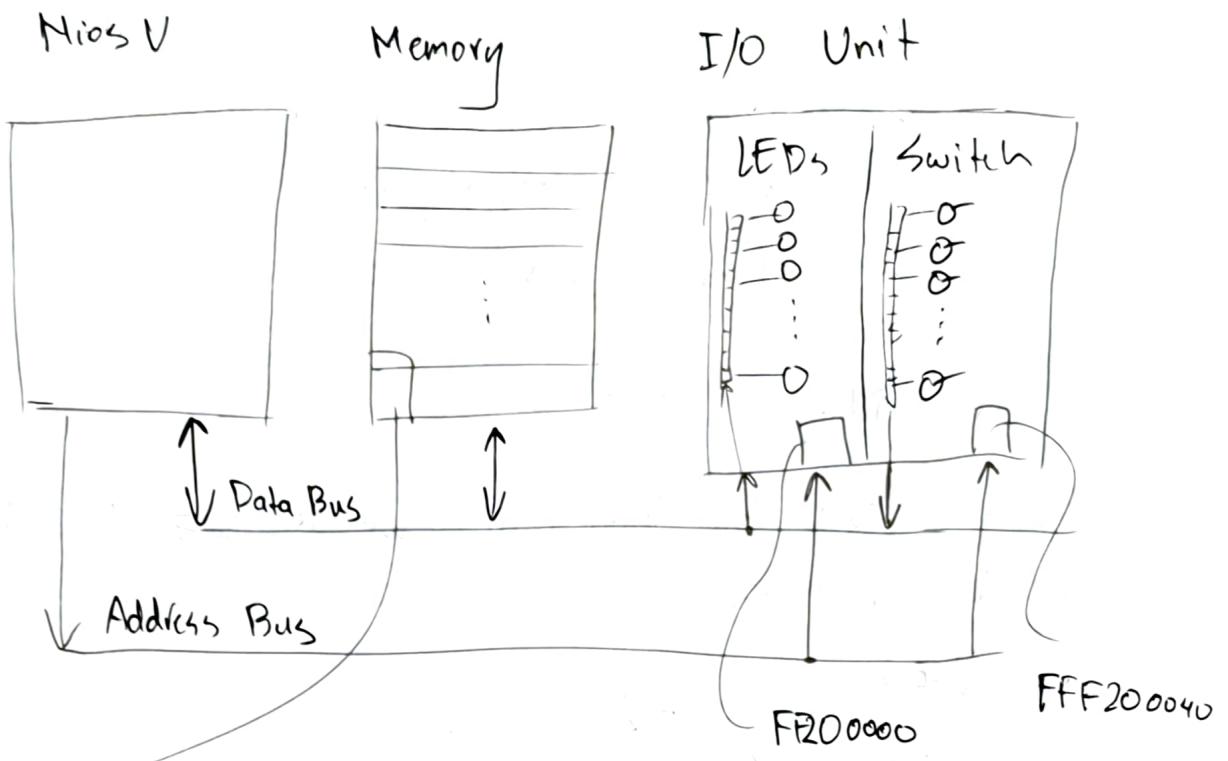
4 Buttons (Key 3:0)

6 Seven-segment Displays (SEG7 5:0)

• these are also available to your program running on the Nios V

• Recall the computer's organization





Digital CKT to decide if address is in the range $0 \rightarrow 1B$

here is code to read the switches and write that number to LEDs:

```
.equ LEDs, 0xFF200000 } these, set the labels
.equ SW, 0xFF200040 } to be equal to those
                      addresses for program
                      readability
```

-start: la t0, LEDs

la t1, SW

*t0 \leftarrow 0xFF200000

loop: lw t2, (t1) ** loads the switches

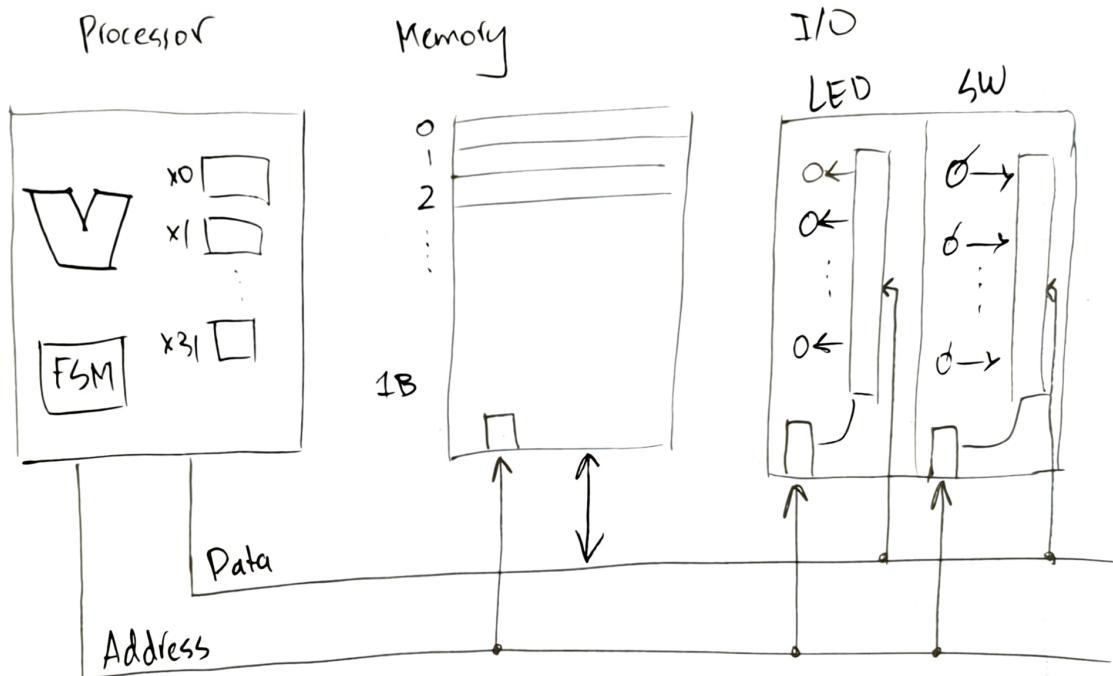
sw t2, (t0) ** stores that value into LEDs

j loop

Today: I/O: parallel ports and The Polling Method of Synchronization

Recall Memory Mapped I/O

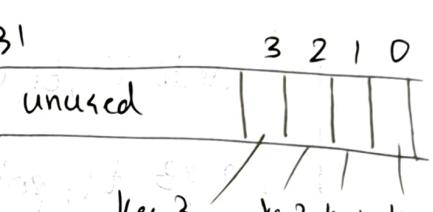
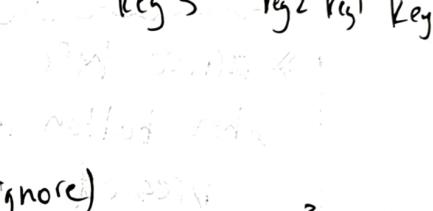
Lat 11
Check
if the
button
got
released



- + the I/O devices have their own registers to hold and synchronize the data transfers
- + these are distinct from the processor registers
- + there are typically more registers that have a variety of functions
- + these are collected into a parallel port.

- These weeks, we're covering synchronization b/n computer and I/O devices
 - ↳ Polling method (ask a lot)
 - ↳ Interrupt Method
- Let's illustrate polling using the pushbuttons (KEY0-3)
- they are connected to a memory-mapped parallel port

Memory Mapped Address Register

- 0xFF200050 Data Register → 
- + bit i goes to 1 when key i pressed
 - + goes back to 0 when key i released
- 0xFF200054 Director Register (ignore) 
- 0xFF200058 Interrupt Mask Register → 
- 0xFF20005C Edge Capture Register 
- + bit i is set to 1 when key i is released
 - + bit i is reset to 0 if a 1 is stored into it (sent to it)
 - + bit i stays the same if 0 is stored into it

- Let's write a program that to use the data register and poll it to see if key1 is pressed.

```
.equ KEY_BASE, 0xFF200050
```

```
la t0, KEY_BASE
```

* polling loop

poll: lw t1, (t0) * load data register ↑ to bank
to find if Bit #1 is 1

```
andi t1, t1, 0x2
```

```
beqz t1, poll
```

* arrive here
when button is
pressed

Data Register B3 B2 B1 B0
AND and 0 10 of t1, 0

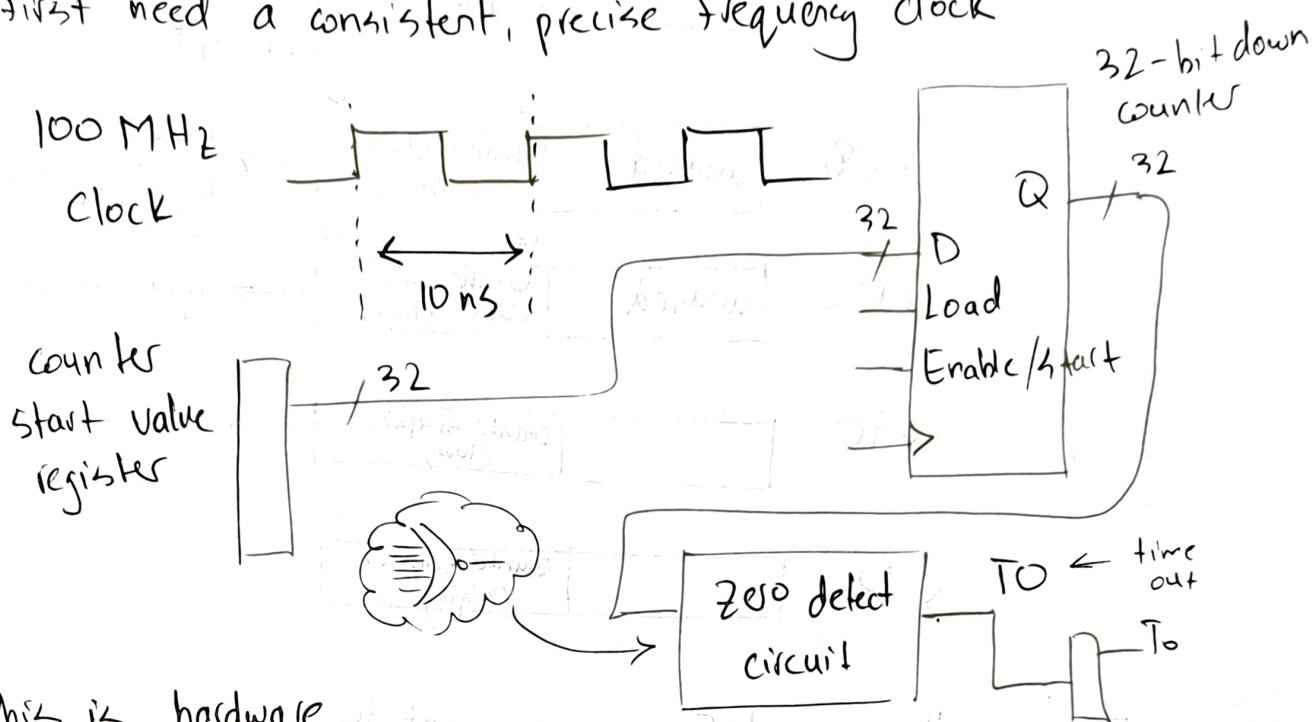
→ zeros out all bits

except bit 1
if word of t1 = 0,
button is not pressed

Today: The timer I/O Device
How software knows what time it is

Hardware Based Timer

- First need a consistent, precise frequency clock



- this is hardware
- the timer I/O device in our DE1-SoC is just like this
- but with memory mapped registers that connect to this ckt

Here is the Interval Timer

- has 6 memory mapped registers

Mem. Mapped Address

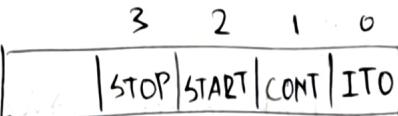
0xFF202000



load
(read only)

Status register

x04

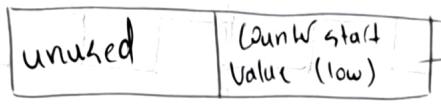


Control register

(write
or store)

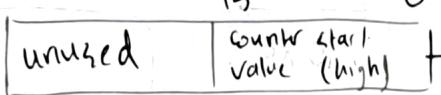
15 0

x08



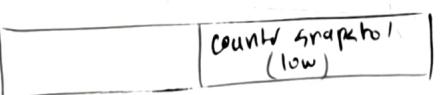
lower 16 bits
of the counter
start register

x0C

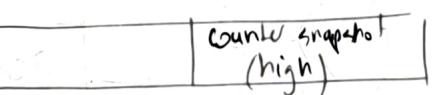


upper 16 bits

x10



x14



To use the counter

to count to

- ~~write~~
- 1) Put the number you want into the count start val register
→ to wait 1 second, 32-bit number needs to be 100 million
 - 2) To count down once, turn on start bit only in control register
to ~~to~~ count down continuously, turn on start bit and cont bit
 - 3) To determine if the most recent count down cycle has been reached, must "poll" the TO (time out) bit in status register.
if TO == 1, it means the count down to zero has been reached.

Program to turn LED0 on and off every second

```
.equ TIMER_BASE, 0xFF202000
.equ COUNTER_DELAY, 100,000000 * the number
.equ LED, 0xFF200000 of cycles of
                           100 MHz clock
                           in one second
_start
    la t5, TIMER_BASE
    sw zero, @t5 * clearing the TO bit in the
                   status register in case its on
    li t0, COUNTER_DELAY
    sw t0, 0x8(t5) * low 16 bits of starting count
    srl t1, t0, 16 * shift upper 16 bits of the
                   delay into the lower 16 bits of t1
    sw t1, 0xc(t5) * put that into the counter start
                      high register
* set up control bits → step 2, turn on cont and start
    li t0, 0b0110
           ^cont
           ^start
    sw t0, 4(t5) control register * store 0110 in to the
                                   timer control register
                                   * causes the to start
                                   and go into continuous
                                   mode
```

➤

Write
holes

la t6, LED } * used to toggle LED on/off

li t2, 1

tloop: sw t2(t6) * store the current value of light
into LED0 → "set LED0"

xori t2, t2, 1 * invert bit t2 → turns off LED0
~~nextime~~ next time

* now wait for timer to count down the next time in
* a polling loop

ploop: lw t0, (t5) * read the time status register AKA
* load timer status register

andi t0, t0, 1 * isolate bit 0 → "TO bit"

beqz t0, ploop * if TO bit is 0, keep asking/checking it

sw zero, (t5) * reset to zero (clear the TO bit)

j tloop

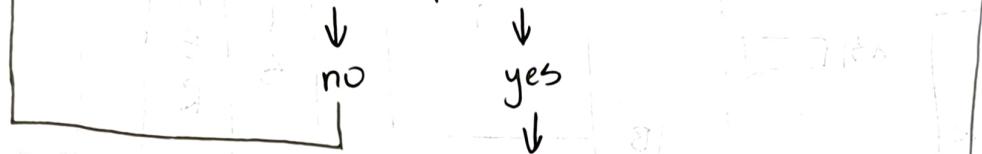
Week 4: Lecture 3

Jan 29, 2024

Today: Interrupt-driven I/O - Part 1 intro

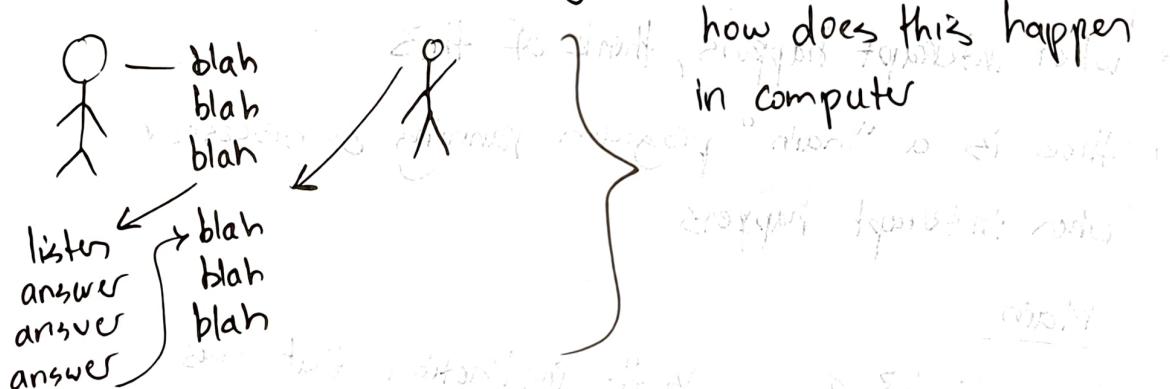
* Polling, though, is a mindless loop:

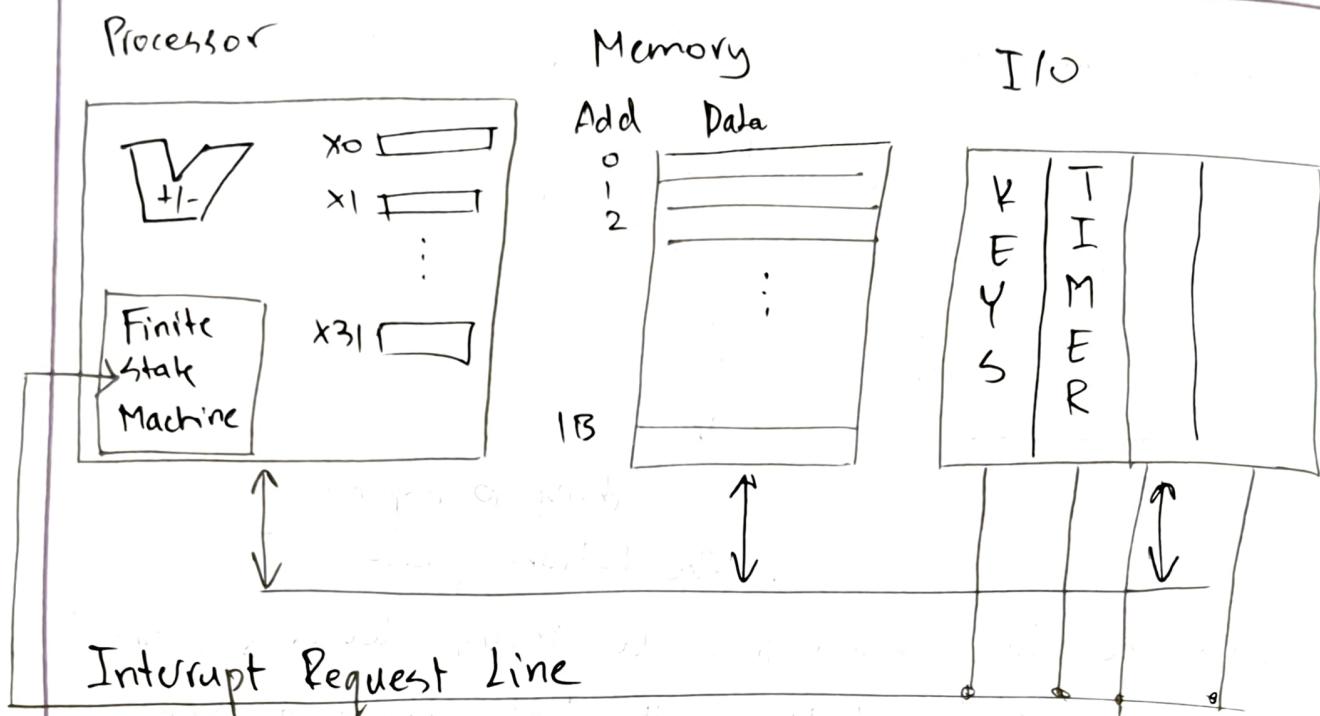
→ was the button pressed?



do something to respond
to the button press

- * can't use this when there are many devices connected to the computer + dealing with bad address accesses
- * need a better mechanism
- * interrupt driven I/O
- * like in a class, me teaching





- with interrupts, when a device needs attention (i.e. the timer timed out : TO, or button press) — it activates the interrupt request line
- when interrupt happens, think of this
- there is a "main" program running on processor when interrupt happens

Main

loop: li t3, 1

add ...

... : *

... :

j loop

is the instruction that was executing when, say, a button that is interrupting, is pressed

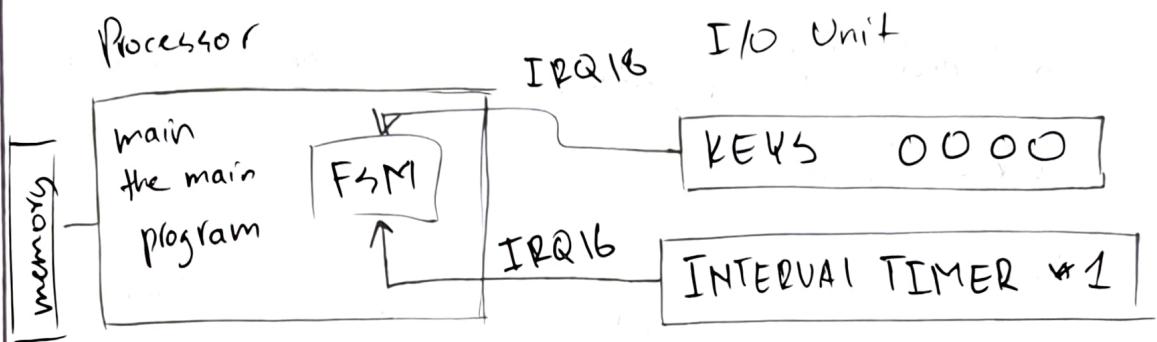
Interrupt Request Handler



- + when an interrupt happens, the code being executed stops
 - by the FSM in the processor
- execution is transferred to the interrupt request handler, which is code that you write.
- that code is executed
- execution returns to where it was stopped in the interrupted main program AS IF NOTHING HAPPENED.

Today: Interrupts part 2:

- need to synchronize external devices/signals with the software running on the processor



IRQ = interrupt request

- when someone pushes and releases a key that you want (and have set up) to cause an interrupt, let's say it happens at instruction $\textcircled{*}$, then
 - processor stops (doesn't execute $\textcircled{*}$)
 - processor starts executing interrupt request handler code
 - when finished the handler, execution returns to $\textcircled{*}$ with none of its registers "disturbed"

- * To make that happen and control it properly, there are two sets of registers —
 - ① in the processor
 - ② in the devices themselves (memory mapped registers)
- * Here are processor control and status registers:
 - (look at figure 5)

Address		Name
0x300	12 11 MPP MIE : 31 ... 16 7 3	mstatus
0x304	IRQ enable MIE1 - MIE1 31 ... 16 7 3	mie
	31 ... 16 7 3 + say handler address mode mtvec	
	exception program counter 31 ... 30	mePC
	interrupt exception code / IRQ	meCause

- * mstatus reg: MIE bit = 1 processors are enabled
= 0 " " disabled
 - global interrupt
- * mie reg: IRQ enable bits 16 → 31
 - each bit individually enables/disables a specific IRQ line

- when an interrupt occurs, the mcause register lower bits encode the number of the IRQ line that is interrupting
(eg: if it was the timer #1, the number 16 would be in mcause)
 - here is how you can modify or read these control and status registers in the processor
 - whole set of instructions on pg. 15-16 of Hobs - V - Intro
 - eg: to set the interrupt handler address into mtvec


```
la t0, interrupt_handler    * t0 ← ...
      csrw mtvec, t0    * mtvec ← t0
      :
      interrupt_handler : add sp, sp, -12
```

"CSRw" control
status register
write instruction
 - general form of control status

<u>CSRw</u>	<u>CSR</u> , <u>R/S</u>
-------------	-------------------------

CSR
 control
status
registers ↑ which of
the above
registers regular
processor
register

can also use an immediate constant
 - there is also $\text{CSRR} \leftarrow \underline{\text{read}}$
 - also get individual bits

 " clear " "
 - now: write, read, set, clear

 rw sr

- + there are other devices (pg. 11 of DE1-SoC-manual)
manual from lab 1
 - interval time : 16
 - Pushbutton KEY : 18
 - mouse / keyboard : 22, 23
- finally, recall the registers for KEYS parallel port:
 - Turn on Interrupt Mask Register bit i to have edge capture register bit i request an interrupt after key i is pressed and released
 - to turn off the request, you must turn off the edge capture bit in the usual way.

Today: Interrupts Part 2

Example program to ~~turn~~ toggle LED3 when Key 1 is pushed and released with interrupts

- global _start
- equ LEDs, 0xFF200000
- equ KEY_BASE, 0xFF200050

_start: la sp, 0x20000 * init stack - use for saving and restoring registers in the interrupt handler
 CSRw mstatus, zero

* tell the KEYS (Key 1) to request an interrupt when pressed and released

* use the KEY parallel port interrupt mask register

la t1, KEY_BASE

li t0, 0b00010

sw t0, 8(t1) * sets bit 1 of interrupt mask reg

sw t0, 12(t1) * turn off any previous edge capture on key 1

* enable specific IRQ line within the processor's

li t0, 0x40000 * this turns bit 18 on

CSRRS mie, t0 * turns bit 18 of mie on → enabling specific line IRQ18

* CSRRS: set specific bits, don't change the others

2nd level
enables

* tell the processor where the handler is in memory

la t0, interrupt_handler

csrw mtvec, t0 * mtvec ← t0



* finally, enable processor interrupts by turning on the MIE (bit 3) of the mstatus register

li t0, 0b1000

csrs mstatus, t0

main-loop: li t0, 1

li t1, 2

j main-loop

* here is where execution "goes" upon an interrupt

interrupt-handler :

↳ before this, mstatus
MIE (bit 3) is set

addi sp, sp, -12

to 0 — preventing

sw t0, 0(sp)

further interrupts

sw t1, 4(sp)

push the

lw ra, 8(sp)

registers used in the

handler onto the stack

* check to see if the cause of the interrupt was what we expected

(next page)

read

li t0, 0x7FFFFFFF
CSRr H, because * because encodes the number of the
and t1, t1, t0 IRQ line that caused the interrupt
* zeroes out bit 31 → used for something

li t0, 18

bne t1, t0, end-interrupt * if not equal, something's
call KEY-ISR wrong, return from interrupt

* KEY interrupt service routine

end-interrupt : lw t0, 0(sp)

lw t1, 4(sp)

lw ra, 8(sp)

add sp, sp, 12

mvf * return from interrupt: pc ← mepc

MIE (bit 3) is

set back to 1

* check Prof notes handout

for full note (the one that
talks about subroutine KEY-ISR)

Today: Connecting the C language to Assembly Language

Consider the compilation of a C-language program

C program file.c

```
int main() {
    while ~
    ~
}
}
```

→ gcc file.c →

file.s

```
.global _start
_start:
    li ...
    ~
    ~
main: ~
```

assembly
language
statement

file.o

- encoded instructions
- space in memory for data (full contents of memory prior to execution of program)

assemble
using
(gnu as)
assemble

- if there are multiple source files → file1.c, file2.c, file3.c

↑
code

Using Memory Mapped I/O Devices in C

- recall the first memory-mapped program (I/O) from Lee 9
 - copy the switches into the LED

.equ LEDs, 0xFF200000
.equ SW, 0xFF200020

-start: la t0, LEDs
la t1, SW

» infinite loop to copy over, over

loop: lw t2, (t1) ①
sw t2, (t0) ②
j loop

```
int main(void){  
    volatile int *LED_ptr =
```

:

value = *SW_ptr;
 *LED_ptr = value;

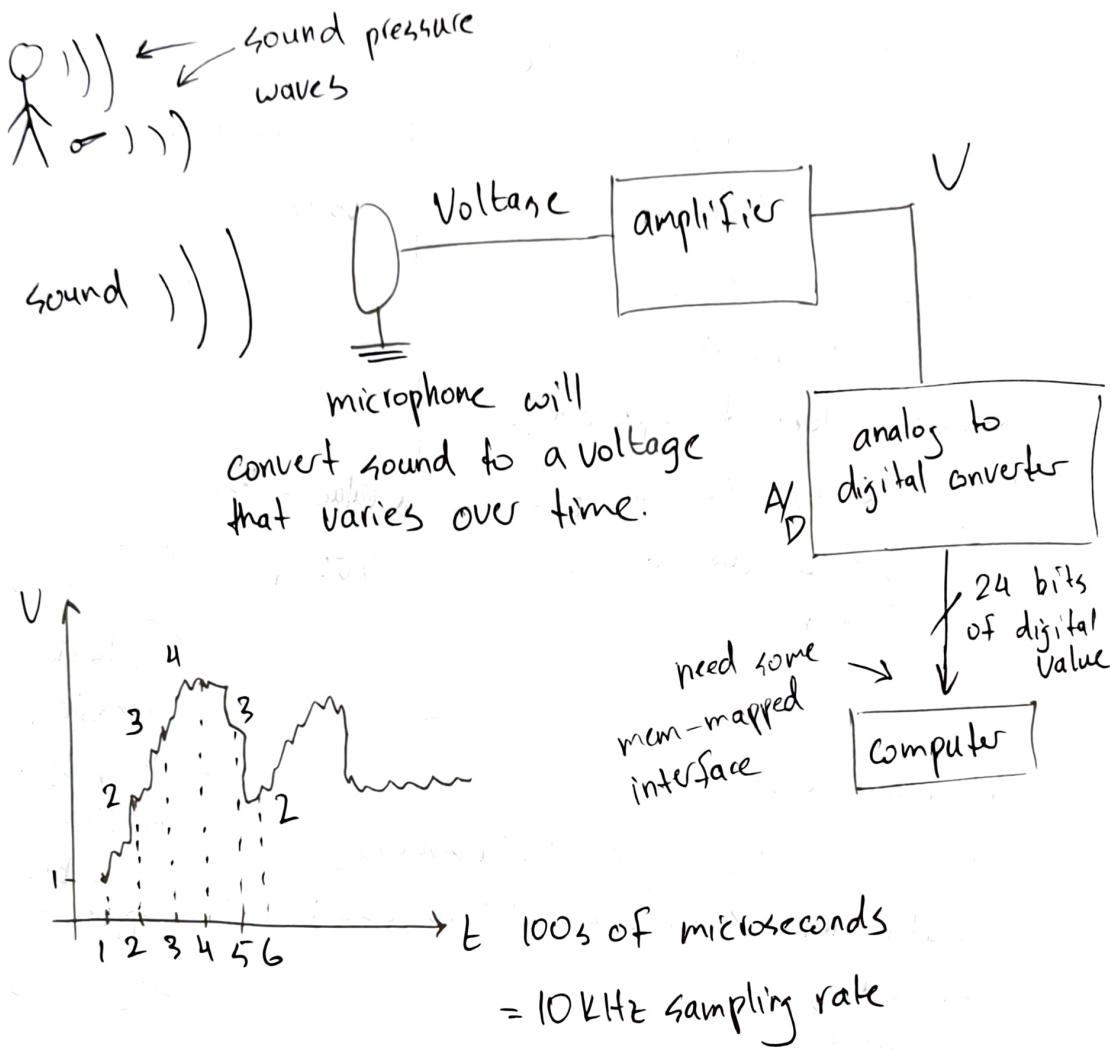
• Volatile means, to the compiler, "don't do fancy optimizations on this variable that would break the memory mapped I/O → X"

```
int main(void){  
    volatile int *LED_ptr = 0xFF200000;  
    volatile int *SW_ptr = 0xFF200020;  
    int value;  
    while(1){  
        value = *SW_ptr; //derefencing a pointer, does what lw statement ①  
        *LED_ptr = value; // same as the sw statement ②  
    }  
}
```

• an optimization that puts it in a register
• TBD: put it in cache

Today: Audio (sound) I/O - building on this course + others

Music / Audio

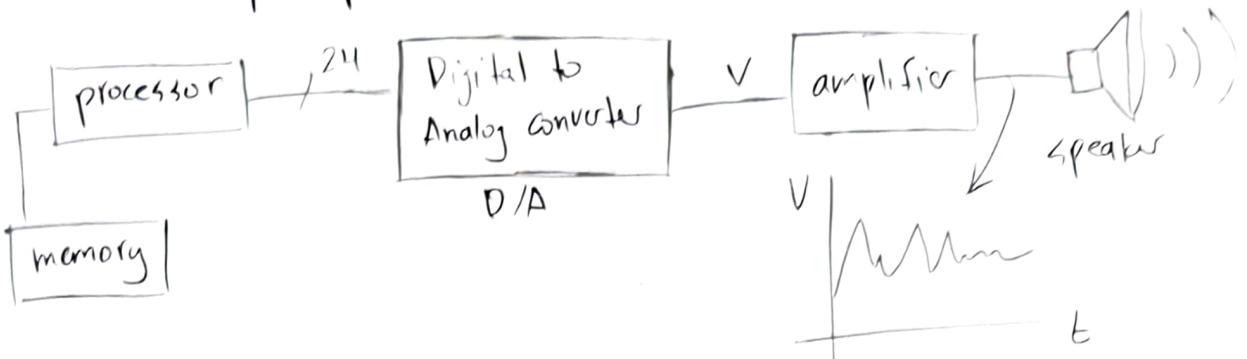


- the ADC converts the time varying voltage into a sequence of digital (binary) values
→ sampled at some rate: eg 10kHz (10K samples/second)

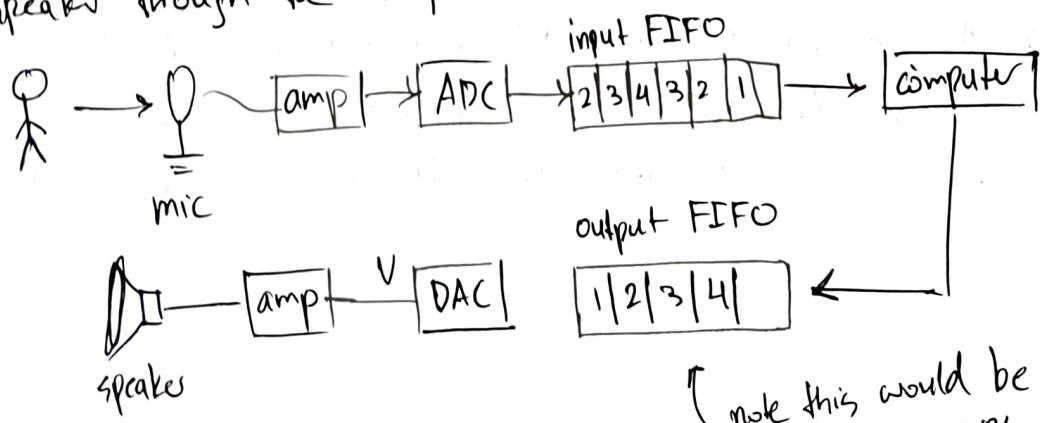
Coming Out of ADC

time	Value		
100μs	1	numbers encoded in binary	can store these in memory
200μs	2		→ permanent memory
300μs	3		
400μs	4		
500μs	3		

- the sound output proceeds in the other direction



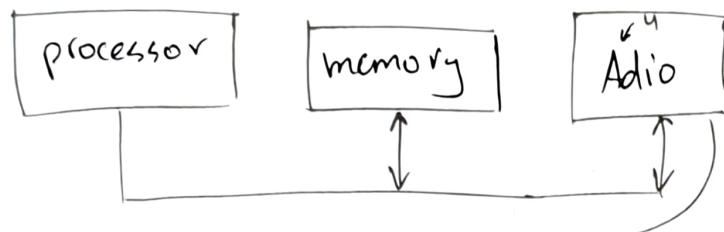
- the DE1-SOC has a microphone input jack and speaker output jack
- there is a CODEC chip connected to the jack that performs all the amplification, and A/D, D/A
- it samples sound input at 8 kHz (125 μs period)
- the output samples also come out @ 8 kHz
- the processor might be busy and miss an input sample (missing up sound) or be late sending an output sample to the D/A
- for both these reasons, the input and output need to be buffered with hardware memory that acts like a queue → a first-in-first-out structure (FIFO)
- called a hardware FIFO
- consider the system that simply connects the microphone to the speaker through the computer



Note this would be mirrored if diagram was drawn in 1 line

Today: Audio Part 2: register level control and 'sound' data flow

- Computer connects to the audio input/output unit through a memory-mapped interface



here are the specific memory-mapped registers that connect to the input and output FIFO's of the audio unit

@ address 0xFF203040

31	9	8	3	2	1	0
WI REI	CW CF	WE PE				

control/
status

31	24	23	16	15	8	7	0
WSLC	WSRC	RALLC	RARC				

FIFO space

31	24	23	Left Data				
31	24	23	Right Data				

Two channels,
Two output FIFO's
Two input FIFO's

1) Left Data and Right Data

- if you load/read from these registers, it gets the next sample from the associated (left or right) FIFO; when that happens, that sample is removed from the input FIFO (so that you get the next sample when you next read)
- if you store/write into these registers, this sends that value as a sample into the left or right output FIFO

- 2) FIFO space H fields tell you how occupied the H FIFO's are
- * RALC is how many of the left input FIFO's entries are occupied
 - * RAPC " " " " right → (of 128)
 - * WSLC " " " " left output FIFO's are empty
 - * WSRC " " " " right " " " "

3) The Control/Status register bits work as follows

- * control : • write a "1" into CR causes the whole (both) input FIFO's to be cleared/emptied
- write a "D" into CW " " " " output " "
- * status: RI goes to 1 if the input FIFO's are $\geq 75\%$ occupied
WI " " " " output FIFO's are $\leq 25\%$ occupied
- * control : RE = 1 means request interrupt when RI = 1
WE = 1 means request an interrupt when WI = 1

Program in C to connect mic to speaker

※ define AUDIO-BASE 0xFF203040

```
int main(void) {
```

```
    volatile int* audio_ptr = (int*)AUDIO-BASE
```

```
    int left, right, fifospace;
```

```
    while(1) {
```

```
        fifospace = *(audio_ptr + 1)
```

```
:
```

adds 4 with pointer arithmetic

cont'd on next page

if ((fifoSpace & 0x000000FF) > 0) { right input
fifo is
not empty

// load the two input channels

left = *(audio_ptr + 2); } same as

right = *(audio_ptr + 3); } lW t0, B(t1)

* (audio_ptr + 2) = left; } 12(t1) contains

* (audio_ptr + 3) = right; } 0xFF203040
store into left and right
data → output FIFO

* can check just the right one (or left)

since they're filled in at ~~the~~ the

same rate

Today: Lab 6 → Making Square Waves + Echo Box + C structure

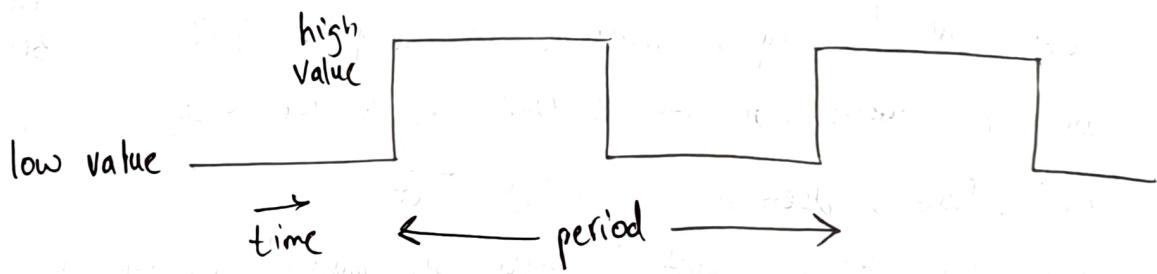
Lab 6:

Part 1: *get used to C in memory-mapped I/O
 • dereferencing pointers + pointer arithmetic
 +1 integer pointer is +4 bytes/addresses

Part 2: * microphone → computer → speaker (code given)

Part 3: * synthesize sound: square wave whose frequency is adjustable with the switches from 100 Hz to 2 kHz (cycles/second)

What is a Square Wave?



low: 0 }

high: 1 1 1 1 ... } 24-bits you figure it out

* need to know the output sample rate to know how many samples to generate in each half period

→ sample rate is 8 kHz

→ period = 125 μs

* how to know when it's okay to send the next sample to the output FIFO → look at WSRC or WSLC in the fifospace register

Part A: *make an echo box

Echo → to both output the current input but also the input that came in some time ago

→ suggest 0.1 seconds ago

→ time is measured in samples → each sample comes in every 125 ms

$$\text{eg: } 1000 \text{ samples} = 125 \text{ ms} = 0.125 \text{ s}$$

→ when echo "comes back" after the delay it is quieter

→ we say it is "damped" by a damping factor D , where $0 < D < 1$

Let's call the input sequence of samples $I(t)$ and the output samples we want to compute $O(t)$

1st $I(0)$
2nd $I(1)$
3rd $I(2)$
 \vdots
100,000 $I(N)$

Aside: Part 2 does this : $O(t) = I(t)$

Now with echo, we want to have the output to include the input + the damped sequence from N samples ago

$$O(t) = I(t) + D \times I(t-N)$$

↖ # of samples
to get 0.1s delay

↖ must store the previous N samples in an array

to get the damped echo $N, 2N, 3N, \dots$ samples ago,

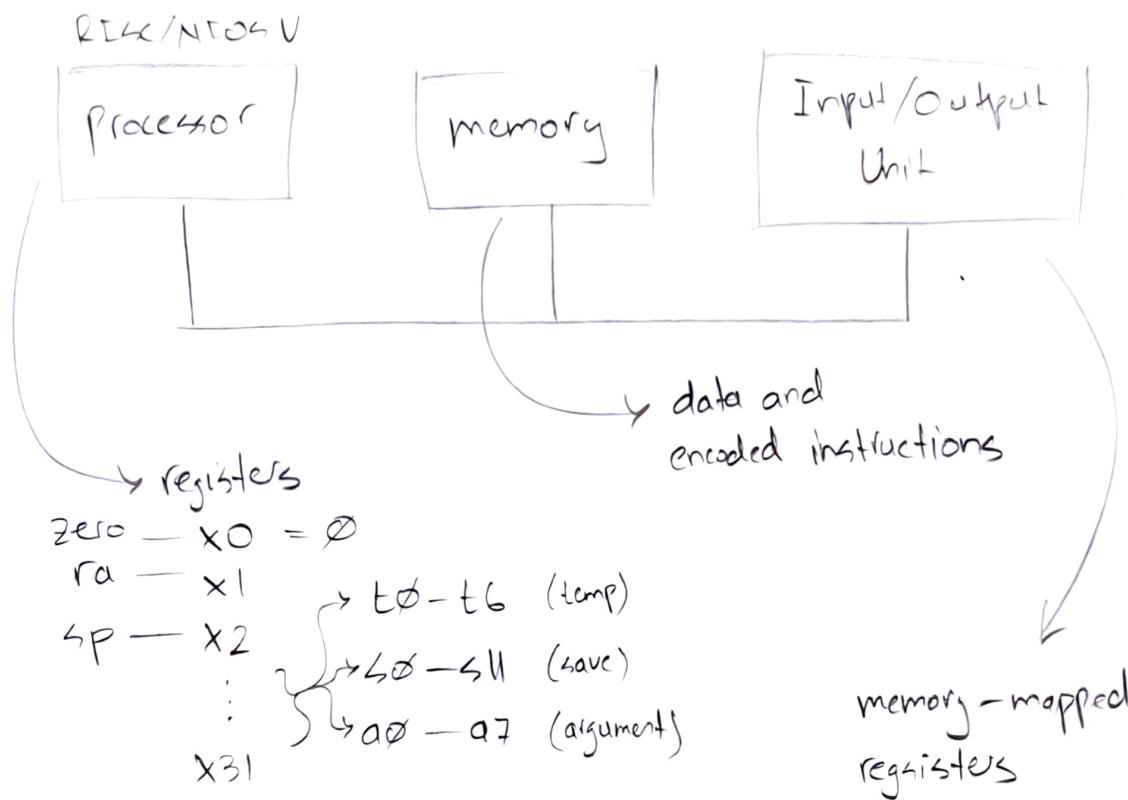
$$\text{store the output } O(t) = I(t) + D \times I(t-N)$$

- only need to store N values of

- + labs 3, 4, 5, 6 → labs 1/2 pretty basic

Review

- Computer Architecture



Memory Organization

- * byte addressability
- + words, half-words, bytes
- order of bytes within a word (little Endian)
- * translate binary to hexadecimal

Instructions for RISC/NIOS V

- + Loads (read) and Stores (write)
 - ↳ lw, lh, lb
 - ↳ sw, sh, sb
 - + indexed addressing
 - ↳ eg: lw \$0, (\$1)
 - + indexed with offset
 - ↳ eg: lb t0, 3(t1)
 - + mv, li, la
 - ↳ can turn into 2 instructions
depending on size of constant
 - + add, addi, sub → no "subi"
 - + shift instructions: srl, srlc, sll, slli
 - + conditional branches
 - ↳ eg: bge x1, zero, loop
 - bxx t4, t3, elsewhere
 - ↳ xx is the condition
 - + unconditional jump: "j"
 - + call } subroutine
 - + ret }
- + cannot access word
at address not divisible
by 4
- + cannot access half word
@ address not divisible
by 2
- + if you jump to
a subroutine instead
of call, it won't
know where to return

Assembler Directives

- * .word .byte .hword .skip (.text .data)
(setup things in memory)

{ know well }

Subroutines

- * how call and ret work
- * how to transmit and receive information to/from a subroutine → MIOS V calling convention.

Convert this program to assembly (obey MIOS V calling conv.)

```
main() {  
    int N = 10;  
    int answer = findsum(N);  
}  
  
int findsum(int N) {  
    int sum = 0;  
    while (N != 0) {  
        sum = sum + N;  
        N = N - 1;  
    }  
    return (sum);  
}
```

- start: la \$t0, 0x20000 * need stack

la \$t0, H

lw \$t0, (\$t0) * t0 holds the 1st argument

call findsum

sw \$t0, H(\$t0)

end: j end

findsum: addi \$sp, \$sp, -8 * push ra & t0 onto stack

lw \$t0, H(\$sp)

lw \$t0, 0(\$sp)

li \$t0, 0 * sum=0

while: beqz \$t0, endloop

add \$t0, \$t0, \$t0 * sum = sum + H

addi \$t0, \$t0, -1 * H = H - 1

j while

endloop: mv \$t0, \$t0 * ans goes into t0

lw \$t0, H(\$sp)

lw \$t0, 0(\$sp)

addi \$sp, \$sp, 8

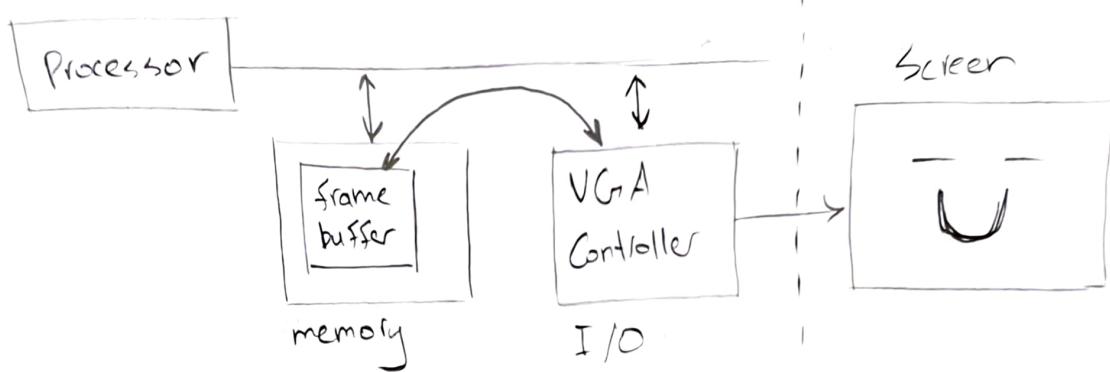
M .word 10

Answer .word 0

rct

Graphics

- * drawing on the screen
- * here is the DE1-SOC graphics system

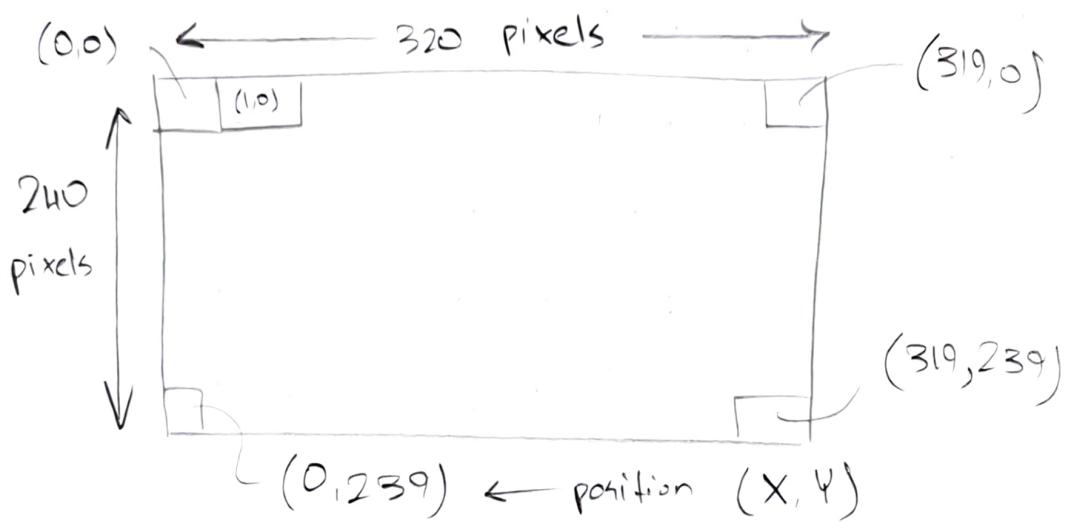


- * software draws the picture you want into the frame buffer, which is specific set of memory locations
- * the VGA controller reads the frame buffer and renders it onto the screen
- * To use this system, you need to know
 - 1) what is a picture made of?
 - * pixels - picture elements
 - * many, all across the screen at different (x,y) location
 - * each pixel's colour is determined by a mixture of red, green, and blue (RGB) - diff. amounts of each get you diff. colours
 - * in the DE1-SOC graphics, each pixel is 16 bits

Blue - 5 bits } organized like this 15 11 10 5 4
 Green - 6 bits }
 Red - 5 bits }

white - all 1's
 black - all 0's
 green - 0x07E0
 blue - 0x001F
 red - 0xF8CO

2) Coordinate System of the screen



3) Need to map (X, Y) location of a pixel to a specific $(1D)$ address in the frame buffer

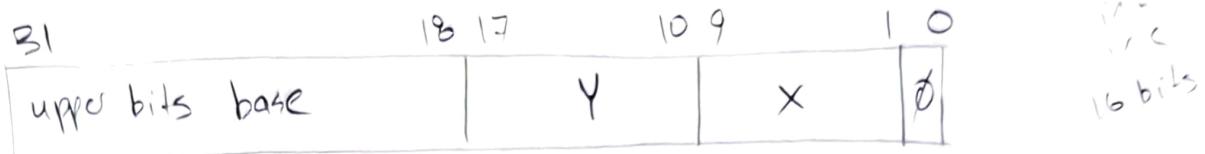
- * Given the (X, Y) , find address

$$0 \leq X \leq 319 \quad 0 \leq Y \leq 239$$

- * need to know where frame buffer begins
(called base address)

- * DE1-SOC Default base: 0x08000000

- + to convert (X, Y) to the address, as follows



- + can compute address given X and Y and base :

$$\text{address} = \text{base} + (Y \ll 10) + (X \ll 1);$$

eg: (X, Y) $\rightarrow @ 0x08000000$
 $(0, 0)$

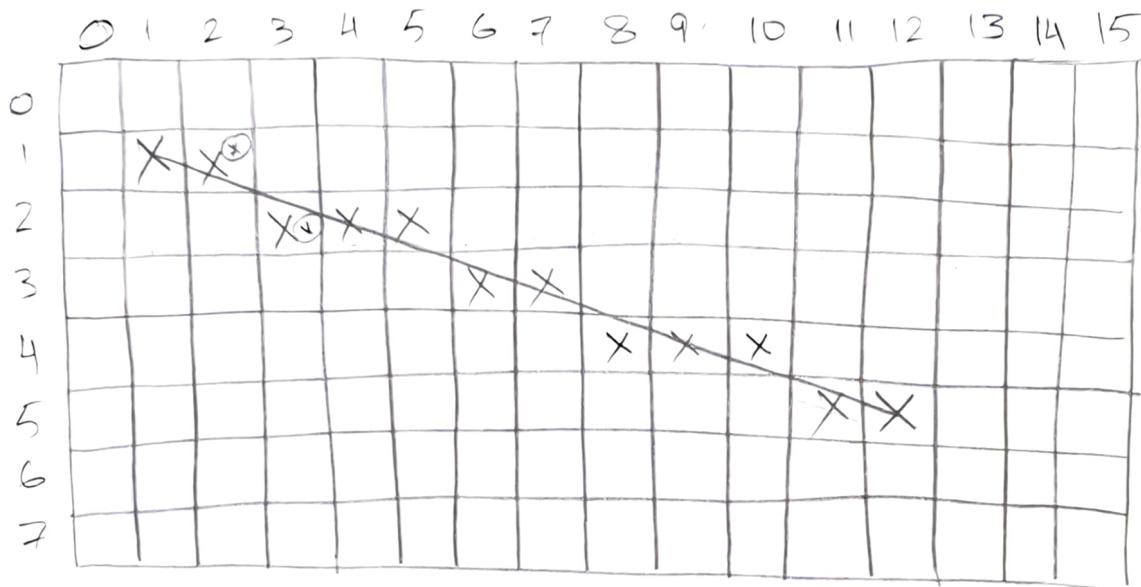
$$(1, 0) \rightarrow @ \text{base} + (0 \ll 10) + (1 \ll 1) = \text{base} + 0 + 2 \\ = 0x08000002$$

$$(1, 1) \rightarrow \dots (1 \ll 10) + (1 \ll 1) \\ = \text{base} + 0x400 + 2 = 0x08000402$$

- + showed the conversion of an (X, Y) location on screen to address in the frame buffer where that pixel's 16-bit colour is stored.

Drawing A Line

- * is harder than you think
- * because of discrete resolution of the screen
- * consider a line from $(X_1, Y_1) = (1, 1)$ to $(X_2, Y_2) = (12, 5)$



- * must choose pixels along the line to turn on
 - ↳ the closest ones
 - ↳ first pixel is given $(1, 1)$ — turn on
 - ↳ next pixel: need to determine the slope of the line
from $(X_1, Y_1) \rightarrow (X_2, Y_2)$: slope = $\frac{Y_2 - Y_1}{X_2 - X_1}$
slope = $\frac{5-1}{12-1} = \frac{4}{11}$ for our line
 - ↳ $X = X_1 = 1$ } turn on pixel $(1, 1)$
 $Y = Y_1 = 1$

↳ $x = 2$ ("walking" one pixel to the right)

$$y = y_1 + \text{slope}(x - x_1) = 1 + \frac{4}{11}(2 - 1) = 1 \frac{4}{11}$$

⇒ round down $1 \frac{4}{11}$ to 1 ⇒ $x = 2, y = 1$ ✘

↳ $x = 3$

$$y = y_1 + \text{slope}(x - x_1) = 1 + \frac{4}{11}(3 - 1) = 1 \frac{8}{11}$$

⇒ round up $1 \frac{8}{11}$ to 2 ⇒ $x = 3, y = 2$ ✘

* there's a problem: steep lines will not be filled in

↳ switch to walking/stepping with Y

↳ called Bresenhan's Algorithm

Lab 7 Part 2: Animation

* horizontal "moving" down, up, down, up ...

* to animate the line

↳ draw white on the line

↳ wait for some ? time

↳ erase line (draw black on it)

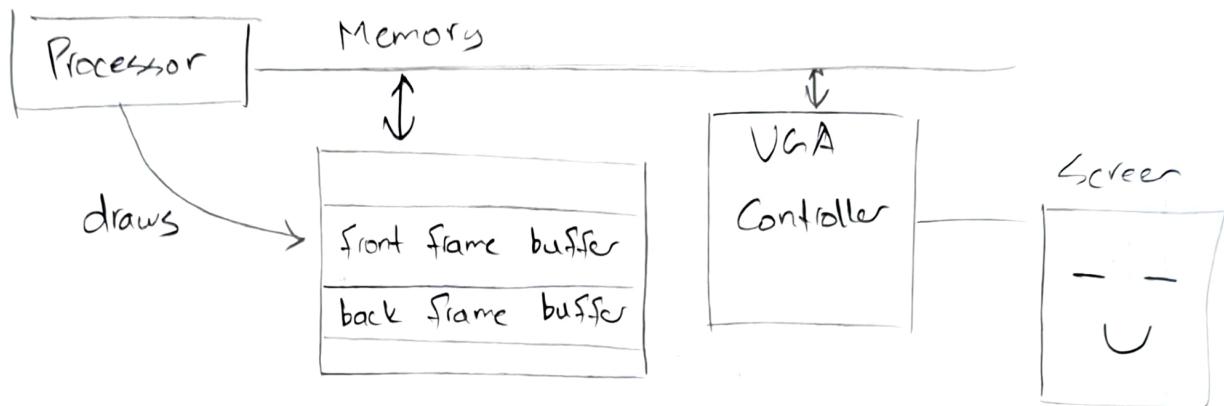
↳ draw white line in a new location

* the VGA controller renders the frame buffer every $\frac{1}{60}$ th second

↳ that line should be moved 1 pixel down/up every $\frac{1}{60}$ second

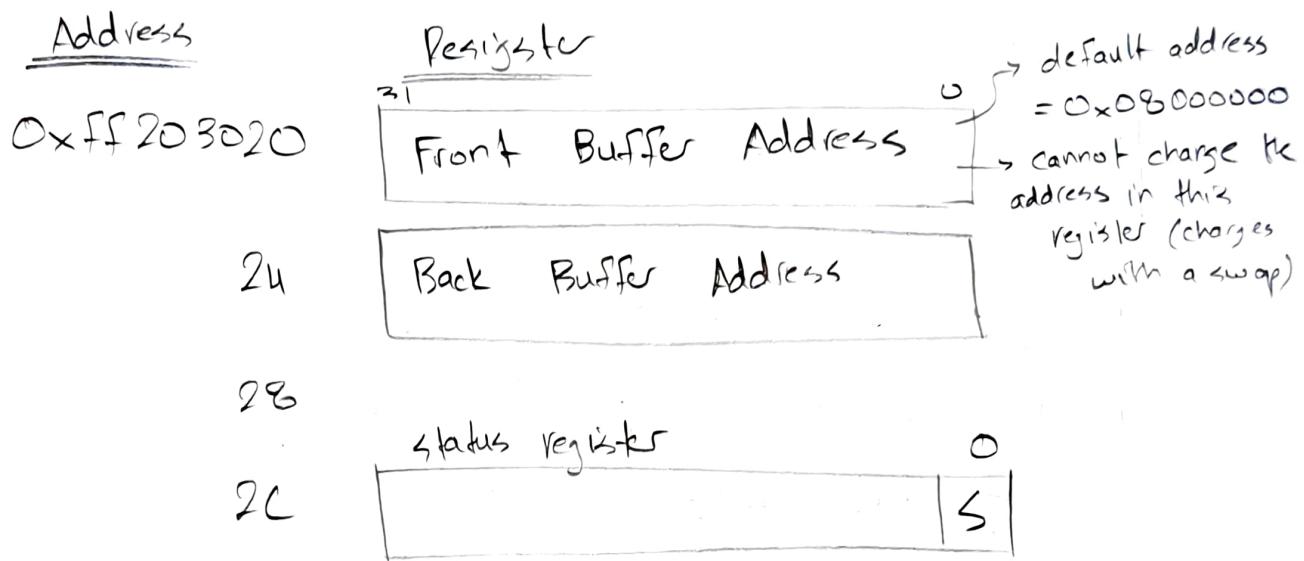
↳ we'll wait for this timing to happen to draw the next line

↳ VGA standard when the rendering is finished a "vertical sync" happens



- * animation is several frames/pictures in sequence
- * Software draws picture/frame into frame buffer
- * VGA controller renders from frame buffer on to screen
-
- * With just one frame buffer, the rendering is disturbed by the drawing
- * per last lecture, we need two frame buffers
 - ↳ One that is being used to do the rendering of the current frame → front buffer
 - ↳ One that is being used to draw the next frame → back buffer
- * Once the rendering is done and the drawing is done, we swap the front buffer and the back buffer
 - ↳ then render the front, then draw the back

- The processor talks to the VGA controller through these mem-mapped registers



- the front and back registers both start = 0x08000000
 - ↳ single buffering
 - to use the second buffer, you must put a different address into the back buffer register.
 - ↳ address must be in range 0x0 → 0xBFFFFF
 - ↳ other addresses are not accessible to VGA controller
- KEY Point: the swap back \leftrightarrow front is just the swap of the addresses in these two registers
 - swap done by VGA controller itself when the time comes
 - The front buffer is rendered every 1/60 of a second

* to cause a swap you do following:

1) Software writes a 1 into the front buffer register
↳ this tells the controller that you want to
do the swap process

↳ causes the S bit in status register (bit 0)
to be set to 1

2) swap doesn't happen right away; controller
waits until rendering done

↳ that point in time is called vertical sync (Vsync)

3) Your code must poll the S bit of status
reg to wait for it to turn into 0

↳ Vsync render done ←

* function to wait for the Vsync

```
Void wait_for_Vsync()
```

```
Volatile int* fbuf = (int*) 0xFF203020 # base  
int status;
```

* fbuf = 1; * start swap cycle per 1)

status = * (fbuf + 3) * read status register

```
while((status & 0x01) != 0) } {  
    status = *(fbuf + 3) } }  
    ← bit = 0
```

- * Lab 7 Part 2 - single buffer Use to know when to change single frame
3 - double buffering * draw new frame

Lab 7 Part 3 goals:

- * draw $N=8$ boxes of size 2×2 pixels in random initial positions on screen
- each box has random colour
- * then:
 - make the boxes move in a random direction and bounce off edges of screen
 - each box should be joined to another (linked in a chain)

Guidance

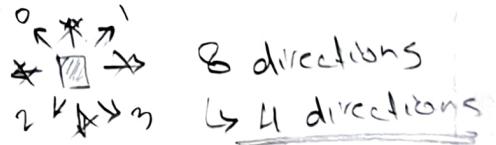
- 1) Use arrays for everything associated with a box

↳ location on screen: int x-box[N], y-box[N];

↳ colour of each box: short int colour-box[N];

↳ direction and amount box moves each frame:

- * int dx[N], int dy[N];
- * only use diagonals
- * only move one pixel per frame in x and y
- * $dx \in \{-1, +1\}$ $dy \in \{-1, +1\}$
- * $2 \times 2 = 4$ directions, randomly choose -1 or +1 for each dx and dy



↳ 8 directions

↳ 4 directions

↳ use rand() to make all the "random" decision

- * part of stdlib → include <stdlib.h>

↳ rand() returns a number 0 and $2^{32}-1$

$$dx = \left(\underbrace{(\text{rand()} \cdot .2)}_{0 \text{ or } 1} \right) \times 2 - 1$$

$\underbrace{\hspace{1cm}}_{0 \text{ or } 2}$

$\underbrace{\hspace{1cm}}_{-1 \text{ or } 1}$

↳ random initial x and y

↳ random colour

* choose carefully

Animation Loop

- + recall: always "draw" (aka set the colour of specific x,y pixels) into the back buffer

// Major outer loop

while () {

 draw(); // draw the frame into the back buffer

 wait-for-Vsync();

 back-buffer = *(buf + 1); // gets the value in back buffer
 register = address of back buffer

}

```
draw() {
```

// erase the previous content of back buffer

① Easy but slow : erase every pixel 320×240

② Harder but faster : erase only what was drawn
into this buffer 2 frames ago

for each box i

```
    draw_box(i);
```

```
    draw_line(i, (i+1) % N);
```

// update box position for next frame

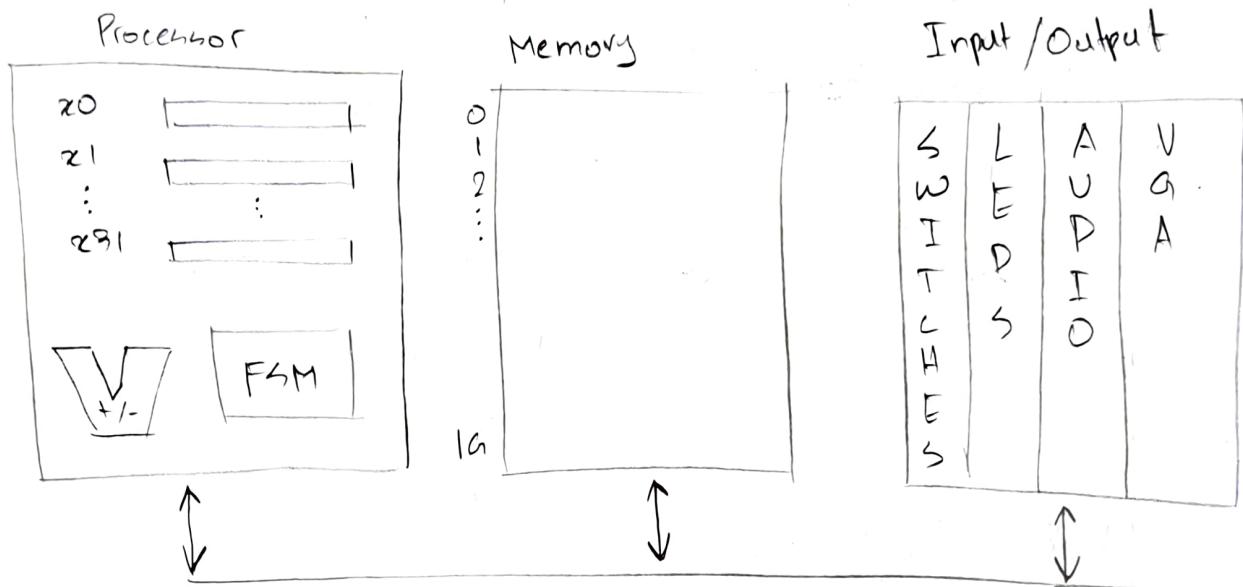
```
    x_box[i] += dx;
```

```
    y_box[i] += dy;
```

// plus check for the screen edges

Digital Design of a Processor

- dive into the processor
- then jump out and talk about some aspects of computer architecture
- computer architecture so far



what
why
how

ECE334

Software
algorithm
hardware
architecture

- discuss digital design of simple processor

↳ has 4 registers - r_0, r_1, r_2, r_3 (V.S. NEON, 32)

↳ each register is 8 bits wide (word size is 8) (V.S. NEON, 32)

↳ has only 256 bytes of memory (V.S. > 1 Gbyte NEON)

↳ the program counter is 8 bits because $2^8 = 256 \leftarrow$ all the memory

↳ 10 instructions : will show the encoding

↳ SIMD instructions have 2 operands

* Z bit = 1 if the result of most recent instruction = \emptyset

$=\emptyset \quad " \neq \emptyset$

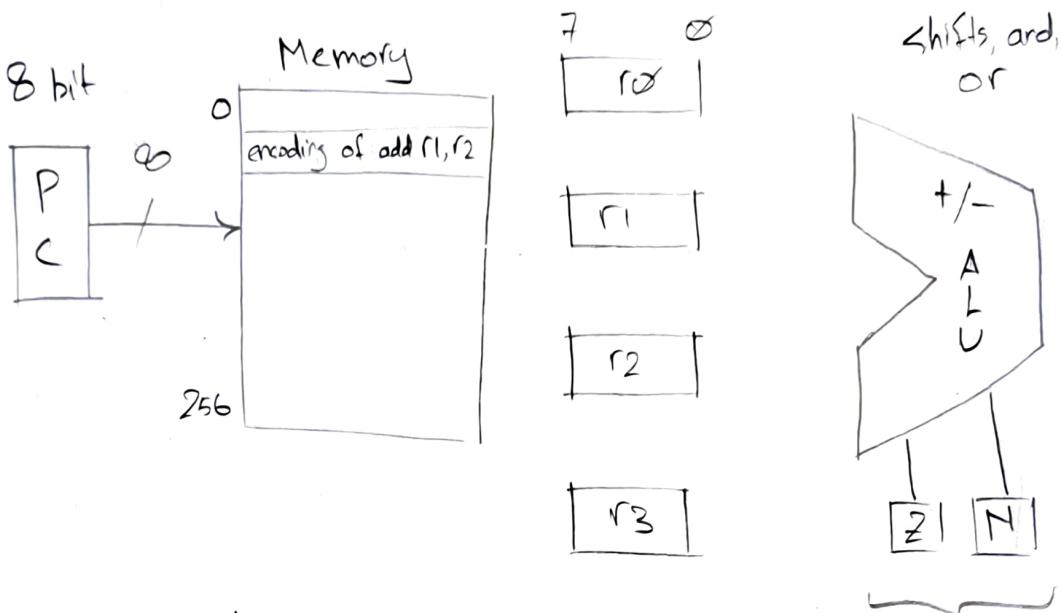
↳ N bit = 1 if result is negative

$=\emptyset \quad " \quad " \quad > \emptyset$

Z and N bit are used in conditional branch instructions

Consider SimP's add instruction

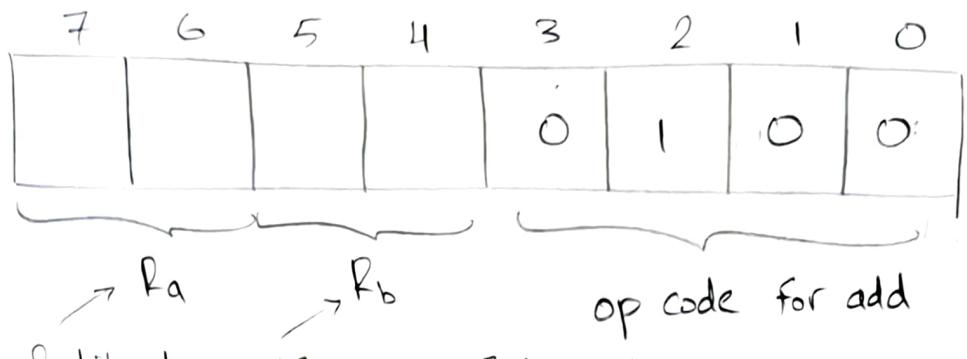
- * eg: add r1, r2 $\Rightarrow r1 \leftarrow r1 + r2$
- * how to build that?



Back to Add

- * general form: add Ra, Rb $\Rightarrow Ra \leftarrow Ra + Rb$
 - * both Ra & Rb can be any of r0, r1, r2, r3
- 1) Fetch instruction from memory
The encoded
 - 2) Figure out what it's supposed to do (in this case, add)
 - 3) Compute $Ra \leftarrow Ra + Rb$
 - 4) if ($Ra == 0$) Z=1 else Z=0
 - 5) if ($R < 0$) N=1 else N=0
 - 6) add 1 to PC: $PC \leftarrow PC + 1$

Instruction Encoding of add Ra, Rb into 8 bits



- * eg: add r1, r2 →

0	1	1	0	0	1	0	0
---	---	---	---	---	---	---	---

 ⇒ 0x64
(hex)

Subtract

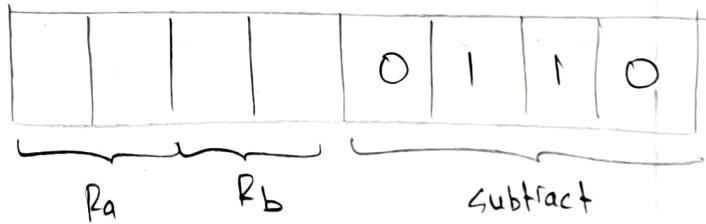
- * sub Ra, Rb ⇔ Ra ← Ra - Rb Encoding

1) Fetches instruction

2) Computes Ra ← Ra - Rb

3) Set Z and N

4) PC ← PC + 1



- * eg: sub r3, r2 →

1	1	1	0	0	1	1	0
---	---	---	---	---	---	---	---

 ⇒ 0xe6
(hex)

NAND Instruction

- * Similar to add, sub

nand Ra, Rb | where $R_a, R_b \in \{R_1, R_2, R_3, R_4\}$

- * Execution of this

- 1) Fetch of instruction from memory

- 2) Compute: $R_a \leftarrow (R_a \& R_b)$ ~~bitwise~~ bitwise NAND

- 3) Condition codes are set : $Z = 1$ if $res = \emptyset$
 $H = 1$ if $res < \emptyset$

- 4) $PC \leftarrow PC + 1$

- * The encoding of NAND is

7	6	5	4	3	... 0
R	a	R	b	1	0 0 0

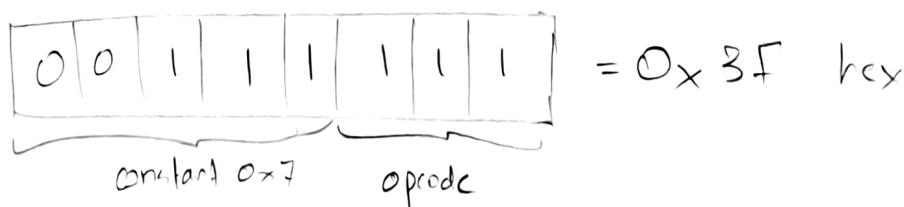
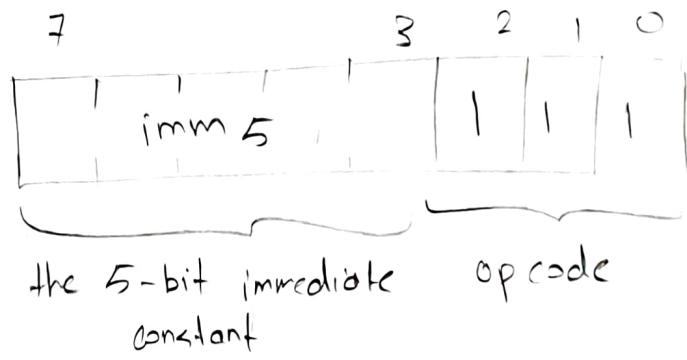
which of the 4 registers are R_a, R_b opcode

eg: nand R2, R1

$$\boxed{1|0|0|1|1|0|0|0|} = 0 \times 98 \text{ (hex)}$$

The OR immediate Instruction

- recall that immediate operands/constants are encoded in the actual instruction
- general form: ori imm5 ← strange, no room to specif reg, so we use only R1 by definition
- Execution:
 - Fetch
 - $R1 \leftarrow R1 \text{ OR } ZE(\text{imm5})$ ← ZE means zero-extend... add 3 zeroes to the left of imm5
 - Set Z and H
 - $PC \leftarrow PC + 1$
- Encoding of ori:
- eg: ori 0x7



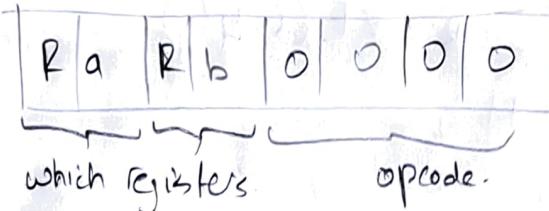
Memory Instructions

→ load Ra, (Rb) | ← similar to lw Ra, (Rb)

Execution:

- 1) $R_a \leftarrow \text{Mem}(R_b)$ ⇒ R_a gets contents of memory address given by R_b
- 2) $PC = PC + 1$ ⇒ Z and H are not set

Encoding:

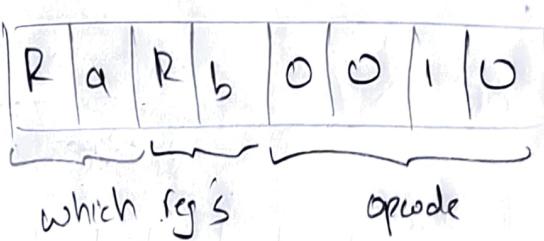


→ store Ra, (Rb) | ← similar to sw Ra, (Rb)

Execution:

- 1) $\text{Mem}(R_b) \leftarrow R_a$
- 2) $PC = PC + 1$

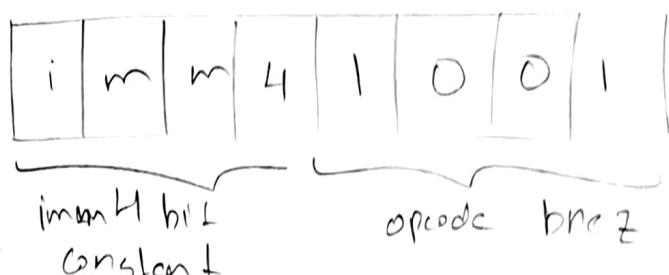
Encoding:



proj.
use of
a proj.

Conditional Branch Instructions

- * Recall in NIOS V bnc x5, x0, loop
- * SimProc: "branch if not zero" (bnz)
 - * "goes" to a diff. address if "z" is not zero
 - ↳ the result of the previous operation
 - ↳ (know this if $z = \emptyset$)
- * general form : 
- * Execution:
 - 1) if ($z == 0$) $PC \leftarrow PC + SE(imm4)$
 - else $PC \leftarrow PC + 1$

SE means sign extend the 4 bit imm4 constant to 8 bits
- * 
- * bpz opcode is 1101 (branch if positive)

1
2 }
3
4 }
5
6 }
7 }
8 }

assemble figures out
what imm4 should
be

$$\text{brnz 2} \leftarrow \text{PC} = 8, \Rightarrow 8 + \text{imm4} = 2 \\ \underline{\text{imm4}} = -6$$

$$P_{\text{new}} = P_{\text{old}} + \Delta E(\text{imm})$$

$$\Theta = 3 + \Delta E(\text{imm})$$

$$\Delta E(\text{imm}) \Rightarrow -3$$

$$(\text{imm}) = -3 = 1101 \Rightarrow \Delta E(\text{imm}) = 1111$$

$$3 \rightarrow -3 \rightarrow \Delta E(-3)$$

$$0011 \rightarrow 1101 \rightarrow 1111 \boxed{1101}$$

Let's Write + Assemble Small Program

- + Subtracts 2 numbers in memory and branches back to beginning if result of sub is not zero

Address	Instruction	Encoding Ra Rb opcode	assume R0=ord R1=rc and mem are set already
0	loop: load R2, (R0)	10 00 0000 = 0x80	
1	load R3, (R1)	11 01 0000 = 0x90	
2	sub R2, R3	10 11 0110 = 0x...	
3	bnez loop	11 01 1001 = 0xD9	

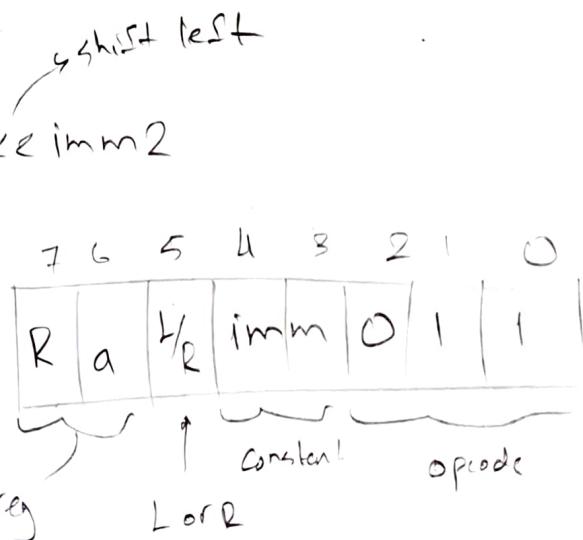
Last Instruction: Shift Left / Right

Shift L/R Ra imm2

- + Execution

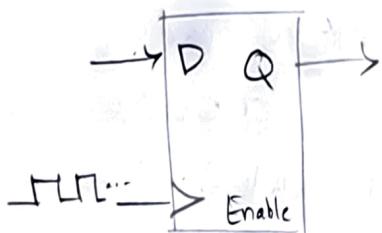
1) if ($L/R == 1$) $Ra = Ra \ll imm2$ else $Ra = Ra \gg imm2$

2) Set Z, H, as usual

3) $PC \leftarrow PC + 1$
 $L/R = 1$ shift left
 $= 0$ " right


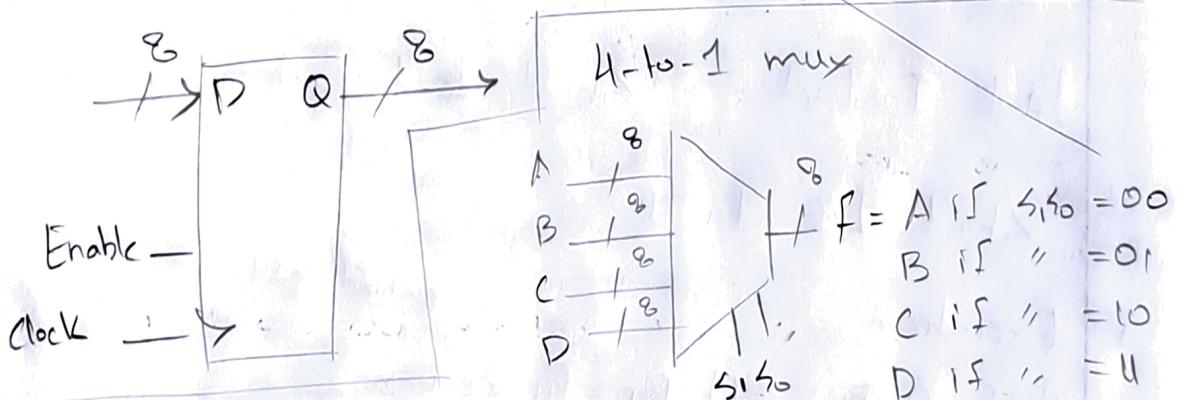
Recall from ECE241

- 1) Recall the flip-flop : if $\text{Enable} = 1$, on a posedge of clock, $Q \leftarrow D$



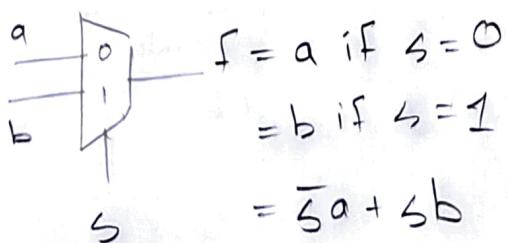
if $\text{Enable} = 0$, Q doesn't change

- * can make an 8-bit register with 8 of these and common Clock and Enable

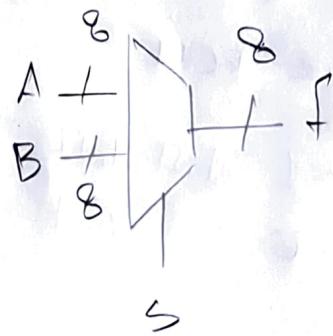


- 2) Recall Multiplexer \rightarrow mux

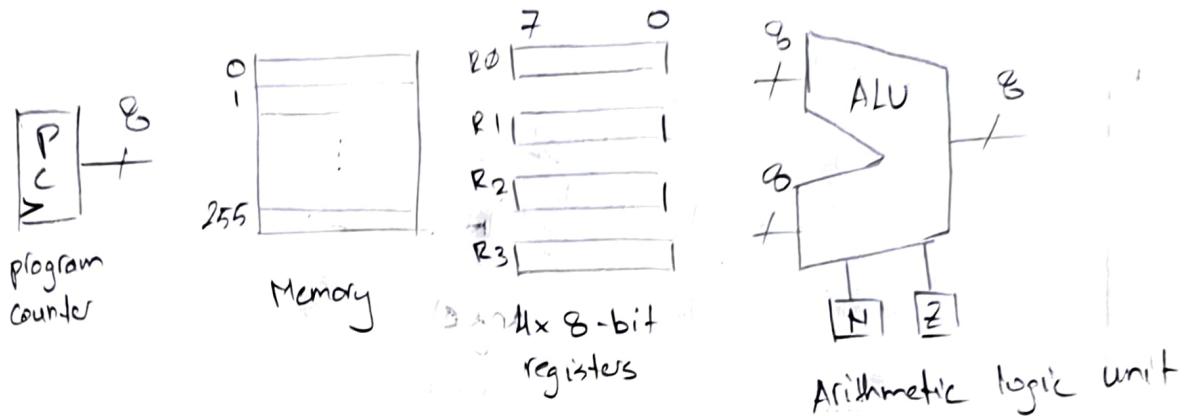
1-bit 2-to-1 mux



8-bit 2-to-1 mux



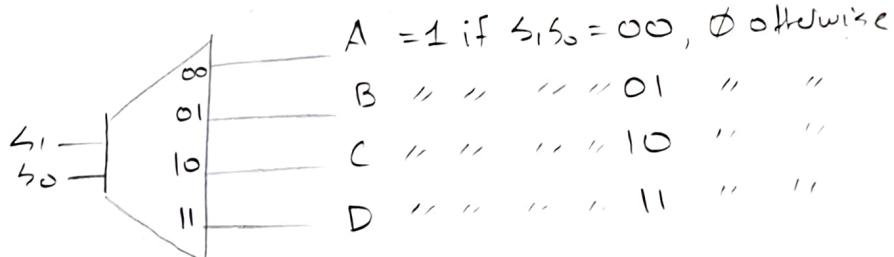
all 8 wires
on 8 2-to-1
1 bit
muxes
selected
by one
signal s

Recall "elements" of ProcessorInstructions :

- + add R_a, R_b * ori imm5
- + sub " " * bnez imm4
- * nand " " * shift R_a, imm2

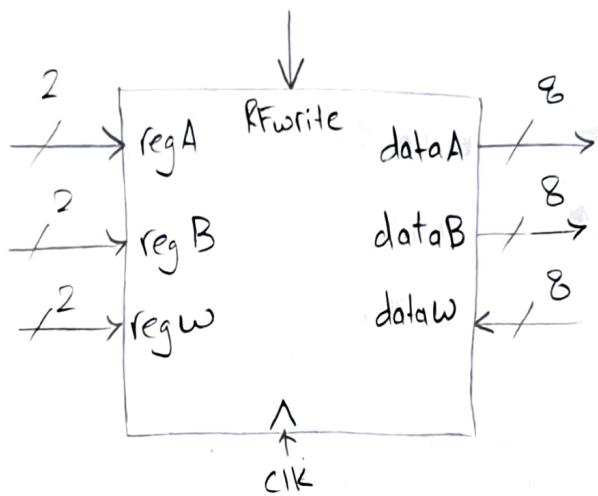
Decoder Ckt

- + here is a 2-to-4 decoder



- + Let's build more detail for the elements
- + begin with registers.

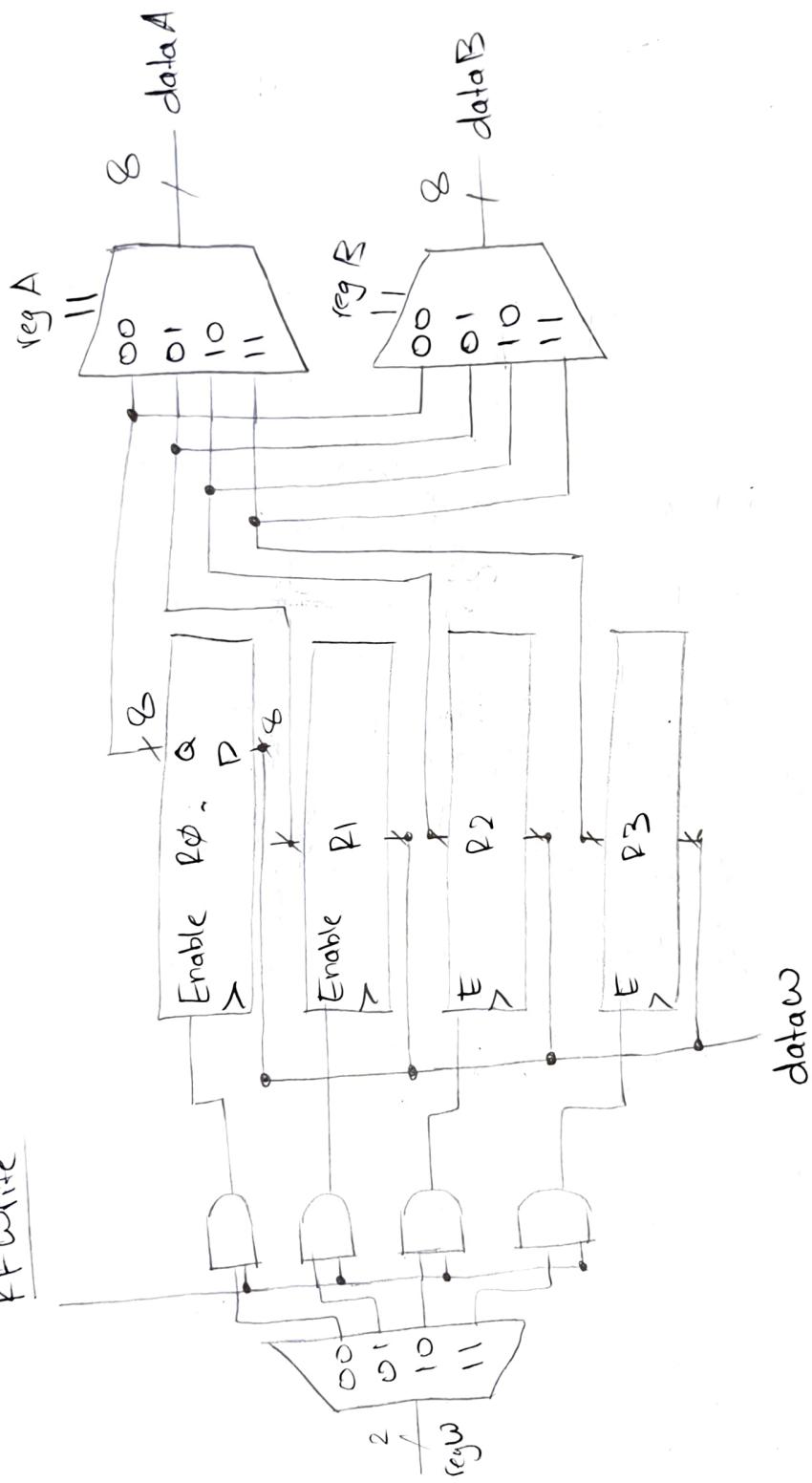
1) Registers \rightarrow Register File (RF)



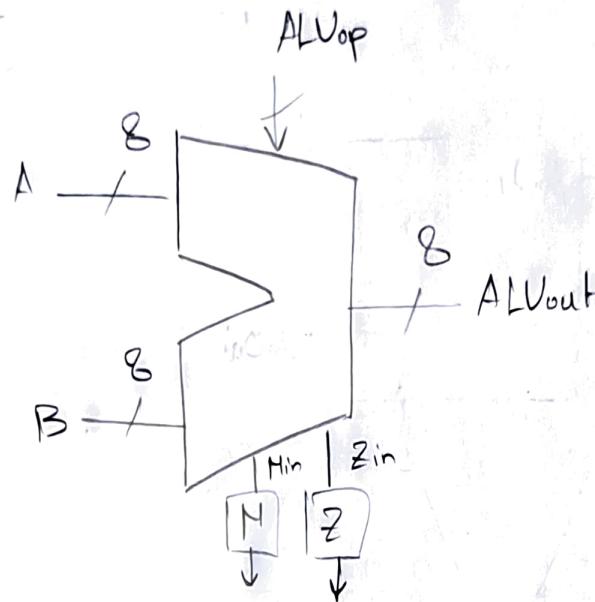
- * regA tells register file which register's data to "put" on dataA
- * regB " " " " " " " " on dataB
- * regW " " " " register to change based on dataW only when RFwrite = 1
- * be careful not to change registers accidentally

Building CKT

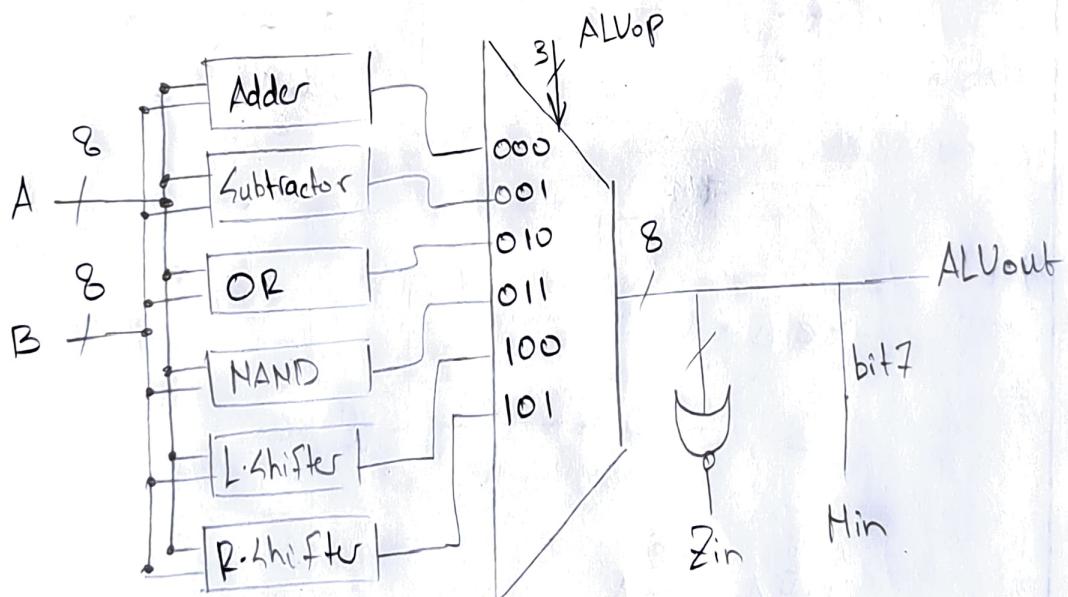
FF write



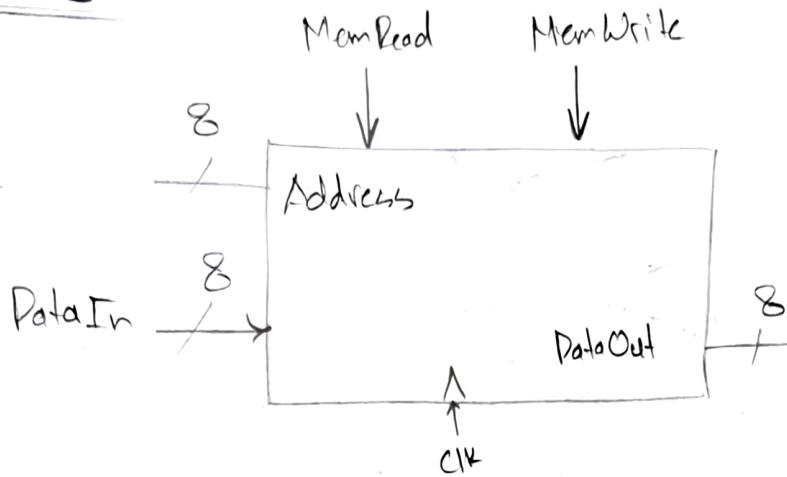
2) ALU : one view



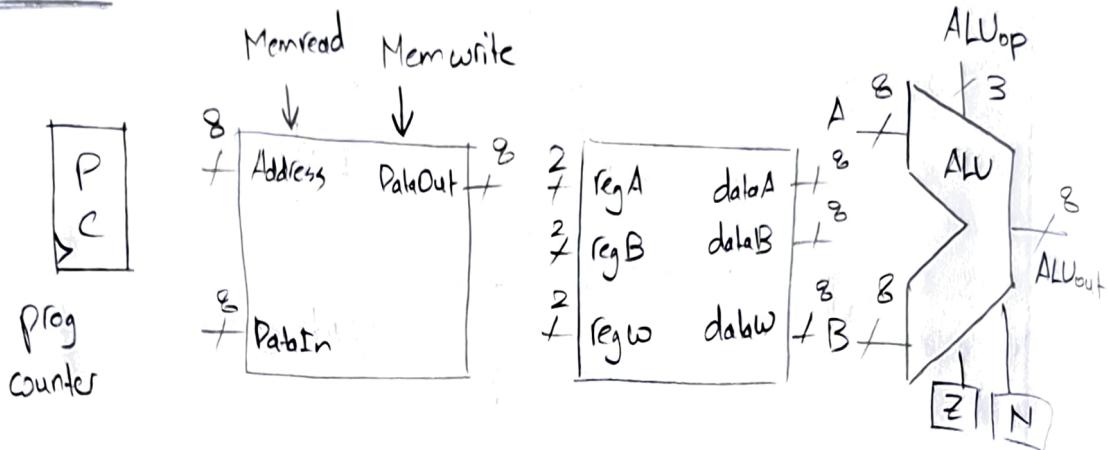
ALUop	ALUout
000	A + B
001	A - B
010	A or B
011	A and B
100	A << B
101	A >> B



3)

Memory

- if MemRead = 1 : . set address
 - on next cycle DataOut gives contents @ that address
- if MemWrite = 1 : . set address , and by end of a cycle
 - contents @ address in memory changes to DataIn

Data PathRecall add instruction

- * add R_a, R_b * R_a ← R_a + R_b

7 6 5 4 3 2 1 0

R _a	R _b	0 1 0 0
----------------	----------------	---------

opcode for "add"

- * here is what happens to do the full execution of this instruction

- 1) Fetch the instruction (encoding) from memory

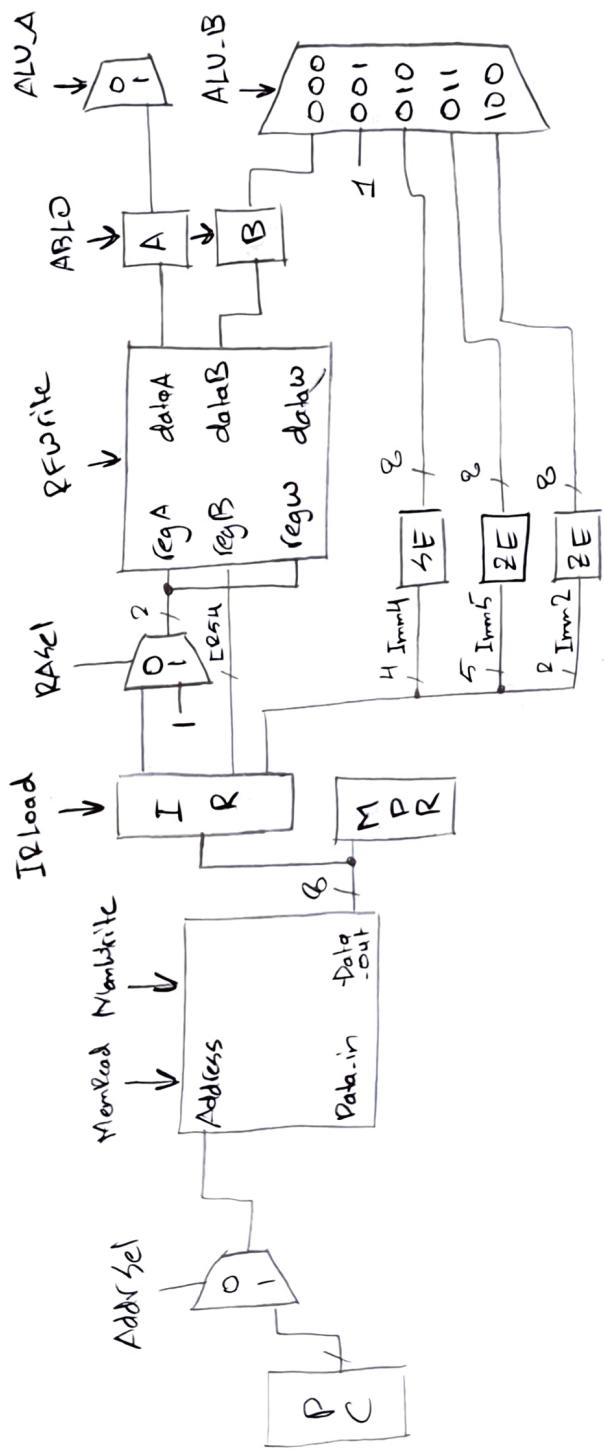
- * put this into the Instruction Register (IR)

- 2) Decode (new word): figure out what the instruction is supposed to do.

- 3) Read R_a and R_b (whichever they are) from the register file → encoded in the instruction

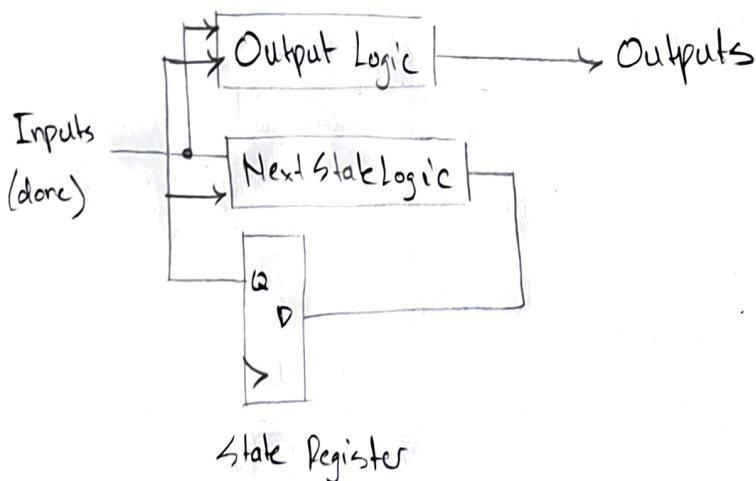
- 4) Compute $ALUout \leftarrow R_a + R_b$ ($ALUout$ is a register)
- 5) Write $ALUout$ back into R_a
- 6) Set Z and N based on $ALUout$
- 7) Compute $PC \leftarrow PC + 1$

incomplete processor
 (didn't get time
 to fully copy
 from board)



Finite State Machine Control of Simple Processor

- * The states of the state machine are just called:
- Cycle 1, Cycle 2, Cycle 3, Cycle 4, Cycle 5
- * This means the next state logic is simple: just proceeds $1 \rightarrow 2 \rightarrow 3 \rightarrow \dots$ until Done
- * General FSM



Encoding of add

add Ra, Rb	7 6 5 4 3 2 1 0
	Ra Rb 0 1 0 0

- * Outputs of FSM:
 - every enable on a register (eg. PC write)
 - every mux select
 - every memory control + ALU op
- * Inputs to FSM:
 - instruction opcode \rightarrow IR [3,2,1,0]
 - the condition codes Z, N (to make branch decision)

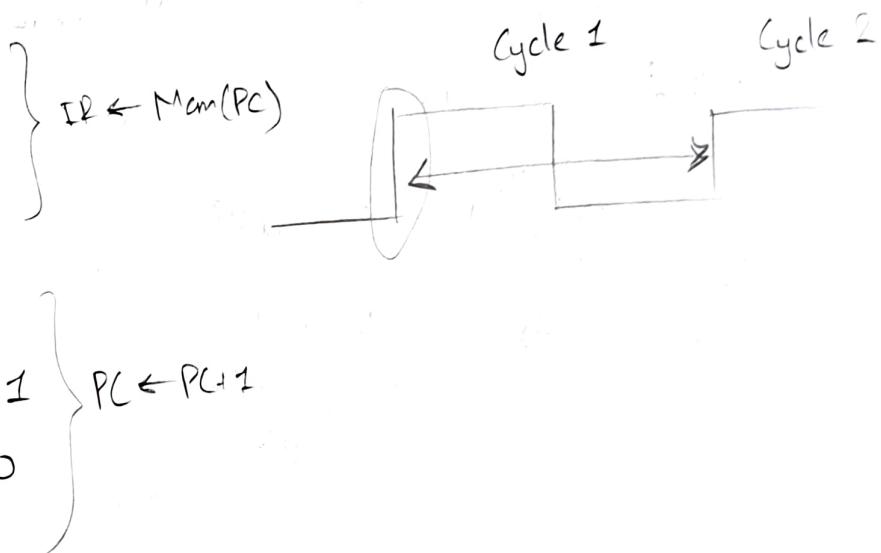
- * only show enables, that we care about
- * In foregoing : only state value of outputs that "are activated"

FSM Control of add, sub, hard

- * During cycle 1: Fetch, $IR \leftarrow \text{Mem}(PC)$ and $PC \leftarrow PC + 1$
 - $\text{AddrSel} = 1$
 - $\text{Mem Read} = 1$
 - $\text{IR Load} = 1$
 - $\text{ALU-A} = 0$
 - $\text{ALU-B} = 001$
 - $\text{ALU op} = 000$
 - $\text{PC Write} = 1$
- * During cycle 2: Decode - there is nothing to do

(but we can do the $R_a R_b$ read from RF $\rightarrow A, B$
because most instructions do this and it will
save a whole cycle of time for them)

- $\text{RA Sel} = \emptyset$
 - $\text{AB LD} = 1$
- $\nwarrow_{\text{AB LD}}$



- * Puring Cycle 3: (assume Add, Sub, Nand)

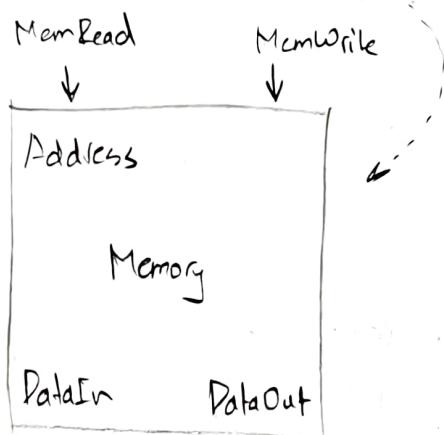
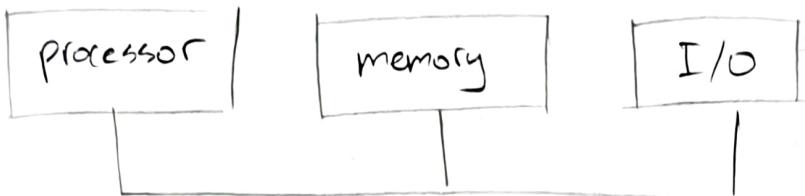
- ALU-A = 1 } selects A
- ALU-B = 000 } selects B
- ALUop = Add or Sub or Nand } depending on instruction
- ALUout_LD = 1 } captures ALUoutput
- FlagWrite = 1 } sets Z and M

- * Puring Cycle 4: (for Add/Sub/Nand)

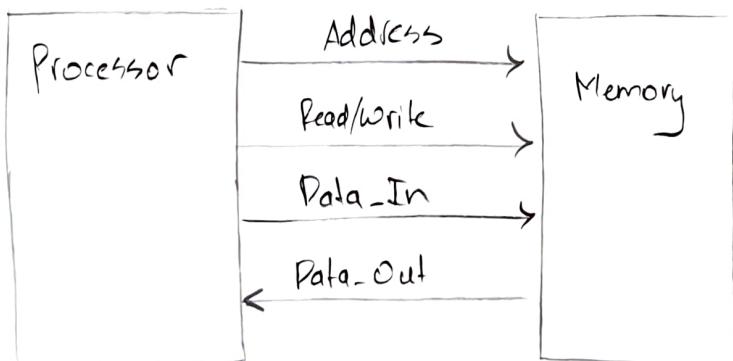
- RegIn = 0 } connects ALUout to dataw
- RFWrite = 1 } change/write into RF
- done = 1

Today: Intro to Cache Memory

- * memory so far: $lwo \quad t1, (t0) \quad t1 \leftarrow (t0)$
 $sws \quad t2, (t3) \quad (t3) \leftarrow t2$



- * Stepping back, the processor connection to memory is this:



- * processor makes 2 kinds of requests to memory
 - ① Read - load instruction or instruction fetch

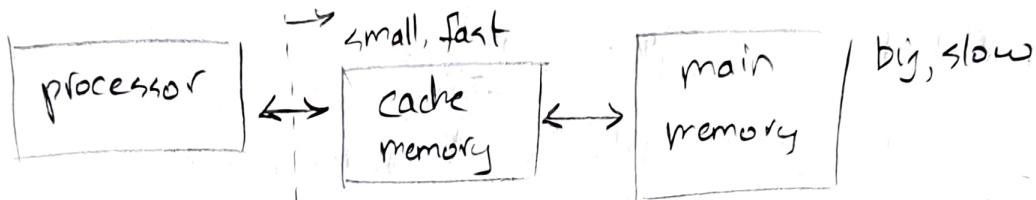
- Processor makes 2 kinds of requests to memory
 - ① Read — load instruction or instruction fetch
 - ② Write — store
- both called a "memory access"
- * we cannot tell from the memory side if a read is a load or an instruction fetch.
- * now, some Physics
 - larger memories (with more bits) are slower. ←
smaller memories are faster (but big memories needed)
 - why? ∵ longer wires \equiv more resistance, capacitance
 - more $R \times C \rightarrow$ more delay
 - also limited by speed of light
 - decoder time for address will get slower

Processors

- * 2-4 GHz clock \rightarrow 0.5 ns - 0.25 ns
- * 16 GB memory system takes \approx 300 cycles to respond to a memory access request

How to solve of fast processor but slow big memory

- * put a small fast memory in between the processor and main memory and keep the (hopefully, much fewer) active memory locations in the small fast memory
- * which is called a Cache Memory



- * need digital "smarts" to manage what is in the small memory
- * works well if memory accesses have "locality of reference"
(use the same memory locations multiple times before moving on) → both instructions and data.
eg: consider instructions in loop:

for ($i=1$; $i \leq 100$; $i++$) {

} } These instructions will be
 executed 100 times

2 kinds ←

- * spatial (addresses are near to each other)
- * temporal (addresses are the same over a time period)

Operation of a Memory Read using Cache

- 1) Processor requests to read data from address A
- 2) Cache (somehow) looks inside to see if it has a copy of the data in main memory address A
 - ↳ if it doesn't: called cache miss
 - ↳ cache requests data @ A from memory
 - ↳ also requests A+1, A+2, A+3 ... (block size)
 - ↳ memory responds slowly and when A's data arrives, cache grabs it and sends on to the processor
 - ↳ if it did contain the data @ A ⇒ called cache hit
 - ↳ can quickly send it to the processor in just a few cycles → good ☺
- 3) If processor then requests A+1, it will come quickly after A because of fetching more than A from memory in a block

hit ratio

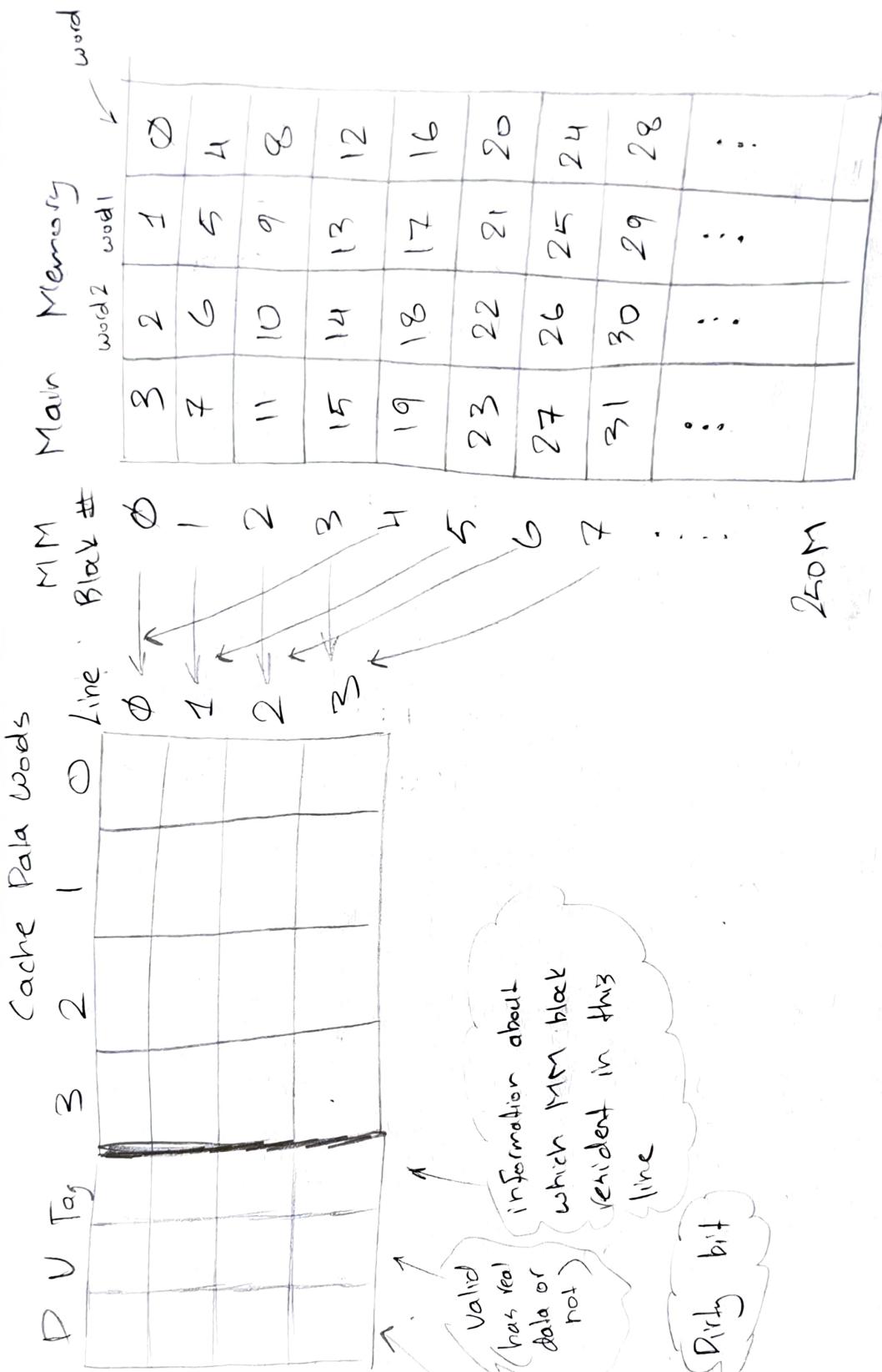
- * when a request from the processor (as an address) is in the cache, called a hit

Bookkeeping and Logistics of Caching

- * First - because the cache asks for more than one memory location at a time, we will think of main memory as organized into blocks of (eg) 4 words (can be any power of 2)
- * aside: for this material, assume every word has an address (v.s. HIOSV-byte addressable)

The Direct-Mapped Cache

- * Cache has rows and columns like MM.
- * each row is called a cache line
- * cache line size = MM block size because we will copy MM blocks into the cache lines
- * consider a (very small) cache with 4 lines and 4 words per line



Valid
has real
information about
data or
which men block
not)
resident in this

Direct-Mapped Cache Cont'd

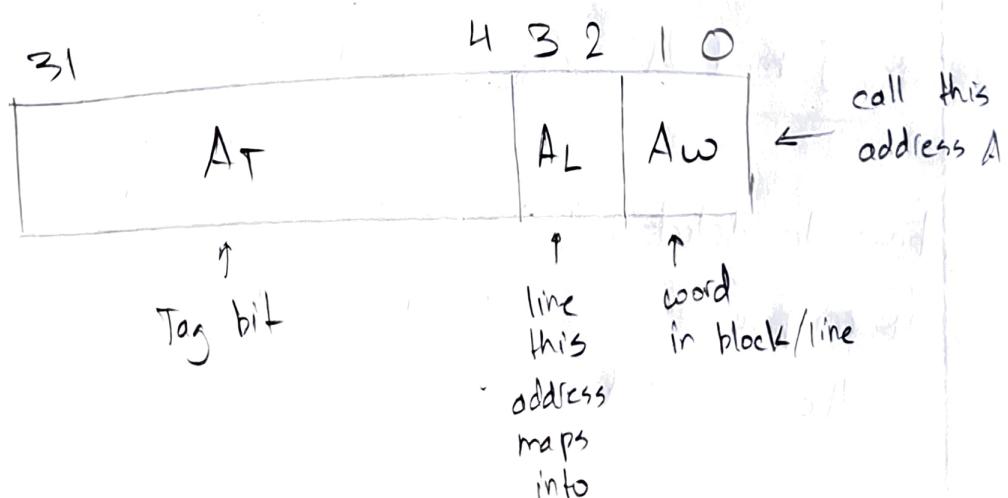
- * now consider when an address (say j) is requested by the processor \rightarrow MM block 2
 - \hookrightarrow where can MM blocks go in cache?
(just 1 line? all lines? some lines?)
- * 1 line \rightarrow "Direct Mapped Cache"
- * good \rightarrow only have to look at one tag/line to see if \exists a hit or miss
- * bad \rightarrow 2 or more MM blocks that use the same line will conflict and evict each other from the cache
- * which line does MM block j "map" into? ($0 \leq j \leq 250M$)
 - $\hookrightarrow j \bmod 4$ (ie: $j \bmod$ # lines in cache)
 - \hookrightarrow remainder after division by 4
- * how to detect hit/miss?
 - \hookrightarrow compare tag of cache line (if valid) with tag of the requested address's block
- * example: assume 32-bit address

$$\text{MM block } j \rightarrow \text{cache line } j \bmod 4$$

of lines
in cache

- * Consider these possible mem. access Request addresses of the first 9 blocks in MM

<u>Block #</u>	<u>Bits of Addr.</u>	<u>00</u>	<u>10</u>	<u>4 words in each block</u>
0	31 30, ... 5, 4, 3, 2, 1, 0 00 ... 0000 XX			
1	00 ... 0001 XX			
2	00 ... 0010 XX			
3	00 ... 0011 XX			
4	00 ... 0100 XX			
5	00 ... 0101 XX			
6	00 ... 0110 XX			
7	00 ... 0111 XX			
8	00 ... 1000 XX			



Direct-Mapped Cache

- * Tag tells us which MM block is in a cache line
- * Valid bit tells us if a cache is actually useful
- * D - Dirty bit
- * AL: the line in the cache that this word / block would reside in if it's in cache
- * Tag: that would be in cache Tag field - $\text{Tag}(A_L)$ if this block was in the cache
- * Consider what happens when the processor makes a read request to the cache at memory address A

↳ Cache looks inside to see if $\text{Tag}(A_L) = A_T$

↳ If $\text{Tag}(A_L) = A_T$ AND $V(A_L) = 1$ (valid), then → hit
(cache sends that word quickly to the processor)

↳ otherwise, → miss

↳ If the line in A_L is valid, it must be evicted because direct-mapped cache only lets a MM block be in one place in cache

↳ The new block is retrieved (slowly) from MM and copied into cache line A_L

↳ The $\text{Tag}(A_L)$ is set to become A_T

↳ Set valid bit = 1 — $V(A_L) = 1$

↳ Set dirty bit = 0 — $D(A_L) = 0$ (not yet dirty)

Eviction of Cache Line

→ Method 1: Write-Back Cache

- * what we've been talking about
- * if there was a write/store to any word in the line, its dirty bit will be set. $D(A_L) = 1$
- * when $D(A_L) = 1$, upon eviction, the line must be "written back" to MM so that MM is correct.
- * multiple writes to a line only results in 1 write to MM upon eviction → goal \checkmark

→ Method 2: Write-Through Cache (aside)

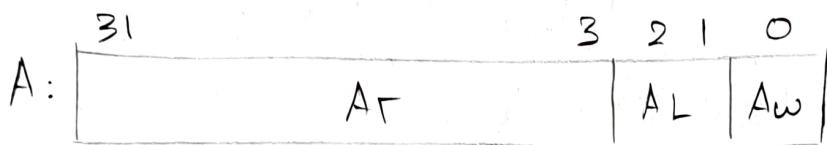
- * does not have a dirty bit
- * all writes to cache are also sent to MM
- * mostly bad because you do more traffic to MM than needed → slow

Example Sequence of Reads for Slightly Diff. Cache

- * 4 lines \rightarrow 2 words /line
- * 32 bit word addressable
- * Cache contains, already

V	Tag	word 1	word 0	line #
1	00 ... 00	0xCA	0xFE	0
1	00 ... 01	0x0E	0xAD	1
1	00 ... 00	0xBE	0xEF	2
0	00 ... 01	0xFE	0xED	3

- * processor read address request



→ different because

- * Binary read address

2 words /lines
words /block

1)

0.....01	01	0
2	1	3
Ar	AL	Aw

 → tag match in line $AL = 01$,
 $V(01) = 1 \rightarrow$ hit \rightarrow get $0xAD$ (word 0 in line 0)

2)

0.....01	11	0
2	3	1
Ar	AL	Aw

 → tag match in line $AL = 11$,
but $V(11) = 0$ (not valid) \rightarrow miss

3)

0....00	00	1
2	1	0
Ar	AL	Aw

 → tag match in line $AL = 00$,
 $V(00) = 1 \rightarrow$ hit \rightarrow get $0xCA$

4)

0....11	01	1
2	1	0
Ar	AL	Aw

 → tag mismatch \rightarrow miss

Recall: Direct Mapped Cache

- * each MM block can only reside in 1 cache line
- * Good: fast to determine if specific block in cache
- * Bad: if blocks want same line at same time \rightarrow conflict
 - \hookrightarrow will evict each other

New Way: Fully Associative Cache

- * every MM block can map into every cache line
- * good: solves the conflict problem in the best way possible
- * Bad: need a lot of digital circuits (comparators) to simultaneously check all tags of all lines

Compromise Between These Two

- * each main memory block can reside in a small set of cache lines
 - \hookrightarrow say "M" cache lines
 - \hookrightarrow called an M-way Set-Associative Cache
 - \hookrightarrow typical values of M are 2-16 (there are typically 1000's of cache lines)

Ex 1: 2-Way Set-Associative Cache

- * allows a main memory to reside in 2 different cache lines \rightarrow those 2 lines are called a set

Ex 2: MIOS-V Cache

- + more realistic
- + for a MIOS-V processor — with byte addressability
- + cache size is 32 Kbytes \rightarrow 32,768 bytes (power of 2)
- + cache has 32 bytes per line = 8 words per line
- + it is an $M=4$ -way set-associative cache
- + To determine the set associative cache mapping / hit or miss, need to compute:

\hookrightarrow Number of lines in the cache: * bytes in cache
* bytes / line

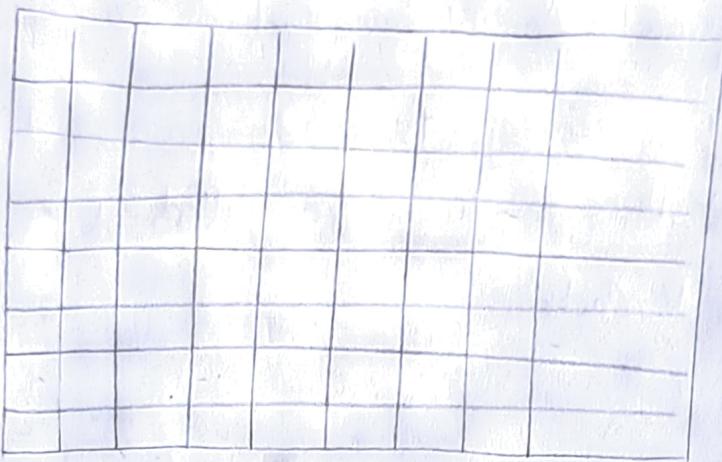
$$\Rightarrow \frac{32\text{Kbytes}}{32\text{ bytes/line}} = 1\text{K} = 1024 \text{ lines}$$

$$\hookrightarrow \text{Number of sets: } \frac{\text{* lines}}{\text{Set size } M} = \frac{1024}{4} = 256$$

\hookrightarrow need $\log_2(256) = 8$ bits to select which set

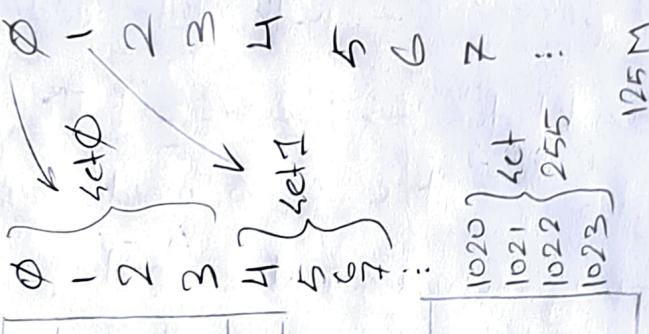
\hookrightarrow 8 words/line \Rightarrow need $\log_2(8) = 3$ bits to determine which word in a line/block is being accessed.

Main Memory



Cache

D N Tag

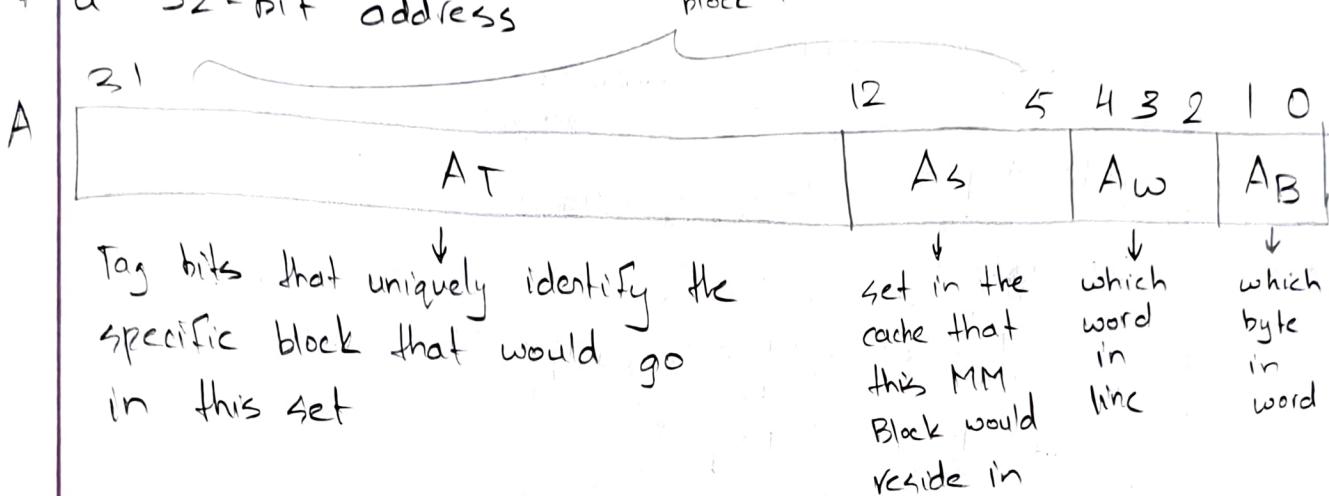


• block $\emptyset \rightarrow$ line $\emptyset \rightarrow$ block $j \rightarrow j \bmod 256$
 $1 \rightarrow \dots$

• random is almost
as good

Ex 2 Cont'd

- * this is in context of the processor requesting to access a memory address A
- + a 32-bit address



- * MM block j can reside in any line in set $j \bmod 256$
- * So, an access of address A to MM works like this:
 - ↳ the cache looks at all 4 tags in set (As) to see if any match AT
 - ↳ if match & V=1 for that line (ie valid), then \Rightarrow hit
 - ↳ else its a miss and must go get the block from MM
- * when a new block is brought in \rightarrow put it in any line that has V=0 (not valid)
 - ↳ if all valid? \rightarrow least recently used line (this is one of many methods)
 - ↳ keep a 2-bit counter for each line & {each line accessed has its counter set to 0; other lines are incremented; if set is full, the lines will have all counters different; evict line with counter = 3 if it is the least recently used one}

Analyzing a Program's Cache Behaviour

- * Consider some contrived code example w.r.t their cache behaviour

Example 1: Conflicting accesses

```
int A[1000], B[1000], sum;  
sum = 0;  
for (int i=0; i<1000; i++) {  
    sum += A[i] + B[i];  
}
```

- * assume a direct-mapped cache with block/line size = 8
- * suppose we're unlucky: every element $A[i]$ maps to the same cache line as $B[i]$.
 - ↳ accessing $B[i]$ after accessing $A[i]$ will kick 8 integers $A[i] \rightarrow A[i+7]$ out of cache
 - ↳ will happen every iteration, destroying the cache's performance
- * you can rewrite code to prevent the problem
 - ↳ eg: sum up the $A[i]$ separately from $B[i]$

Example 2: Accessing 2-Dimensional Arrays (eg. frame buffer)

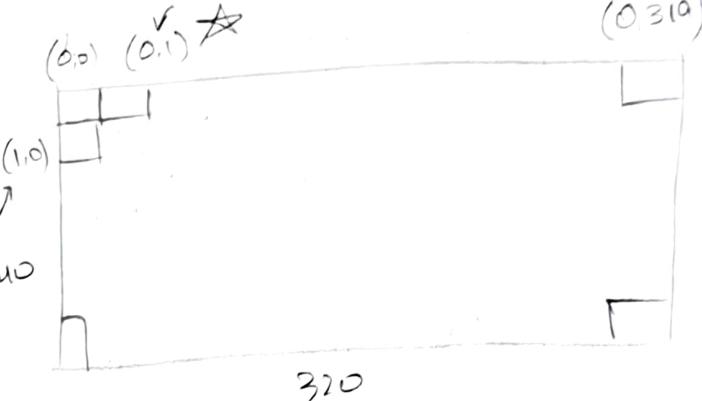
short buf[240][320];

(0,0) is @ address 1000

(0,1) " " " 1002

(0,2) " " " 1004 240

(0,319) " " " 1638



Void bad_clear() {

```
for (int c=0; c<320; c++) {
```

```
    for (int r=0; r<240; r++) {
```

```
        buf[r][c] = 0x0;
```

```
}
```

+ bad because the sequence of addresses is not
in sequential order

+ lots more misses

+ can rewrite to make better

↳ reverse the loop order (next page)

```
void cache-friendly-clear() {  
    for (int r=0; r<240; r++) {  
        for (int c=0; c<320; c++) {  
            buf[r][c] = 0x0;  
        }  
    }  
}
```

- + does access memory sequentially
↳ much better performance

Week 12: Lecture 3

Performance Analysis of Processors with Caches (running a program)

- * processor emits a request to access memory
- * if a memory address hits in the cache, the cache can respond quickly
- * let's say "quickly" means that the cache takes c cycles to respond to a hit (processor cycles)
- * $c \approx 4$ in a modern processor (of modern processors running at 2-5 GHz)
- * if the memory request is a miss, let's say takes m cycles to get the block from main memory and the word requested to the processor.
- * $m \approx 300$ cycles on a modern processor
- * m is called a "miss penalty"
- * lots of misses is bad

Define : hit rate, $h = \frac{\text{# accesses that hit in cache}}{\text{total # of accesses}}$

$$0 \leq h < 1$$

- * Using h, c, m , we can compute the average number of cycles it takes to fulfill a memory request.

- How many cycles would 1000 accesses take?

$$\frac{1000 \times h \times C + 1000 \times (1-h) \times M}{1000}$$

miss rate

- Average Access

Cycles (AAC)

$$AAC = h \times C + (1-h)M$$

- guess $h = 0.9$ (90% hit rate) $C=4$, $M=300$

$$AAC = 0.9 \times 4 + 0.1 \times 300 = 33.6 \text{ cycles} \gg 4$$

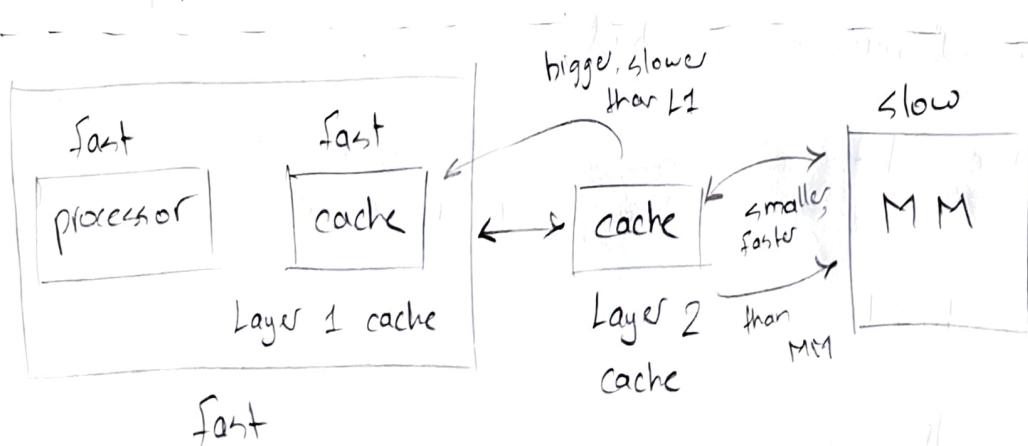
\Rightarrow :: slow!

h	AAC
0.85	48
0.90	34
0.95	19
0.98	9.9
0.99	7.0

hit ratio?

3% difference in miss penalty

is 2x performance



- the layer 2 cache will have ~~x~~ cycles to access faster than MM, slower than cache
 - bigger, better hit rate → speed up things overall
 - "effective" miss penalty seen by L1 cache will be lower
 - can have L3 as well, L4 not common
 - because instruction memory accesses have different "pattern" than data accesses, the L1 cache is often split into
 - ↳ L1 instruction cache
 - ↳ L2 data cache
- CE 453
- architect car design
different caches for each
(diff size, associativity, etc...)

Final Exam Practice

Allan Compiler Questions

1) .global add

```
add: add a0, a0, a1
      ret
```

2) .global sum-to-ten

sum-to-ten: la t0, sum

lw t1, (t0) * t1 < sum

li t2, 1 * int i = 1

li t3, 10

my-loop: add t1, t1, t2 * sum = sum + i

addi t2, t2, 1

ble t2, t3, my-loop

mv a0, t1

ret

at bottom

sum: .word 0

3) .global check-number

check-number: bgz a0, ret-pos

ret-pos: la t0, pos

mv a0, (t0)

ret

blz a0, ret-neg:

la t0, zero

mv a0, (t0)

ret

ret-neg: la t0, neg

mv a0, (t0)

ret

at bottom

.data

pos: .byte 70

neg: .byte 6E

zero: .byte 7A

4)

.global main

main: li a0, 3

li a1, 4

call multiply

ret

multiply: li t0, 0

myloop: add t0, t0; a0

addi a1, a1, -1

bneq a1, myloop

mv a0, t0

ret

* only works if int a and b
are (+)ve integers

Simple Processor

10 Instructions

- + add Ra, Rb * store Ra, (Rb)
- + sub Ra, Rb * shift R/L Ra imm2
- * nand Ra, Rb * bnez imm4
- * ori imm5 * bnez imm4
- * load Ra, (Rb) * bz imm4

Cycle - by - Cycle Control and Datapath

add / sub / nand

Cycle 1	Cycle 2	Cycle 3	Cycle 4	Cycle 5
Addressel = 1	RASel = 0	ALU-A = 1	RegIn = 0	
MemRead = 1	ABLD = 1	ALU-B = 000	RASel = 0	
IRload = 1		ALUop = 000	RFWrite = 1	
ALU-A = 0		ALUoutLD = 1		
ALU-B = 001		Flagwrite = 1		
ALUop = 000				

PCWrite = 1

the above is for add

for sub and nand, exactly the same except

ALUop changes value for sub or nand in cycle 3 //

ori imm5

Cycle 1	Cycle 2	Cycle 3	Cycle 4	Cycle 5
//	//	RASel = 1	ALU-A = 1	RASel = 1
		ABLD = 1	ALU-B = 011	RFWrite = 1
			ALUop = 010	RegFn = 0
			ALUoutLD = 1	
load	//	//	Addressel = 0	RASel = 0
			MemRead = 1	RFWrite = 1
			MDRload = 1	RegIn = 1
store	//	//	Addressel = 0	
			MemWrite = 1	
shift	//	//	ALU-A = 1	RASel = 0
			ALU-B = 100	RFWrite = 1
			ALUop = 100	RegIn = 0
			ALUoutLD = 1	
bnz / bnez / bz	//	//	ALU-A = 0	
			ALU-B = 010	
			ALUop = 000	
			PCWrite = 2	
			= N	
			= Z	
			← bnz	
			← bnez	
			← bz	

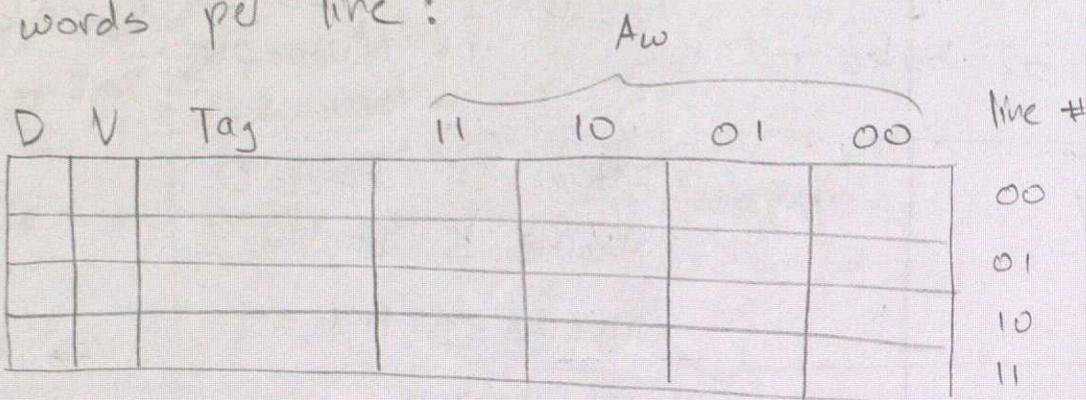
depends
on L or
R bit

Flagwrite = 1

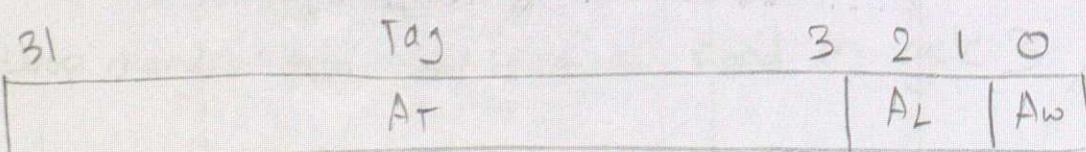
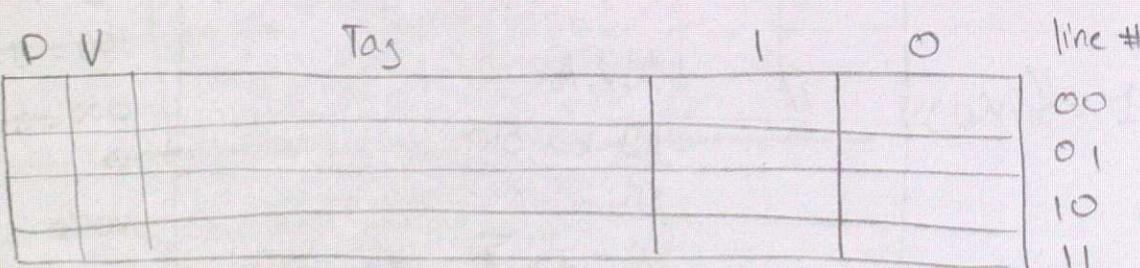
Introduction to Caches

Direct Mapped Cache

Consider a computer system that is word-addressable and has a direct-mapped cache with 4 lines and 4 words per line:



Now consider a 32-bit word-addressable system with a 4-line cache that has 2 words per line



M-Way Set-Associative Cache

Consider a practical, 4-way set associative cache implemented for the MIOS V computer system, a 32-bit byte addressable system. The 32 kbyte cache has 32 bytes per line. To design the cache:

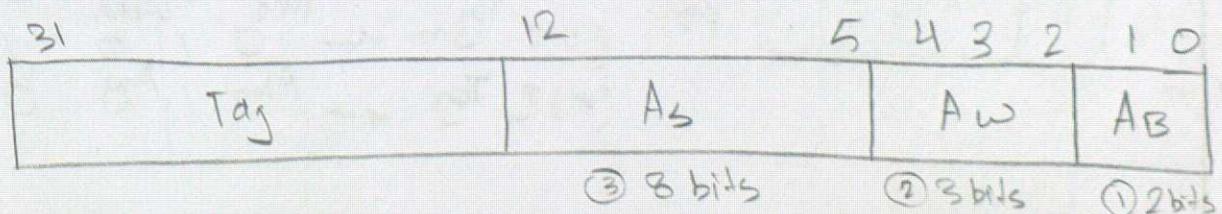
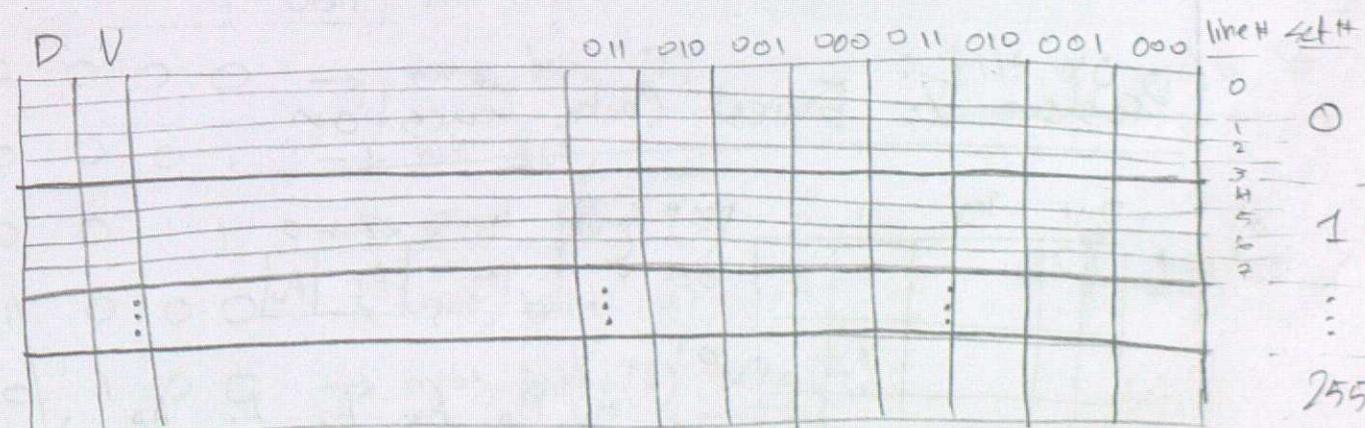
$$\# \text{ words per line} : 32 \text{ bytes} \times \frac{8 \text{ bits}}{\text{byte}} = 8 \text{ words/line} \quad (2)$$

$$\# \text{ lines in cache} : \frac{32 \text{ kbytes}}{32 \text{ bytes/line}} = 1 \text{ K lines} = 1024 \text{ lines}$$

$$\# \text{ sets in cache} : \frac{1024 \text{ lines}}{4 \text{ lines/set}} = 256 \text{ sets in cache} \quad (3)$$

$$\text{already known that } \exists \frac{32 \text{ bits/word}}{8 \text{ bits/byte}} = 4 \text{ bytes/word} \quad (1)$$

$$(1) : \log_2(4) = 2 \quad (2) \log_2(8) = 3 \quad (3) \log_2(256) = 8$$



(3) 8 bits (2) 3 bits (1) 2 bits

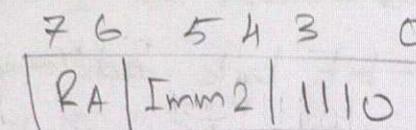
Prof. Rose Final Lecture Questions

Problem 1: Processor

a) MAC RA Imm2 $\Rightarrow RA = RA + RI \times 2^{\text{Imm2}}$

$$\textcircled{1} \quad RA = RA + (RI \ll \text{Imm2})$$

$$\textcircled{2} \quad PC = PC + 1$$



Cycle 1 | Cycle 2 | Cycle 3 | Cycle 4 | Cycle 5 | Cycle 6

Address = 1	RA sel = 0	RA sel = 1	ALU-A = 1	ALU-B = 100	RA sel = 0
MemRead = 1	ABLD = 1	ABLD = 1	ALU-A = 1	RFWrite = 1	
IR Load = 1			ALU op = SH.LF.	ALU-B = 101*	RegIn = 0
ALU-A = 0			ALUout LD = 1	ALUop = ADD	
ALU-B = 001			RA sel = 0	ALUout LD = 1	
ALUop = ADD					
PCWrite = 1					

- b) We need wire going from ALUout to ALU-B mux
and call those select bits 101* (new*)

Problem 2: Easiest Cache Question

a)

D	V	Tag	1	0	line#	4	3	2	1	0
			A ₄	A ₃	A ₂	A ₁	A ₀			
			00			Tag	A ₂	A ₀		
			01							
			10							
			11							

*: $\underbrace{A_4}_{\text{Tag}} \underbrace{A_3}_{\text{line}} \underbrace{A_2}_{\text{word}}$

b) 10 00 0 → miss bring (1) into cache

10 00 1 → hit $\uparrow(1)$

10 01 1 → miss // (3) // //

11 00 0 → miss → evict and bring (4)

00 10 0 → miss bring (5) // //

00 10 1 → hit $\uparrow(5)$

11 01 0
10 01 0

∴ hits are 2 and 6

c) 10 01 0 → this will hit since ref # 3 brought that block
 $\therefore H \rightarrow 10010$

d)

P	V	Tag	1	0	line#	set#	A ₄	A ₃	A ₂	A ₁	A ₀
			5	1		0					
			4	1		0					
			3	2							
			2	3							
			1	3							
			0	3							

Tag set word ✓

e) 10 00 0 → miss, bring (1)

10 00 1 → hit $\uparrow(1)$

10 01 1 → miss, bring (3)

11 00 0 → miss, bring (4)

00 10 0 → miss, evict (1), bring (5)

00 10 1 → hit $\uparrow(5)$

11 01 0 → miss, bring (7)

10 00 1 → hit $\uparrow(3)$...

∴ hits are 2, 6, 8

f) 11110 ✓

Let 1, 1
word 0 ...
empty, go
miss

3. Hardest Cache Question

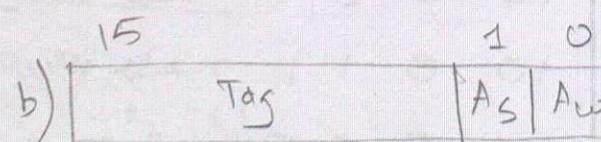
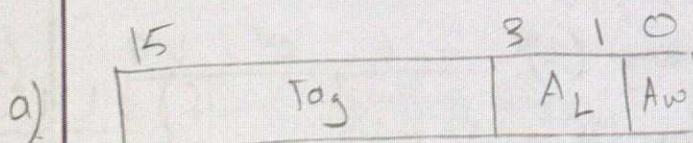
- * address bus 16 bits
- * word size 8 bits (1 byte)
- * cache size 16 bytes (16 words)
- * words per line is 2 bytes (2 words)
- * write-back cache
-
- * 4-way set associative (2 sets total)
- * eviction via LRU model

Direct-Mapped

D	V	Tag	1	0	L
0					0
1					1
2					2
3					3
4					4
5					5
6					6
7					7

4-way Set-Associative

D	V	Tag	1	0	S	L
0					1	0
1					0	1
2					2	2
3					3	3
4					4	4
5					5	5
6					6	6
7					7	7



Assuming word-addressability

(i)

1: Aw=0, AL=000, miss, bring (1)

2: hit \uparrow (1)

3: Aw=1, AL=011, miss, bring (3)

4: Aw=0, AL=111, miss, bring (4)

5: Aw=0, AL=101, miss, bring (5), edit so D=1.

6: Aw=1, AL=111, miss, bring (6), evict (4)

7: Aw=0, AL=001, miss, bring (7), edit so D=1

8: Aw=1, AL=111, miss, bring (8), evict (6)

9: Aw=0, AL=110, miss, bring (9)

10: Aw=0, AL=100, miss, bring (10), edit so D=1

11: Aw=0, AL=010, miss, bring (11)

12: Aw=1, AL=010, hit \uparrow (11)

13: Aw=1, AL=100, hit \uparrow (10)

14: Aw=0, AL=011, miss, bring (14), evict (3)

15: Aw=0, AL=100, miss, bring (15), evict (13)

(ii)

$$\text{hit rate} = \frac{\text{hits}}{\text{total}} = \frac{3}{15} = 0.2$$

hits = 3

misses = 12

	A ₀₀	A ₀₁	
1:	0	0	miss, bring (1)
2:	1	0	hit \uparrow (1)
3:	1	1	miss, bring (3)
4:	0	1	miss, bring (4)
5:	0	1	miss, bring (5), edit so D=1
6:	1	1	miss, bring (6)
7:	0	1	miss, evict LRU(1) = (3), bring (7), edit so D=1
8:	1	1	hit \uparrow (8)
9:	0	0	miss, bring (9)
10:	0	0	miss, bring (10), edit so D=1
11:	0	0	miss, bring (11)
12:	1	0	hit \uparrow (11)
13:	1	0	hit \uparrow (10)
14:	0	1	miss, evict LRU(1) = (5), bring (14)
15:	0	0	miss, evict LRU(1) = (1), bring (12)

ii) $\times \text{ hits} = 4$ hit ratio $= \frac{\text{hits}}{\text{total } \times} = \frac{4}{15} = \underline{\underline{0.267}}$

$\times \text{ misses} = 11$