

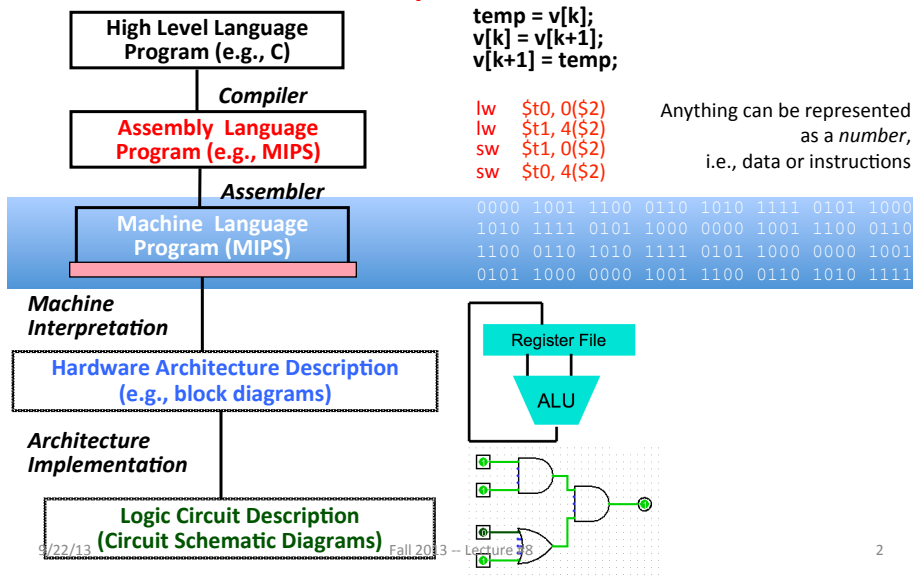
CS 61C: Great Ideas in Computer Architecture *Instructions as Numbers and Floating Point Numbers*

Instructor:

Randy H. Katz

<http://inst.eecs.Berkeley.edu/~cs61c/fa13>

Big Idea #1: Levels of Representation/ Interpretation



Agenda

- Signed vs. Unsigned Numbers
- Instructions as Numbers
- Administrivia
- Instructions as Numbers
- Technology Break
- Floating Point Numbers
- And in Conclusion, ...

9/22/13

Fall 2013 -- Lecture #8

3

What If Operation Result Doesn't Fit in 32 Bits?

- *Overflow*: calculate too big a number to represent within a word
- Unsigned numbers: $1 + 4,294,967,295$ ($2^{32}-1$)
- Signed numbers: $1 + 2,147,483,647$ ($2^{31}-1$)

9/24/13

Fall 2013 -- Lecture #8

4

Depends on the Programming Language

- C unsigned number arithmetic ignores overflow (arithmetic modulo 2^{32})

$$1 + 4,294,967,295 =$$

Depends on the Programming Language

- C unsigned number arithmetic ignores overflow (arithmetic modulo 2^{32})

$$1 + 4,294,967,295 = \text{FFFFFFFF}_{\text{hex}} + 1 = \mathbf{0}$$

Depends on the Programming Language

- C signed number arithmetic also ignores overflow

$$1 + 2,147,483,647 (2^{31}-1) =$$

9/24/13

Fall 2013 -- Lecture #8

7

Depends on the Programming Language

- C signed number arithmetic also ignores overflow

$$1 + 2,147,483,647 (2^{31}-1) = 1 + 7FFFFFFF_{\text{hex}} = 80000000_{\text{hex}} = -2,147,483,648$$

9/24/13

Fall 2013 -- Lecture #8

8

Depends on the Programming Language

- Other languages want overflow signal on signed numbers (e.g., Fortran)
- What's a computer architect to do?

MIPS Solution: Offer Both

- Instructions that can trigger overflow:
 - `add`, `sub`, `mult`, `div`, `addi`, `multi`, `divi`
- Instructions that don't overflow are called "unsigned" (really means "no overflow"):
 - `addu`, `subu`, `multu`, `divu`, `addiu`, `multiu`, `diviu`
- Given semantics of C, always use unsigned versions
- Note: `slt` and `slti` do signed comparisons, while `sltu` and `sltiu` do unsigned comparisons
 - Nothing to do with overflow
 - When would get different answer for `slt` vs. `sltu`?

MIPS Solution: Offer Both

- Instructions that can trigger overflow:
 - add, sub, mult, div, addi, multi, divi
- Instructions that don't overflow are called "unsigned" (really means "no overflow"):
 - addu, subu, multu, divu, addiu, multiu, diviu
- Given semantics of C, always use unsigned versions
- Note: `slt` and `slti` do signed comparisons, while `sltu` and `sltiu` do unsigned comparisons
 - Nothing to do with overflow
 - When would get different answer for `slt` vs. `sltu`?
 - $-1 < 0$ signed, but $FFFF_{\text{hex}} > 0$ unsigned!

9/24/13

Fall 2013 -- Lecture #8

11

Impact of Signed and Unsigned Integers on Instruction Sets

- What (if any) instructions affected?
 - Load word, store word?
 - branch equal, branch not equal?
 - and, or, sll, srl?
 - add, sub, mult, div?
 - `slti` (set less than immediate)?

9/22/13

Fall 2013 -- Lecture #8

12



C provides two sets of operators for AND (& and &&) and two sets of operators for OR (| and ||) while MIPS doesn't. Why?

- Logical operations AND and OR do & and | while conditional branches do && and ||
- The previous statement has it backwards: && and || logical ops, & and | are branches
- They are redundant and mean the same thing: && and || are simply inherited from the programming language B, the predecessor of C

13



C provides two sets of operators for AND (& and &&) and two sets of operators for OR (| and ||) while MIPS doesn't. Why?

- Logical operations AND and OR do & and | while conditional branches do && and ||
- The previous statement has it backwards: && and || logical ops, & and | are branches
- They are redundant and mean the same thing: && and || are simply inherited from the programming language B, the predecessor of C

14

Agenda

- Signed vs. Unsigned Numbers
- **Instructions as Numbers**
- Administrivia
- Instructions as Numbers Continued
- Technology Break
- Floating Point Numbers
- And in Conclusion, ...

9/22/13

Fall 2013 -- Lecture #8

15

Everything in a Computer is Just a Binary Number

- Up to program to decide what data means
 - Example 32-bit data shown as binary number:
0000 0000 0000 0000 0000 0000 0000 0000_{two}
- What does it mean if its treated as
1. Signed integer
 2. Unsigned integer
 3. *(Floating point)*
 4. ASCII characters
 5. Unicode characters
 6. *MIPS instruction*

9/22/13

Fall 2013 -- Lecture #8

16

Implications of Everything is a Number

- *Stored program concept*
 - Invented about 1947 (many claim invention)
- As easy to change programs as to change data!
- Implications?

9/22/13

Fall 2013 -- Lecture #8

17

RISC ISA Design Principles

- #1: Simplicity Favors Regularity
 - Small number of formats, fixed fields
- #2: Smaller is Faster
 - Few instructions, basic primitives
- #3: Good Design Demands Good Compromises
 - Balance usefulness and simplicity
- #4: *Make the Common Case Fast*


9/22/13

Fall 2013 -- Lecture #8

18

MIPS ISA

1. Fold bottom side (columns 3 and 4) together



2. Fold bottom side (columns 3 and 4) together

MIPS Reference Data				①	ARITHMETIC CORE INSTRUCTION SET	②	OPCODE
				/ FUNCT			/ FMT / FT / FUNCT (Hex)
NAME, MNEMONIC	FOR-MAT	OPERATION (in Verilog)	OPCODE (Hex)	NAME, MNEMONIC	FOR-MAT	OPERATION	OPCODE (Hex)
CORE INSTRUCTION SET							
Add	add R	R[rd] = R[rs] + R[rt]	(1) 0/20 _{hex}	Branch On FP True	bclt FI	if(!FPcond)PC=PC+4+BranchAddr	(4) 11/8/1/-
Add Immediate	addi I	R[rt] = R[rs] + SignExtImm	(1,2) 8 _{hex}	Branch On FP False	bclt FI	if(FPcond)PC=PC+4+BranchAddr	(4) 11/8/0/-
Add Imm. Unsigned	addiu I	R[rt] = R[rs] + SignExtImm	(2) 9 _{hex}	Divide	div R	Lo=R[rs] R[rt]; Hi=R[rs]%R[rt]	0/-/-/1a
Add Unsigned	addu R	R[rd] = R[rs] + R[rt]	0/21 _{hex}	Divide Unsigned	divu R	Lo=R[rs] R[rt]; Hi=R[rs]%R[rt]	(6) 0/-/-/1b
And	and R	R[rd] = R[rs] & R[rt]	0/24 _{hex}	FP Add Single	add.s FR	F[fd] = F[fs] + F[ft]	11/10/-/0
And Immediate	andi I	R[rt] = R[rs] & ZeroExtImm	(3) 6 _{hex}	FP Add Double	add.d FR	{F[fd],F[fd+1]} = {F[fs],F[fs+1]} + {F[ft],F[ft+1]}	11/10/-/0
Branch On Equal	beq I	if(R[rs]==R[rt]) PC=PC+4+BranchAddr	(4) 4 _{hex}	FP Compare Single	c.x.s* FR	FPcond = (F[fs] op F[ft]) ? 1 : 0	11/10/-/y
Branch On Not Equal	bne I	if(R[rs]!=R[rt]) PC=PC+4+BranchAddr	(4) 5 _{hex}	FP Compare Double	c.x.d* FR	FPcond = ((F[fs],F[fs+1]) op (F[ft],F[ft+1])) ? 1 : 0	11/11/-/y
Jump	j J	PC=JumpAddr	(5) 2 _{hex}	* (x is eq, lt, or le) (op is ==, <, or <=) (y is 32, 3c, or 3e)			
Jump And Link	jal J	R[31]=PC+8;PC=JumpAddr	(5) 3 _{hex}	FP Divide Single	div.s FR	F[fd] = F[fs] / F[ft]	11/10/-/3
Jump Register	jr R	PC=R[rs]	0/08 _{hex}	FP Divide Double	div.d FR	{F[fd],F[fd+1]} = {F[fs],F[fs+1]} / {F[ft],F[ft+1]}	11/11/-/3
Load Byte Unsigned	lbu I	R[rt]=(24'b0.M[R[rs]+SignExtImm](7:0))	(2) 24 _{hex}	FP Multiply Single	mul.s FR	F[fd] = F[fs] * F[ft]	11/10/-/2
Load Halfword Unsigned	lhu I	R[rt]=(16'b0.M[R[rs]+SignExtImm](15:0))	(2) 25 _{hex}	FP Multiply Double	mul.d FR	{F[fd],F[fd+1]} = {F[fs],F[fs+1]} * {F[ft],F[ft+1]}	11/11/-/2
Load Linked	ll I	R[rt] = M[R[rs]+SignExtImm]	(2,7) 30 _{hex}	FP Subtract Single	sub.s FR	F[fd]=F[fs] - F[ft]	11/10/-/1
Load Upper Imm.	lui I	R[rt] = {imm, 16'b0}	f _{hex}	FP Subtract Double	sub.d FR	{F[fd],F[fd+1]} = {F[fs],F[fs+1]} - {F[ft],F[ft+1]}	11/11/-/1
Load Word	lw I	R[rt] = M[R[rs]+SignExtImm]	(2) 23 _{hex}	Load FP Single	lwc1 I	F[rt]=M[R[rs]+SignExtImm]	(2) 31/-/-/1
Nor	nor R	R[rd] = ~(R[rs] R[rt])	0/27 _{hex}	Load FP Double	ldc1 I	F[rt]=M[R[rs]+SignExtImm]	(2) 35/-/-/1
Or	or R	R[rd] = R[rs] R[rt]	0/25 _{hex}	Move From Hi	mfmhi R	R[rd] = Hi	0/-/-/10
Or Immediate	ori I	R[rt] = R[rs] ZeroExtImm	(3) 4 _{hex}	Move From Lo	mfmlo R	R[rd] = Lo	0/-/-/12
Set Less Than	slt R	R[rd] = (R[rs] < R[rt]) ? 1 : 0	0/2a _{hex}	Move From Control	mfc0 R	R[rd] = CR[rs]	10/0/-/0
				Multiply	mult R	{Hi,Lo} = R[rs] * R[rt]	0/-/-/18
				Multiply Unsigned	multu R	{Hi,Lo} = R[rs] * R[rt]	(6) 0/-/-/19
				Shift Right Arith.	sra R	R[rd] = R[rt] >>> shamt	0/-/-/3
				Store FP Single	swc1 I	M[R[rs]+SignExtImm] = F[rt]	(2) 39/-/-/1
				Store FP Double	swd1 I	M[R[rs]+SignExtImm+4] = F[rt+1]	(2) 3d/-/-/1

http://www.ece.cmu.edu/~ece447/s13/lib/exe/fetch.php?media=mips_reference_data.pdf
 9/22/13 Fall 2013 -- Lecture #8 19

Instructions as Numbers

- Instructions are also kept as binary numbers in memory
 - Stored program concept
 - As easy to change programs as it is to change data
- Register names mapped to numbers
- Need to map instruction operation to a part of number

9/22/13
Fall 2013 -- Lecture #8
20

Names of MIPS fields

op	rs	rt	rd	shamt	funct
----	----	----	----	-------	-------

6 bits 5 bits 5 bits 5 bits 5 bits 6 bits

- *op*: Basic operation of instruction, or *opcode*
- *rs*: 1st register source operand
- *rt*: 2nd register source operand.
- *rd*: register destination operand (result of operation)
- *shamt*: Shift amount.
- *funct*: Function. This field, often called *function code*, selects the specific variant of the operation in the *op* field

9/22/13

Fall 2013 -- Lecture #8

21

Instructions as Numbers

- `addu $t0, $s1, $s2`
 - Destination register `$t0` is register 8
 - Source register `$s1` is register 17
 - Source register `$s2` is register 18
 - Add unsigned instruction encoded as number 33

0	17	18	8	0	33
---	----	----	---	---	----

000000	10001	10010	01000	00000	100001
--------	-------	-------	-------	-------	--------

6 bits 5 bits 5 bits 5 bits 5 bits 6 bits

- Groups of bits call *fields* (unused field default is 0)
- Layout called *instruction format*
- Binary version called *machine instruction*

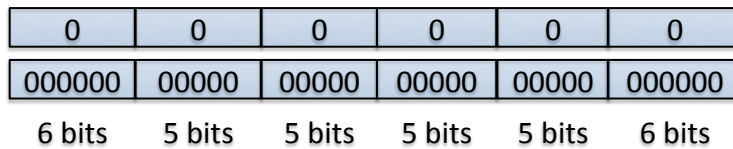
9/22/13

Fall 2013 -- Lecture #8

22

Instructions as Numbers

- `sll $zero, $zero, 0`
 - `$zero` is register 0
 - Shift amount 0 is 0
 - Shift left logical instruction encoded as number 0



- Can also represent machine code as base 16 or base 8 number: `0000 0000hex`, `000000000000oct`

9/22/13

Fall 2013 -- Lecture #8

23

What about Load, Store, Immediate, Branches, Jumps?

- Fields for constants only 5 bits (-16 to +15)
 - Too small for many common cases
- #1 Simplicity favors regularity (all instructions use one format) vs. #3 Make common case fast (justify multiple instruction formats)?
- 4th Design Principle: *Good design demands good compromises*
- Better to have multiple instruction formats and keep all MIPS instructions same *size*
 - All MIPS instructions are 32 bits or 4 bytes

9/22/13

Fall 2013 -- Lecture #8

24

Names of MIPS Fields in I-type

op	rs	rt	address or constant
----	----	----	---------------------

6 bits 5 bits 5 bits 16 bits

- *op*: Basic operation of instruction, or *opcode*
- *rs*: 1st register source operand
- *rt*: 2nd register source operand for branches but register destination operand for *lw*, *sw*, and immediate operations
- *Address/constant*: 16-bit two's complement number
 - Note: equal in size of *rd*, *shamt*, *funct* fields

9/22/13

Fall 2013 -- Lecture #8

25

Register (R), Immediate (I), Jump (J) Instruction Formats

R-type

op	rs	rt	rd	shamt	funct
----	----	----	----	-------	-------

6 bits 5 bits 5 bits 5 bits 5 bits 6 bits

I-type

op	rs	rt	address or constant
----	----	----	---------------------

6 bits 5 bits 5 bits 16 bits

- Now loads, stores, branches, and immediates can have 16-bit two's complement address or constant: $-32,768 (-2^{15})$ to $+32,767 (2^{15}-1)$
- What about jump, jump and link?

J-type

op	address
----	---------

6 bits 26 bits

9/22/13

Fall 2013 -- Lecture #8

26

Encoding of MIPS Instructions: Must Be Unique!

Instruction	Format	op	rs	rt	rd	shamt	funct	address
addu	R	0	reg	reg	reg	0	33 _{ten}	n.a.
subu	R	0	reg	reg	reg	0	35 _{ten}	n.a.
sltu	R	0	reg	reg	reg	0	43 _{ten}	n.a.
sll	R	0	reg	n.a.	reg	constant	0 _{ten}	n.a.
addi unsigned	I	9 _{ten}	reg	reg	n.a.	n.a.	n.a.	constant
lw (load word)	I	35 _{ten}	reg	reg	n.a.	n.a.	n.a.	address
sw (store word)	I	43 _{ten}	reg	reg	n.a.	n.a.	n.a.	address
beq	I	4 _{ten}	reg	reg	n.a.	n.a.	n.a.	address
bne	I	5 _{ten}	reg	reg	n.a.	n.a.	n.a.	address
j (jump)	J	2 _{ten}	n.a.	n.a.	n.a.	n.a.	n.a.	address
jal	J	3 _{ten}	n.a.	n.a.	n.a.	n.a.	n.a.	address
jr (jump reg)	R	0	reg	reg	reg	0	8 _{ten}	n.a.

9/22/13 Fall 2013 -- Lecture #8 27

Converting C to MIPS Machine code

\$t0 (reg 8), &A in \$t1 (reg 9), h=\$s2 (reg 18)

A[300] = h + A[300];



Format?

lw \$t0,1200(\$t1)

addu \$t0,\$s2,\$t0

sw \$t0,1200(\$t1)

Instruction	Format	op	rs	rt	rd	shamt	funct	address
addu	R	0	reg	reg	reg	0	33 _{ten}	n.a.
lw (load word)	I	35 _{ten}	reg	reg	n.a.	n.a.	n.a.	address
sw (store word)	I	43 _{ten}	reg	reg	n.a.	n.a.	n.a.	address
R-type		op	rs	rt	rd	shamt	funct	
I-type		op	rs	rt	address or constant			
J-type		op	address					

9/22/13 Fall 2013 -- Lecture #8 [Student Roulette?](#) 28

Agenda

- Signed vs. Unsigned Numbers
- Instructions as Numbers
- **Administrivia**
- Instructions as Numbers (continued)
- Technology Break
- Floating Point Numbers
- And in Conclusion, ...

9/22/13

Fall 2013 -- Lecture #8

30

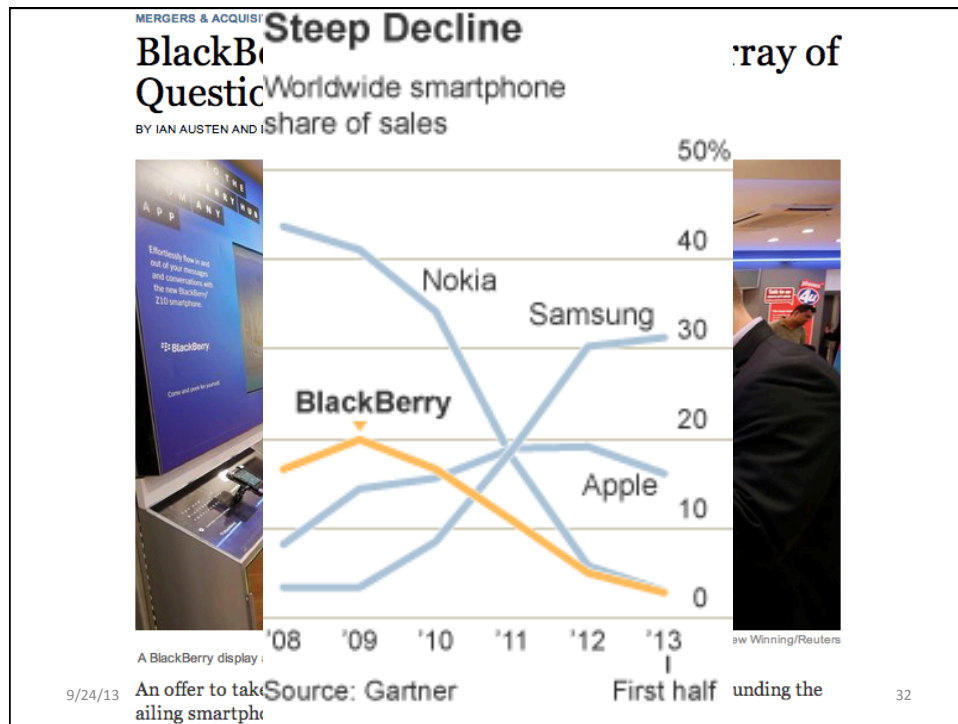
Administrivia

- Tuesday discussion sections henceforth cancelled:
 - Sagar: Tu 9-10 → Friday, 8-9, 405 Soda
 - Sung Roa: Tu 10-11 → Friday, 1-2, 320 Soda
- Piazza Etiquette
 - Please *don't* post your code
 - Please *don't* post copyrighted material
 - Remember *do* look before you post

9/22/13

Fall 2013 -- Lecture #8

31



Agenda

- Signed vs. Unsigned Numbers
- Instructions as Numbers
- Administrivia
- **Instructions as Numbers (cont.)**
- Technology Break
- Floating Point Numbers
- And in Conclusion, ...

Addressing in Branches



- Programs much bigger than 2^{16} bytes, but branch address must fit in 16-bit field
 - Must specify a register for branch addresses for big programs: $PC = \text{Register} + \text{Branch address}$
 - Which register?
- Conditional branching for IF-statement, loops
 - Tend to be near branches; $\frac{1}{2}$ within 16 instructions
- Idea: *PC-relative branching*

9/22/13

Fall 2013 – Lecture #8

34

Addressing in Branches



- Hardware increments PC early, so relative address is $PC = (PC + 4) + \text{Branch address}$
- Another optimization since all MIPS instructions 4 bytes long?
- Multiply value in branch address field by 4!
- MIPS PC-relative branching
 $PC = (PC + 4) + (\text{Branch address} * 4)$

9/22/13

Fall 2013 – Lecture #8

35

Addressing in Jumps



- Same trick for Jumps, Jump and Link
PC = Jump address * 4
- Since PC = 32 bits, and Jump address * 4 = 28 bits, what about other 4 bits?
- Jump and Jump and Link only changes bottom 28 bits of PC

9/22/13

Fall 2013 -- Lecture #8

36

Converting to MIPS Machine code

<i>Address</i>	Loop:	Format?		
800	sll \$t1,\$s3,2	—		
804	addu \$t1,\$t1,\$s6	—		
808	lw \$t0,0(\$t1)	—		
812	bne \$t0,\$s5, Exit	—		
816	addiu \$s3,\$s3,1	—		
820	j Loop	—		

Exit:

R-type	op	rs	rt	rd	shamt	funct
I-type	op	rs	rt	address or constant		
J-type	op address					

9/22/13

Fall 2013 -- Lecture #8

37

32 bit Constants in MIPS

- Can create a 32-bit constant from two 32-bit MIPS instructions
- *Load Upper Immediate* (`lui` or “Louie”) puts 16 bits into upper 16 bits of destination register
- MIPS to load 32-bit constant into register \$s0?
 $0000\ 0000\ 0011\ 1101\ 0000\ 1001\ 0000\ 0000_{\text{two}}$
`lui $s0, 61 # 61 = 0000 0000 0011 1101two`
`ori $s0, $s0, 2304 # 2304 = 0000 1001 0000 0000two`

9/22/13

Fall 2013 -- Lecture #8

39

Agenda

- Signed vs. Unsigned Numbers
- Instructions as Numbers
- Administrivia
- Instructions as Numbers (continued)
- Technology Break
- **Floating Point Numbers**
- And in Conclusion, ...

9/22/13

Fall 2013 -- Lecture #8

40

Goals for Floating Point

- Standard arithmetic for reals for all computers
 - Like two's complement
- Keep as much precision as possible in formats
- Help programmer with errors in real arithmetic
 - $+\infty$, $-\infty$, Not-A-Number (NaN), exponent overflow, exponent underflow
- Keep encoding that is somewhat compatible with two's complement
 - E.g., 0 in Fl. Pt. is 0 in two's complement
 - Make it possible to sort without needing to do floating point comparison

9/22/13

Fall 2013 -- Lecture #8

41

Scientific Notation (e.g., Base 10)

- Normalized *scientific notation* (aka *standard form* or *exponential notation*):
 - $r \times E^i$, E is exponent (usually 10), i is a positive or negative integer, r is a real number ≥ 1.0 , < 10
 - Normalized \Rightarrow No leading 0s
 - 61 is 6.10×10^2 , 0.000061 is 6.10×10^{-5}

9/22/13

Fall 2013 -- Lecture #8

42

Scientific Notation (e.g., Base 10)

- $(r \times e^i) \times (s \times e^j) = (r \times s) \times e^{i+j}$
 $(1.999 \times 10^2) \times (5.5 \times 10^3) = (1.999 \times 5.5) \times 10^5$
 $= 10.9945 \times 10^5$
 $= 1.09945 \times 10^6$
- $(r \times e^i) / (s \times e^j) = (r / s) \times e^{i-j}$
 $(1.999 \times 10^2) / (5.5 \times 10^3) = 0.3634545... \times 10^{-1}$
 $= 3.634545... \times 10^{-2}$
- For addition/subtraction, you first must align:
 $(1.999 \times 10^2) + (5.5 \times 10^3)$
 $= (.1999 \times 10^3) + (5.5 \times 10^3) = 5.6999 \times 10^3$

9/22/13

Fall 2013 -- Lecture #8

43

Which is Less? (i.e., closer to $-\infty$)

- 0 vs. 1×10^{-127} ?
- 1×10^{-126} vs. 1×10^{-127} ?
- -1×10^{-127} vs. 0?
- -1×10^{-126} vs. -1×10^{-127} ?

9/22/13

Fall 2013 -- Lecture #8

44

Floating Point: Representing Very Small Numbers

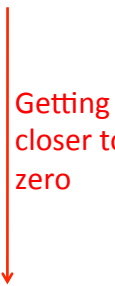
- Zero: Bit pattern of all 0s is encoding for 0.000
 - ⇒ But 0 in exponent should mean most negative exponent (want 0 to be next to smallest real)
 - ⇒ Can't use two's complement ($1000\ 0000_{\text{two}}$)
- *Bias notation*: subtract bias from exponent
 - Single precision uses bias of 127; DP uses 1023
- 0 uses $0000\ 0000_{\text{two}} \Rightarrow 0-127 = -127$;
 ∞ , NaN uses $1111\ 1111_{\text{two}} \Rightarrow 255-127 = +128$
 - Smallest SP real can represent: $1.00\dots00 \times 2^{-126}$
 - Largest SP real can represent: $1.11\dots11 \times 2^{+127}$

9/22/13

Fall 2013 -- Lecture #8

46

Bias Notation (+127)

	How it is interpreted		How it is encoded		
	Decimal Exponent	signed 2's complement	Biased Notation	Decimal Value of Biased Notation	
∞ , NaN 	For infinities		11111111	255	
	127	01111111	11111110	254	
	
	2	00000010	10000001	129	
	1	00000001	10000000	128	
	0	00000000	01111111	127	
	-1	11111111	01111110	126	
	-2	11111110	01111101	125	
	
	-126	10000010	00000001	1	
	For Denorms	10000001	00000000	0	
	Zero				

9/22/13

Fall 2013 -- Lecture #8

47

What About *Real* Numbers in Base 2?

- $r \times E^i$, E where exponent is (2), i is a positive or negative integer, r is a real number ≥ 1.0 , < 2
- Computers version of normalized scientific notation called *Floating Point* notation

9/22/13

Fall 2013 -- Lecture #8

48

Floating Point Numbers

- 32-bit word has 2^{32} patterns, so must be approximation of real numbers ≥ 1.0 , < 2
 - IEEE 754 Floating Point Standard:
 - 1 bit for *sign* (s) of floating point number
 - 8 bits for *exponent* (E)
 - 23 bits for *fraction* (F)
(get 1 extra bit of precision if leading 1 is implicit)
- $$(-1)^s \times (1 + F) \times 2^E$$
- Can represent from 2.0×10^{-38} to 2.0×10^{38}

9/22/13

Fall 2013 -- Lecture #8

49

Floating Point Numbers

- What about bigger or smaller numbers?
- IEEE 754 Floating Point Standard:
 - Double Precision* (64 bits)
 - 1 bit for *sign* (*s*) of floating point number
 - 11 bits for *exponent* (*E*)
 - 52 bits for *fraction* (*F*)
 - (get 1 extra bit of precision if leading 1 is implicit)
- $(-1)^s \times (1 + F) \times 2^E$
- Can represent from 2.0×10^{-308} to 2.0×10^{308}
- 32 bit format called *Single Precision*

9/22/13

Fall 2013 – Lecture #8

50

More Floating Point

- What about 0?
 - Bit pattern all 0s means 0, so no implicit leading 1
- What if divide 1 by 0?
 - Can get infinity symbols $+\infty$, $-\infty$
 - Sign bit 0 or 1, largest exponent, 0 in fraction
- What if do something stupid? ($\infty - \infty$, $0 \div 0$)
 - Can get special symbols NaN for Not-a-Number
 - Sign bit 0 or 1, largest exponent, not zero in fraction
- What if result is too big? ($2 \times 10^{308} \times 2 \times 10^2$)
 - Get *overflow* in exponent, alert programmer!
- What if result is too small? ($2 \times 10^{-308} \div 2 \times 10^2$)
 - Get *underflow* in exponent, alert programmer!

9/22/13

Fall 2013 – Lecture #8

51

Floating Point Add Associativity?

- $A = (1000000.0 + 0.000001) - 1000000.0$
- $B = (1000000.0 - 1000000.0) + 0.000001$
- In single precision floating point arithmetic, A does not equal B
 - $A = 0.000000$, $B = 0.000001$
- Floating Point Addition is not Associative!
 - Integer addition is associative
- When does this matter?

9/22/13

Fall 2013 -- Lecture #8

52

MIPS Floating Point Instructions

- C, Java has single precision (`float`) and double precision (`double`) types
- MIPS instructions: `.s` for single, `.d` for double
 - Fl. Pt. Addition single precision:
Fl. Pt. Addition double precision:
 - Fl. Pt. Subtraction single precision:
Fl. Pt. Subtraction double precision:
 - Fl. Pt. Multiplication single precision:
Fl. Pt. Multiplication double precision:
 - Fl. Pt. Divide single precision:
Fl. Pt. Divide double precision:

9/22/13

Fall 2013 -- Lecture #8

53

MIPS Floating Point Instructions

- C, Java have single precision (`float`) and double precision (`double`) types
- MIPS instructions: `.s` for single, `.d` for double
 - Fl. Pt. Comparison single precision:
Fl. Pt. Comparison double precision:
 - Fl. Pt. branch:
- Since rarely mix integers and Floating Point, MIPS has separate registers for floating-point operations: `$f0`, `$f1`, ..., `$f31`
 - Double precision uses adjacent even-odd pairs of registers:
 - `$f0` and `$f1`, `$f2` and `$f3`, `$f4` and `$f5`, ..., `$f30` and `$f31`
- Need data transfer instructions for these new registers
 - `lwc1` (load word), `swc1` (store word)
 - Double precision uses two `lwc1` instructions, two `swc1` instructions

9/22/13

Fall 2013 – Lecture #8

55

MIPS Floating Point Instructions

- C, Java have single precision (`float`) and double precision (`double`) types
- MIPS instructions: `.s` for single, `.d` for double
 - Fl. Pt. Comparison single precision: `c.lt.s` (`eq`, `ne`, `lt`, `le`, `gt`, `ge`)
Fl. Pt. Comparison double precision: `c.lt.d`
 - Fl. Pt. branch: `bclt`, `bclf`
- Since rarely mix integers and Floating Point, MIPS has separate registers for floating-point operations: `$f0`, `$f1`, ..., `$f31`
 - Double precision uses adjacent even-odd pairs of registers:
 - `$f0` and `$f1`, `$f2` and `$f3`, `$f4` and `$f5`, ..., `$f30` and `$f31`
- Need data transfer instructions for these new registers
 - `lwc1` (load word), `swc1` (store word)
 - Double precision uses two `lwc1` instructions, two `swc1` instructions

9/22/13

Fall 2013 – Lecture #8

56

Peer Instruction Question

Suppose Big, Tiny, and BigNegative are floats in C, with Big initialized to a big number (e.g., age of universe in seconds or 4.32×10^{17}), Tiny to a small number (e.g., seconds/femtosecond or 1.0×10^{-15}), BigNegative = - Big.

Here are two conditionals:

I. $(\text{Big} * \text{Tiny}) * \text{BigNegative} == (\text{Big} * \text{BigNegative}) * \text{Tiny}$

II. $(\text{Big} + \text{Tiny}) + \text{BigNegative} == (\text{Big} + \text{BigNegative}) + \text{Tiny}$

Which statement about these is correct?

Orange. I. is false and II. is false

Green. I. is false and II. is true

Pink. I. is true and II. is false

Yellow. I. is true and II. is true

9/22/13

Fall 2013 -- Lecture #8

57

Pitfalls

- Floating point addition is NOT associative
- Some optimizations can change order of floating point computations, which can change results
- Need to ensure that floating point algorithm is correct even with optimizations

9/22/13

Fall 2013 -- Lecture #8

60

“And in Conclusion, ...”

- Program can interpret binary number as unsigned integer, two's complement signed integer, floating point number, ASCII characters, Unicode characters, ... even instructions!
- Floating point is an approximation of reals