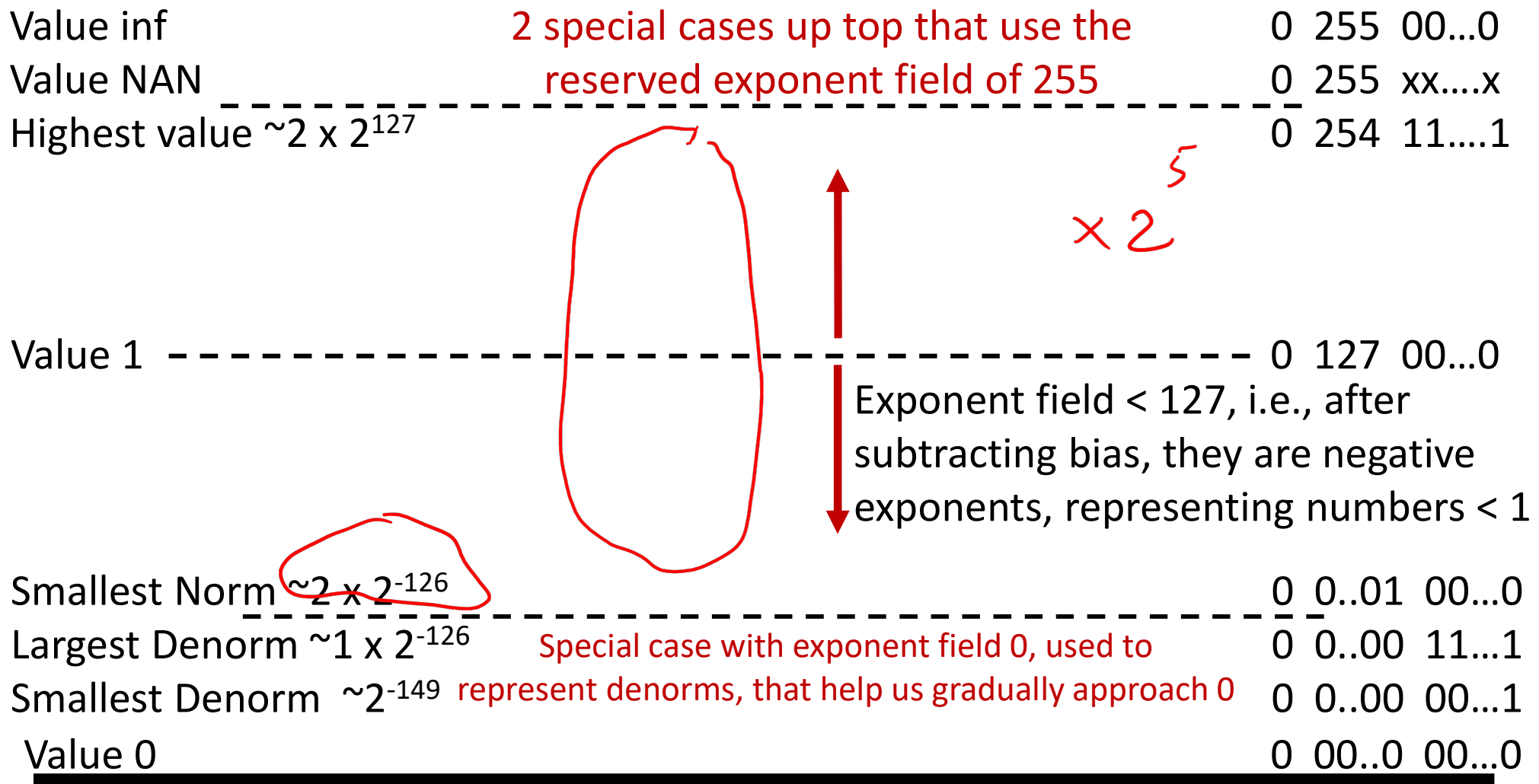


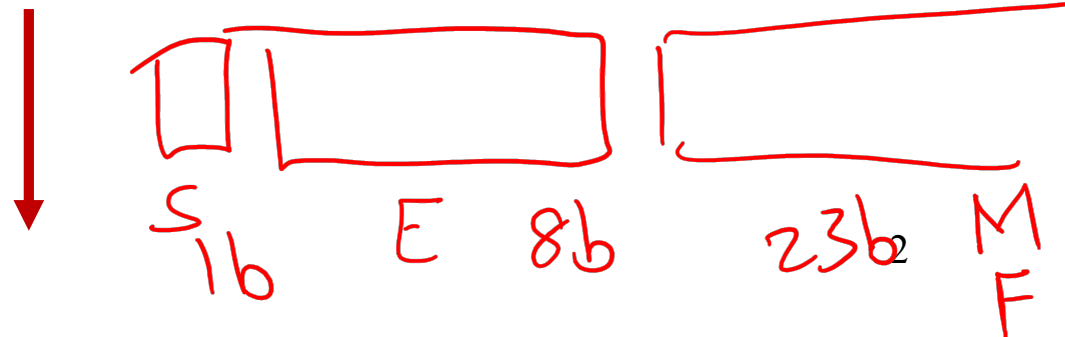
Lecture 11: Floating Point, Digital Design

- Today's topics:
 - FP formats, arithmetic
 - Intro to Boolean functions
- Exam reminders:
 - 1 sheet of notes, plus green sheet
 - No phones, simple calculators ok
 - 10:40 – 12:10, attempt every question
 - Practice exam posted tomorrow
 - Study tips
 - HW 1-4, until Lecture 8 / Slide 6



Same rules as above, but the sign bit is 1

Same magnitudes as above, but negative numbers



$$4.2321 \times 10^1 = 42.321 \quad \text{leftmost}$$

$$42 \div 2 = 21 \text{ rem } 0$$

$$21 \div 2 = 10 \text{ rem } 1$$

$$10 \div 2 = 5 \text{ rem } 0$$

$$5 \div 2 = 2 \text{ rem } 1$$

$$2 \div 2 = 1 \text{ rem } 0$$

$$1 \div 2 = 0 \text{ rem } 1$$

$$0 \div 2 = 0 \text{ rem } 0$$

~~42~~

$$321 \times 2 = 0.642$$

$$0.642 \times 2 = 1.284$$

$$0.284 \times 2 = 0.568$$

$$0.568 \times 2 = 1.136$$

$$0.136 \times 2 = 0.272$$

128 4
↑ ↑
1000 0100

Remember:

5 $\xrightarrow{+127}$ 132
True exponent Exponent in register
 $\xleftarrow{-127}$

0101010.01010.....

1.0101001010..... $\times 2^5$

0 | 1000 0100 | 0101001010.....

Sign

Examples

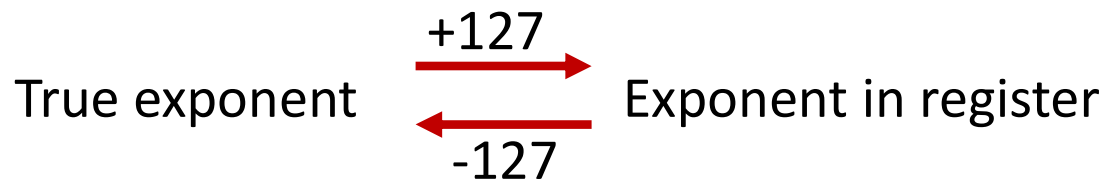
Final representation: $(-1)^S \times (1 + \text{Fraction}) \times 2^{(\text{Exponent} - \text{Bias})}$

- Represent -0.75_{ten} in single and double-precision formats

Single: $(1 + 8 + 23)$

Double: $(1 + 11 + 52)$

Remember:



- What decimal number is represented by the following single-precision number?

1 1000 0001 01000...0000

1 1000 0001

↓
negative

↓ 128

↓ +

↓ 1

= 129

Exp in Reg

↓

-127

↓

2

True exp

0100...0

-1.01 × 2²

↑
implicit
1

FP Register

= -101.00
(binary)
= -5.0

Remember:

True exponent

+127

-127

Exponent in register

1/2

1

1/4

-1.01

= -1.25 × 4 = -5

Examples

Final representation: $(-1)^S \times (1 + \text{Fraction}) \times 2^{(\text{Exponent} - \text{Bias})}$

- Represent -0.75_{ten} in single and double-precision formats

Single: $(1 + 8 + 23)$

1 0111 1110 1000...000

Double: $(1 + 11 + 52)$ *1023*

1 0111 1111 110 1000...000

- What decimal number is represented by the following single-precision number?

1 1000 0001 01000...0000

-5.0

FP Addition

- Consider the following decimal example (can maintain only 4 decimal digits and 2 exponent digits)

$$9.999 \times 10^1 + 1.610 \times 10^{-1}$$

Convert to the larger exponent:

$$9.999 \times 10^1 + 0.016 \times 10^1$$

Add

$$10.015 \times 10^1$$

Normalize

$$1.0015 \times 10^2$$

Check for overflow/underflow

Round

$$1.002 \times 10^2$$

Re-normalize

shrink

$$0.0161 \times 10^{-1} \xrightarrow{+1} 0.016 \times 10^0$$

grow

$$\begin{array}{r} 9.999 \times 10^1 \\ 0.016 \times 10^1 \\ \hline 10.015 \times 10^1 \\ 1.0015 \times 10^2 \end{array}$$

FP Addition

- Consider the following decimal example (can maintain only 4 decimal digits and 2 exponent digits)

$$9.999 \times 10^1 + 1.610 \times 10^{-1}$$

Convert to the larger exponent:

$$9.999 \times 10^1 + 0.016 \times 10^1$$

Add

$$10.015 \times 10^1$$

Normalize

$$1.0015 \times 10^2$$

Check for overflow/underflow

Round

$$1.002 \times 10^2$$

Re-normalize

$$\begin{array}{r} 1.610 \\ 9.999 \\ \hline 10.015 \end{array}$$

If we had more fraction bits,
these errors would be minimized

FP Addition – Binary Example

□ □ □ □ □

$$001.01 \times 2^1$$

- Consider the following binary example

$$0.0101 \times 2^3$$

$$1.010 \times 2^1 + 1.100 \times 2^3$$

Convert to the larger exponent:

$$0.0101 \times 2^3 + 1.1000 \times 2^3$$

Add

$$1.1101 \times 2^3$$

Normalize

$$1.1101 \times 2^3$$

Check for overflow/underflow

Round

Re-normalize

IEEE 754 format: 0 10000010 1101 100000000000000000000000

(room for 4 bits right of binary point)

$$1.1000 \times 2^3$$

$$0.0101 \times 2^3$$

$$1.1101 \times 2^3$$

130 ✓ +127

round
1.1101 → 1.111 9 ✓

FP Multiplication

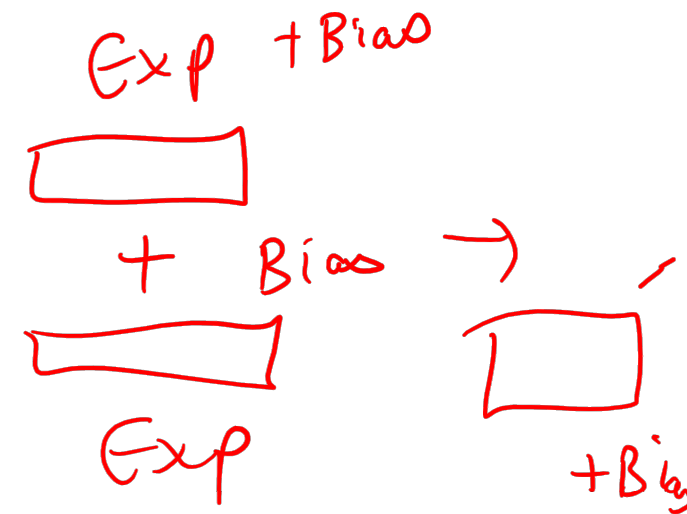
- Similar steps:

- Compute exponent (careful!)
- Multiply significands (set the binary point correctly)
- Normalize
- Round (potentially re-normalize)
- Assign sign

$$\begin{array}{r}
 1.23 \times 10^3 \\
 4.71 \times 10^4 \\
 \hline
 \end{array}$$

$$\begin{array}{r}
 \times 2^3 \\
 \times 2^7 \\
 \hline
 \times 2^{10}
 \end{array}$$

Addition (sub one of the biases)



MIPS Instructions

single-prec 32b

double prec 64b

- add.s, add.d, and similarly for sub, mul, div

add.d \$f0, \$f2, \$f4

- Comparison instructions: c.eq.s, c.neq.s, c.lt.s....

These comparisons set an internal bit in hardware that is then inspected by branch instructions: bc1t, bc1f

- Separate register file \$f0 - \$f31 : a double-precision value is stored in (say) \$f4-\$f5 and is referred to by \$f4

- Load/store instructions (lwc1, swc1) must still use integer registers for address computation

c. eg. s
bc1t \$f3, \$f6
L1

c. eg. s
add
add
sub
bc1t¹¹

Code Example

lwc1 \$f2, 8(\$gp)

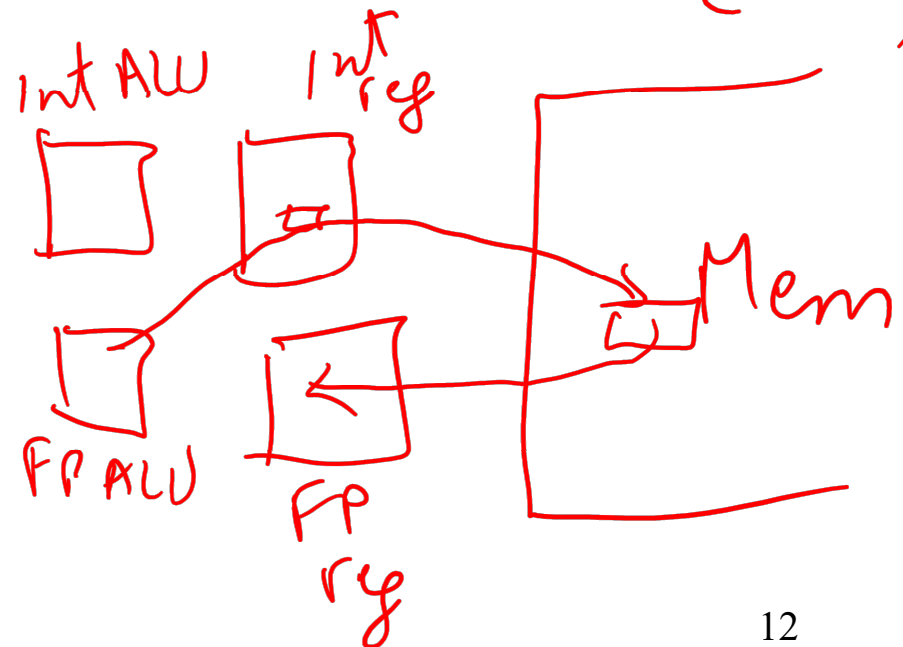
FP dest
reg

Int
src
reg to
calculate
addr
(Int)

```
float f2c (float fahr)
{
    return ((5.0/9.0) * (fahr - 32.0));
}
```

(argument fahr is stored in \$f12, return value in \$f0)

```
lwc1   $f16, const5
lwc1   $f18, const9
div.s   $f16, $f16, $f18
lwc1   $f18, const32
sub.s   $f18, $f12, $f18
mul.s   $f0, $f16, $f18
jr     $ra
```



Fixed Point

Int - Add - 1 cycle
FP - Add - 4 cycles

- FP operations are much slower than integer ops
- Fixed point arithmetic uses integers, but assumes that every number is multiplied by the same factor
- Example: with a factor of $1/1000$, the fixed-point representations for 1.46, 1.7198, and 5624 are respectively 1460, 1720, and 5624000
- More programming effort and possibly lower precision for higher performance

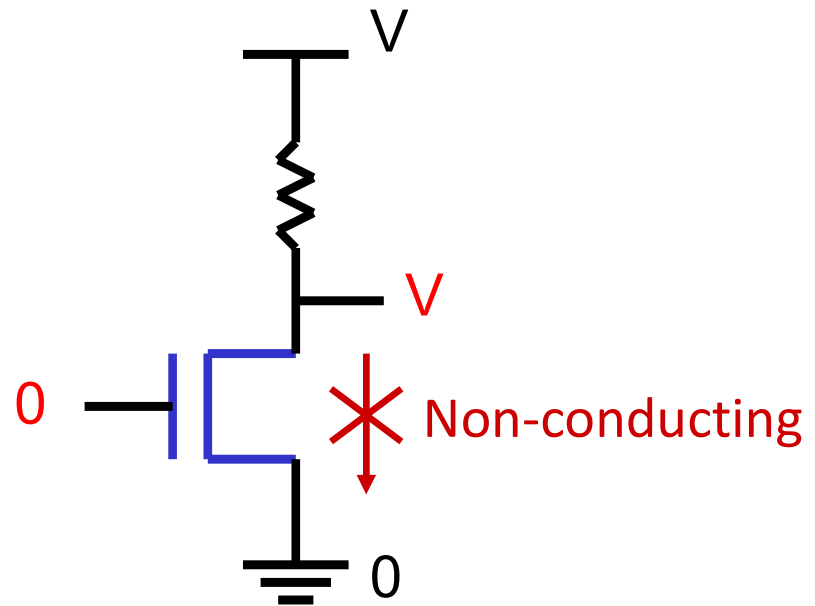
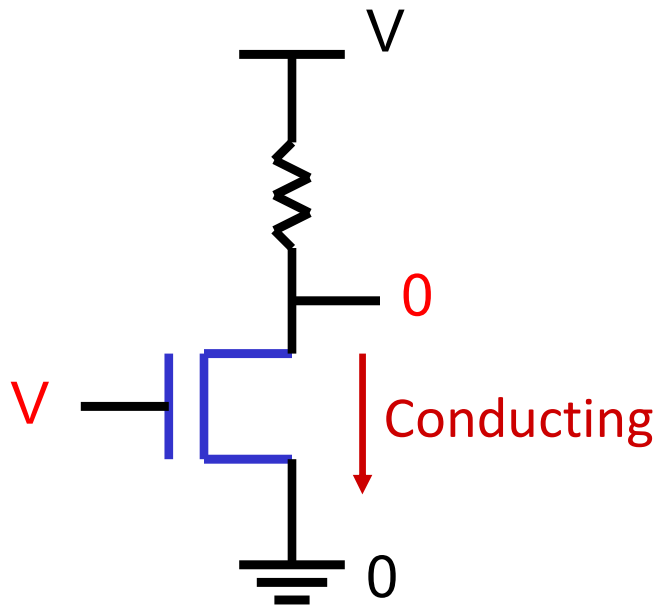
0.00239
2

Subword Parallelism

- ALUs are typically designed to perform 64-bit or 128-bit arithmetic
- Some data types are much smaller, e.g., bytes for pixel RGB values, half-words for audio samples
- Partitioning the carry-chains within the ALU can convert the 64-bit adder into 4 16-bit adders or 8 8-bit adders
- A single load can fetch multiple values, and a single add instruction can perform multiple parallel additions, referred to as subword parallelism

Digital Design Basics

- Two voltage levels – high and low (1 and 0, true and false)
Hence, the use of binary arithmetic/logic in all computers
- A transistor is a 3-terminal device that acts as a switch



Logic Blocks

- A logic block has a number of binary inputs and produces a number of binary outputs – the simplest logic block is composed of a few transistors
- A logic block is termed *combinational* if the output is only a function of the inputs
- A logic block is termed *sequential* if the block has some internal memory (state) that also influences the output
- A basic logic block is termed a *gate* (AND, OR, NOT, etc.)

We will only deal with combinational circuits today

Truth Table

- A truth table defines the outputs of a logic block for each set of inputs
- Consider a block with 3 inputs A, B, C and an output E that is true only if *exactly* 2 inputs are true

A	B	C	E

Truth Table

- A truth table defines the outputs of a logic block for each set of inputs
- Consider a block with 3 inputs A, B, C and an output E that is true only if *exactly* 2 inputs are true

A	B	C	E
0	0	0	0
0	0	1	0
0	1	0	0
0	1	1	1
1	0	0	0
1	0	1	1
1	1	0	1
1	1	1	0

Can be compressed by only representing cases that have an output of 1