

lecture 2

- fixed point
- IEEE floating point standard

Wed. January 13, 2016

For those interested in finding out what research is all about, I encourage you to participate in studies such as these.

Participants needed for Social Psychology research

Principal Investigator:
Dr. John Lydon
Contact:
Thomas Khullar at
lydonlab.coordinator@gmail.com

- If you think you might like to participate in a study and be compensated for your time, you can **enter the participant pool**.
- Researchers in McGill's Social Psychology Labs will then have access to your information and will contact you to let you know about available studies.
- Compensation is **~\$10/hr**, depending on the study.
- All studies are approved by the McGill Ethics Committee.
- There is **no obligation** to participate in a study once you are in the pool, and you are free to withdraw from the pool at any time.

To sign up, simply go to the following website and complete the 1-minute **McGill Social Psychology Recruitment Survey**:

<http://fluidsurveys.com/s/paidpool/>

Note: if you signed up last year and are interested in participating again this year, please complete the online questionnaire again as your information may have changed

Fixed point

Fixed point means we have a constant number of bits (or digits) to the left and right of the binary (or decimal) point.

Examples :

23953223.49 (base 10)

Currency uses a fixed number of digits to the right.

10.1101 (base 2)

Two's complement for fixed point numbers

e.g. 0110.1000 which is 6.5 in decimal

How do we represent -6.5 in fixed point ?

$$\begin{array}{rcl} 0110.1000 & & \\ 1001.0111 & \leftarrow & \text{invert bits} \\ + \underline{0000.0001} & \leftarrow & \text{add .0001} \\ \hline 0000.0000 & & \end{array}$$

Thus,

$$\begin{array}{rcl} 1001.0111 & \leftarrow & \text{invert bits} \\ + \underline{0000.0001} & \leftarrow & \text{add .0001} \\ \hline 1001.1000 & \leftarrow & \text{answer: -6.5 in (signed) fixed point} \end{array}$$

Scientific Notation (floating point)

$$300,000,000 = 3 \times 10^8$$
$$= 3.0 E + 8$$

$$.00000456 = 4.56 E - 6$$



"Normalized" : one digit to the left of the decimal point.

Scientific Notation in binary

$$(1000.01)_2 = 1.00001 \times 2^3$$

$$(0.111)_2 = 1.11 \times 2^{-1}$$

"Normalized" means one "1" bit to the left of the binary point. **(Note that 0 cannot be represented this way.)**

sign

"exponent"

+

|

.

"significand"

(also called
"mantissa")

x

2

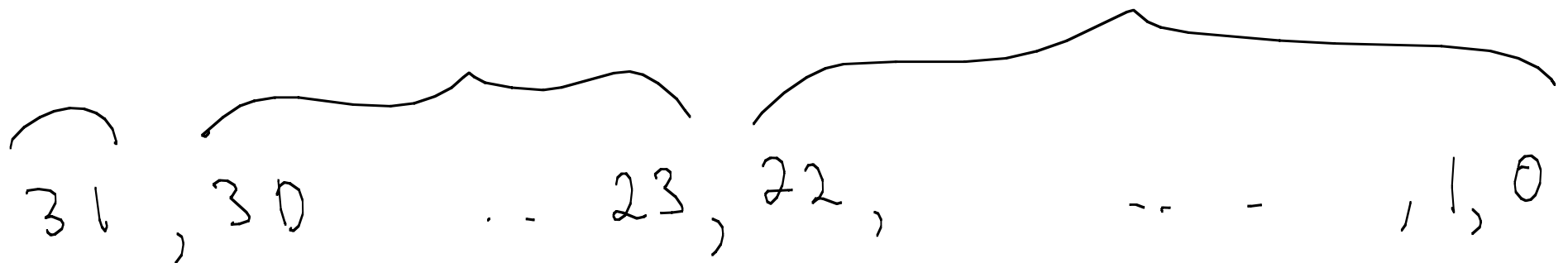
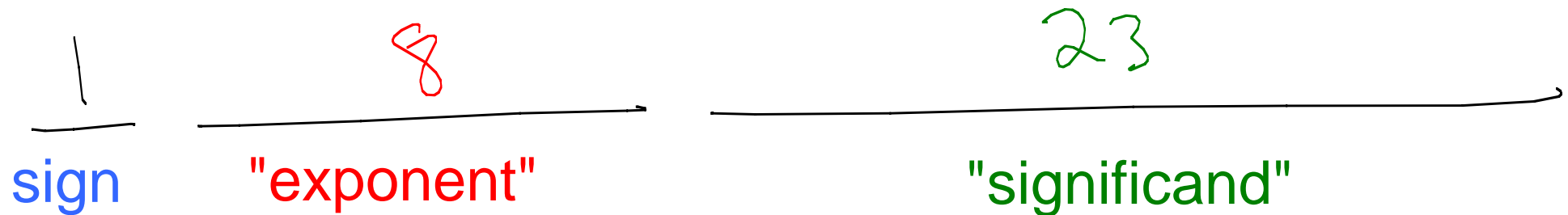
^E

How to represent this information ?

How to represent the number 0 ?

IEEE floating point standard (est. 1985)

case 1: single precision (32 bits = 4 bytes)



Let's look at these three parts, and then examples.

sign 0 for positive, 1 for negative

"significand"



You don't encode the "1" to the left of the binary point.

Only encode the first 23 bits to the right of the binary point.

exponent code

exponent value

00000000

reserved (explained soon)

00000001

-126

00000010

-125

00000011

- 124

:

:

:

:

01111111

0

This is not two's
complement !

10000000

1

10000001

2

:

:

:

:

11111110

127

11111111

reserved (explained soon)

unsigned exponent code = exponent value + "bias"
(for 8 bits, bias is defined to be 127)

Q: What is the largest positive normalized number ?
(single precision)

$$1.\text{ } \times 2^e$$

A:

$$1.\text{ } 1111 \dots 11 \times 2^{127}$$

$$2^{127} \approx 10^?$$

$$2^{10} \approx 10^3$$

$$2^{127} = 2^{120} \cdot 2^7$$

$$= (2^{10})^{12} \cdot 2^7$$

$$\approx (10^3)^{12} \cdot 10^2$$

$$= 10^{38}$$

Q: What is the smallest positive normalized number ?
(single precision)

$$1 \text{ . } \text{-----} \times 2^e$$

A:

$$1 \text{ . } 000000 \dots 0 \times 2^{-126}$$

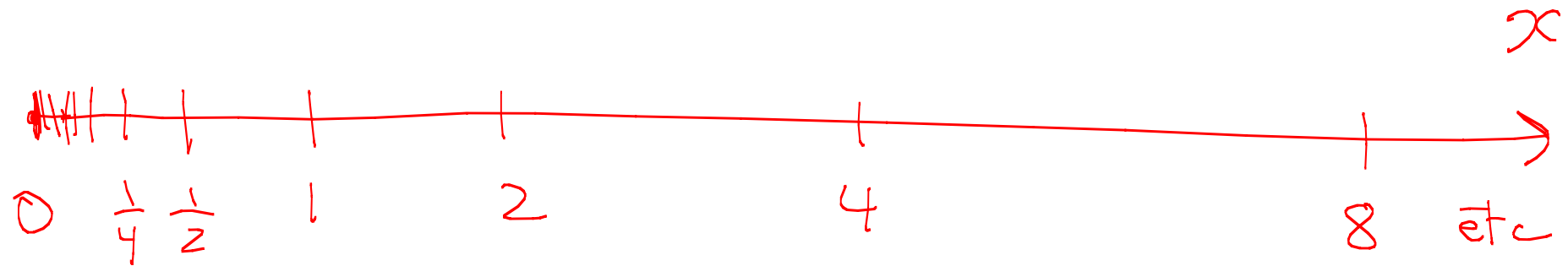
Exponent code 00000000 reserved for
"denormalized" numbers

$$1 \pm 0. \text{-----} \times 2^{-126}$$

• belong to $(-2^{-126}, 2^{-126})$

• includes 0

Dividing each power of 2 interval into 2^{23} equal parts
(same for negative real numbers).



Note the power of 2 intervals themselves are equally spaced on a log scale.



Exponent code 1111111 also reserved.

if significand is all 0's

then value is \pm infinity (depending on sign bit)

else value is NaN ("not a number")

e.g. variable is declared but hasn't been
assigned a value

This is the stuff you put on an exam crib sheet.
(Yes, you can bring a crib sheet for the quizzes.)

Example: write 8.75 a single precision float (IEEE).

First convert to binary.

$$8.75$$

$$= (1000)_2 \cdot (75)_{10}$$

$$= (10001)_2 \cdot (5)_{10} \times 2^{-1}$$

$$= 100011.0 \times 2^{-2}$$

$$= 1.00011 \times 2^3$$

$$(8.75)_{10} = (1.00011)_2 \times 2^3$$

23 bit significand: 000110000000000000000000000

exponent value: $e = 3$

exponent code = exponent value (e) + bias

Thus, exponent code is unsigned $3 + 127$.

$$(130)_{10} = (10000010)_2$$

So, the 32 bit representation is :

$\underbrace{0}_{0 \times 4} \underbrace{10000010}_1 \underbrace{0001}_{0} \underbrace{100000000000000000000000}_{c} \underbrace{0}_{0} \underbrace{0}_{0} \underbrace{0}_{0} \underbrace{0}_{0}$

Recall last lecture: 0.05 cannot be represented exactly.

```
float x = 0;  
for (int ct = 0; ct < 20; ct++) {  
    x += 1.0 / 20;  
    System.out.println( x );  
}
```

0.05

0.1

0.15

0.2

0.25

0.3

0.350000002

0.400000004

0.450000005

0.500000006

etc

Floating Point Addition

$$x = 1.00100100010000010100001 * 2^2$$

$$y = 1.101010000000000000101010 * 2^{-3}$$

$$x + y = ?$$

Floating Point Addition

$$x = 1.00100100010000010100001 \quad * \quad 2^2$$

$$y = 1.10101000000000000101010 \quad * \quad 2^{-3}$$

$$x + y = ?$$

$$x = 1.0010010001000001010000100000 \quad * \quad 2^2$$

$$y = .000011010100000000000101010 \quad * \quad 2^2$$

but the result $x+y$ has more than 23 bits of significand

How many *digits* (base 10) of precision can we represent with 23 *bits* (base 2) ?

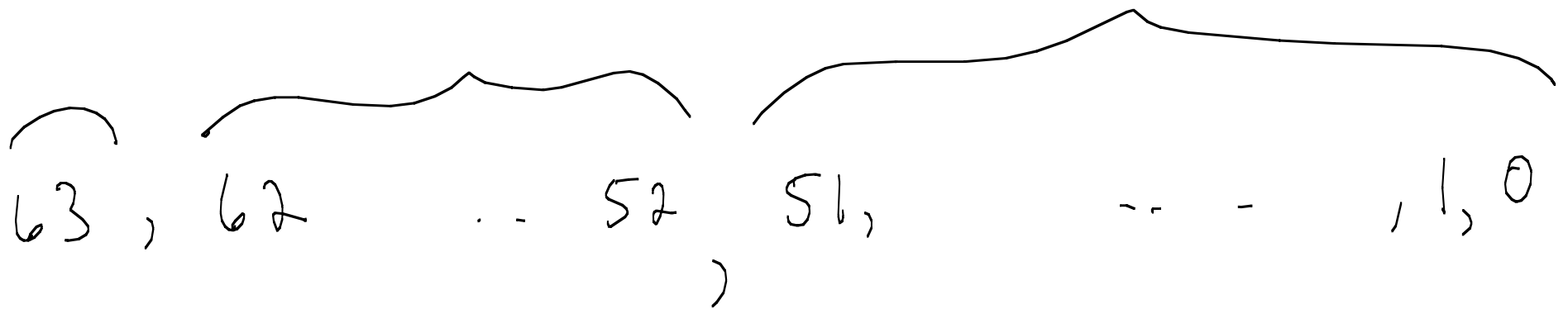
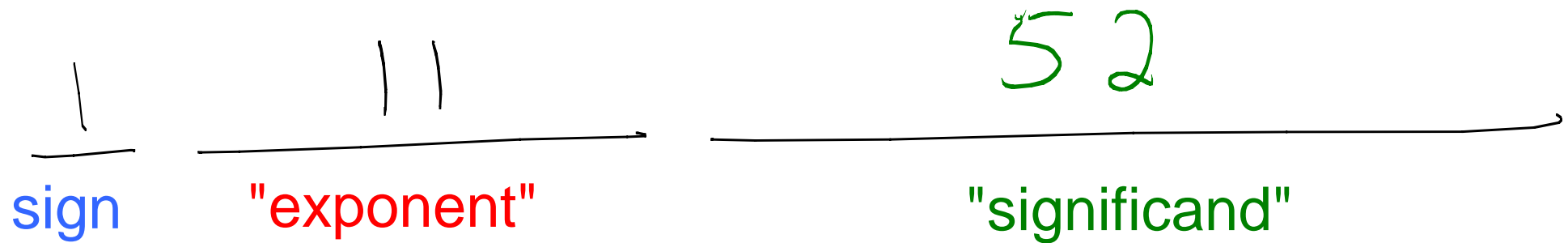
$$2^{23}$$

$$= (2^{10})^2 2^3$$

$$\approx 10^6 \cdot 10^1$$

$$= 10^7$$

case 2: double precision (64 bits = 8 bytes)



exponent code

exponent value

unsigned exponent code = exponent value + bias

For 11 bits, bias is defined to be $2^{10} - 1 = 1023$.

000000000000

000000000001

000000000010

000000000011

⋮

⋮

011111111111

100000000000

100000000001

⋮

⋮

111111111110

111111111111

reserved

-1022

-1021

- 1020

⋮

⋮

0

1

2

⋮

⋮

1023

reserved

Example

$$(8.75)_{10} = (1.00011)_2 \times 2^3$$

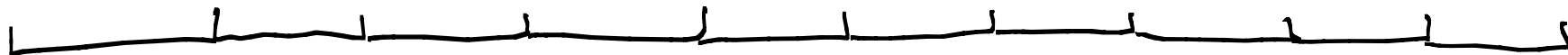
significand (52 bits)

$$= .000110000000000000000000000000000000....$$

exponent = 3, code using 11 bits:

$$3 + 1023 = 1026 = (10000000010)_2$$

double precision float (64 bits)

[illegible]

0 x 4 0 2 1 8 0 0 0 0 0 000000

Q: What is the largest positive normalized number ?
(double precision)

$$1.5999999999999999 \times 2^e$$

A:

$$2^{1023}$$

11

$$(2^{10})^{102}$$

$$2^3$$

22

$$(10^3)^{102}$$

10

11

$$10^{307}$$

Approximation Errors (Java/C/...)

```
double x = 0;  
for (int ct=0; ct < 10; ct++) {  
    x += 1.0 / 10;  
    System.out.println( x );  
}
```

0.1

0.2

0.300000000000000000000004

0.4

0.5

0.6

0.7

0.799999999999999999999999

0.899999999999999999999999

0.999999999999999999999999

How many *digits* of precision can we represent with 52 *bits* ?

$$\begin{aligned} & 2^{52} \\ = & (2^{10})^5 2^2 \\ \approx & (10^3)^5 10 \\ = & 10^{16} \end{aligned}$$

52 bits covers about the same "range" as 16 digits.

That is why the print out on the previous slide had up to (about) 16 digits to the right of the decimal point.

Announcements

- public web page (Course outline etc)
- corequisite courses:
 - COMP 206 (official)
 - COMP 250 (unofficial)

It is not recommended to do 250+206+273 together.
Rather, 250+206 only, or 206+273 only.
- assignments, there will be 4 (not 3), logisim,
each should take ~10 hours (still worth total of 30%)
- waiting list issues $(14 \times 12 + 10 = 178$ seats in room)
- quiz 1: may have to sit on stairs and use a book :/
(only 15 min)