

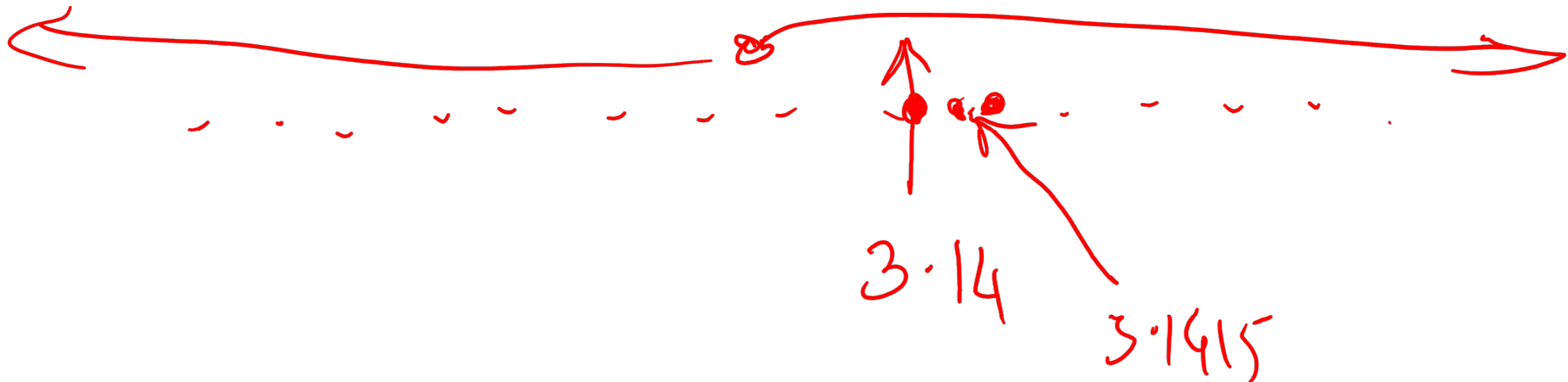
Lecture 10: Floating Point

- Today's topics:

- IEEE 754 representations



32 b
unsigned 0 — 4B
signed -2B — 2B



$$0.5 \Rightarrow 5.0 \times 10^{-1}$$

Floating Point

$$\times 0.234$$

$$\checkmark 2.345821 \times 10^{-17}$$

$$= 2 + 3 \times 10^{-1} + 4 \times 10^{-2} + 5 \times 10^{-3} \dots$$

- Normalized scientific notation: single non-zero digit to the left of the decimal (binary) point – example: 3.5×10^9

$$1.010001 \times 2^{-5}_{\text{two}} = (1 + 0 \times 2^{-1} + 1 \times 2^{-2} + \dots + 1 \times 2^{-6}) \times 2^{-5}_{\text{ten}}$$

- A standard notation enables easy exchange of data between machines and simplifies hardware algorithms – the IEEE 754 standard defines how floating point numbers are represented

$$0.5_{\text{dec}} \rightarrow 0.1 \Rightarrow 1.0 \times 2^{-1}$$

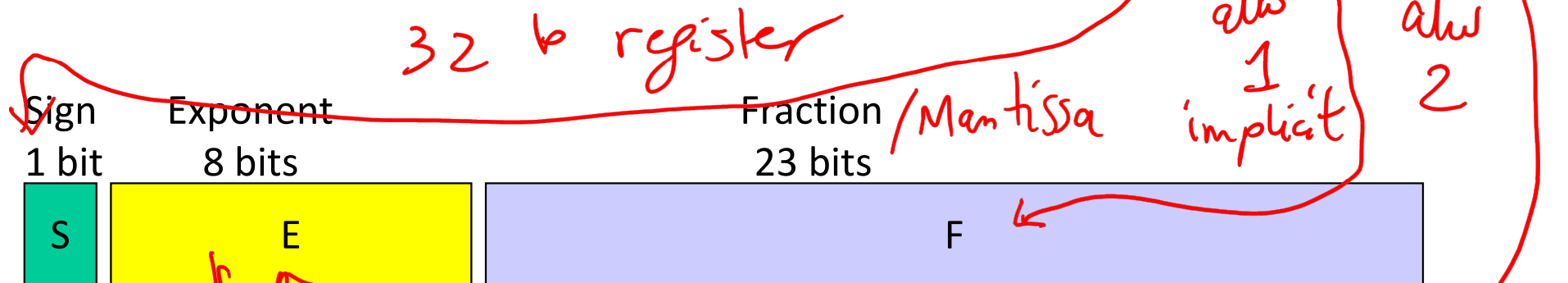
$$1.111$$

$$1 + 1 \times 2^{-1} + 1 \times 2^{-2} + 1 \times 2^{-3} = 1.875$$

$$1 + 0.5 + 0.25 + 0.125$$

decimal $+1.875 \Rightarrow$ bin 1.111×2^0

Sign and Magnitude Representation



- More exponent bits \rightarrow wider range of numbers (not necessarily more numbers — recall there are infinite real numbers)
 $-128 \leftrightarrow 128$

- More fraction bits \rightarrow higher precision

- Register value = $(-1)^S \times F \times 2^E$

- Since we are only representing normalized numbers, we are guaranteed that the number is of the form 1.xxxx..
Hence, in IEEE 754 standard, the 1 is implicit

Register value = $(-1)^S \times (1 + F) \times 2^E$

dec 0.1

0.00...

Exponent Representation

Exp
8b

- To simplify sort, sign was placed as the first bit
- For a similar reason, the representation of the exponent is also modified: in order to use integer compares, it would be preferable to have the smallest exponent as 00...0 and the largest exponent as 11...1
- This is the biased notation, where a bias is subtracted from the exponent field to yield the true exponent
- IEEE 754 single-precision uses a bias of 127 (since the exponent must have values between -127 and 128)...double precision uses a bias of 1023

0 → 255
↓
-127 +128

00010
60111
-127
-126

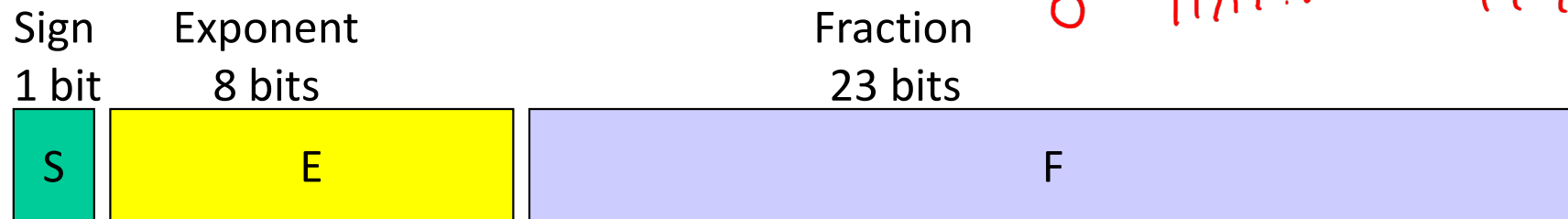
add bias

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

+128 [111111]

Sign and Magnitude Representation

Handwritten notes: $1.11111111 \times 2^{128}$ and $\frac{1}{2}$

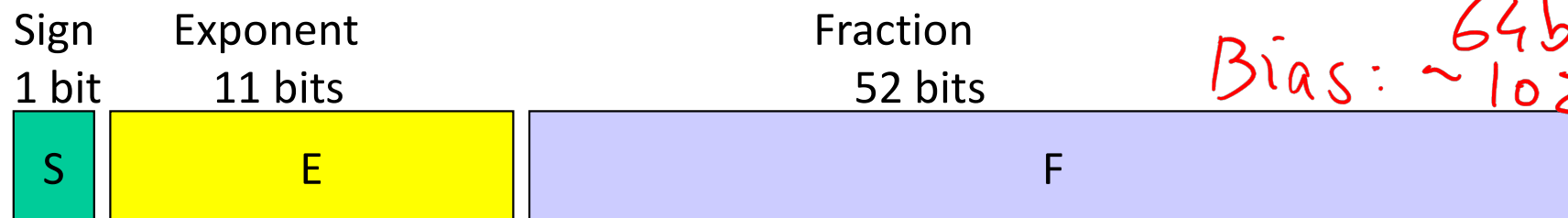


- Largest number that can be represented: $2.0 \times 2^{128} = 2.0 \times 10^{38}$ (not really – see upcoming details)
- Smallest number that can be represented: $1.0 \times 2^{-127} = 2.0 \times 10^{-38}$ (not really – see upcoming details)
- Overflow: when representing a number larger than the max;
Underflow: when representing a number smaller than the min

*Single prec
32b*

*Bias
: -127*

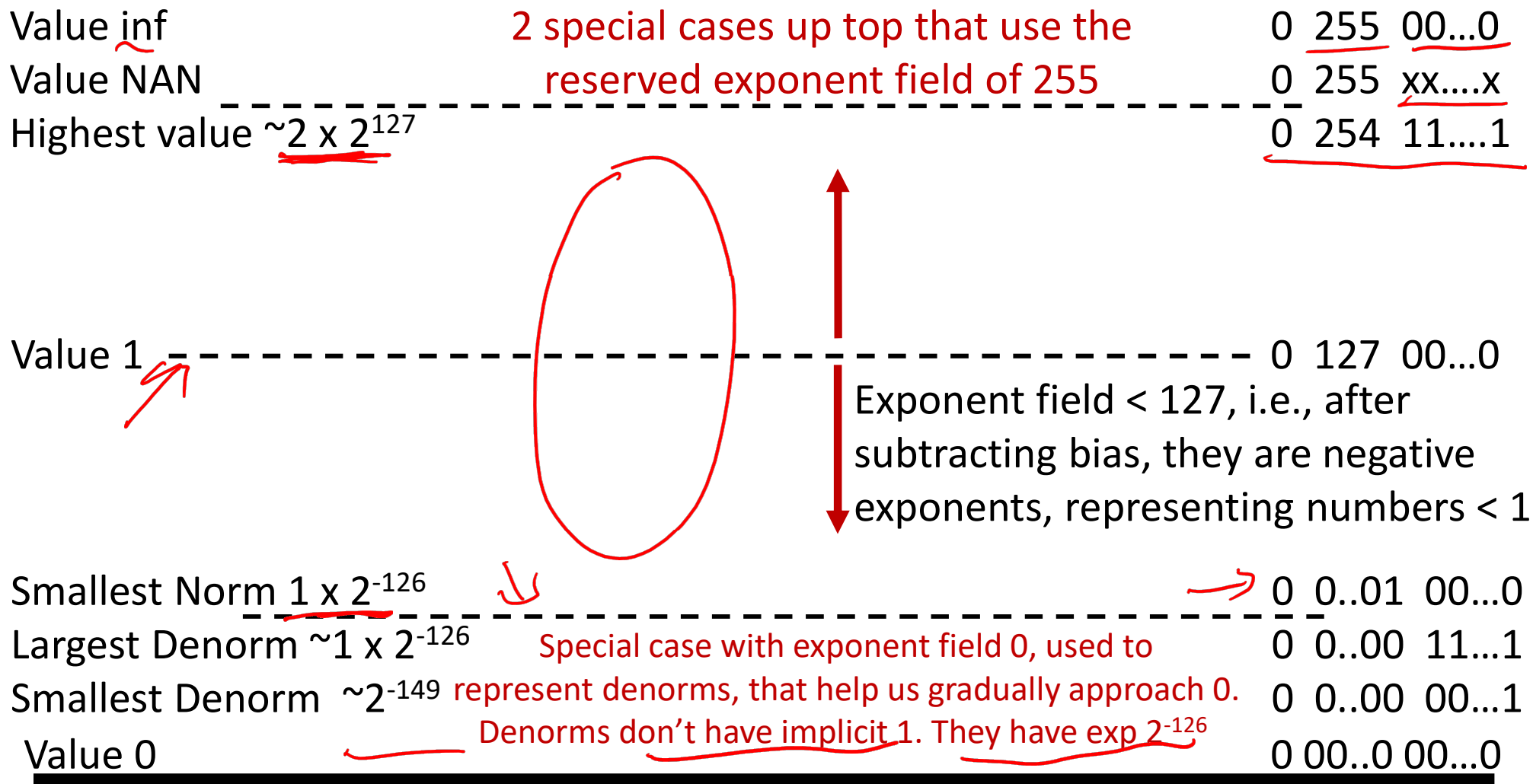
- Double precision format: occupies two 32-bit registers:
Largest: Smallest:



*Double prec
64b
Bias: ~1023*

Details

- The number “0” has a special code so that the implicit 1 does not get added: the code is all 0s
(it may seem that this takes up the representation for 1.0, but given how the exponent is represented, that’s not the case)
(see discussion of denorms in the textbook)
- The largest exponent value (with zero fraction) represents +/- infinity
- The largest exponent value (with non-zero fraction) represents NaN (not a number) – for the result of 0/0 or (infinity minus infinity)
- Note that these choices impact the smallest and largest numbers that can be represented



Same rules as above, but the sign bit is 1

Same magnitudes as above, but negative numbers

0.000...01

Examples

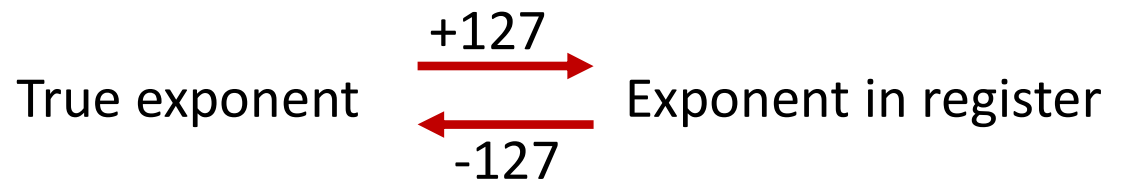
Final representation: $(-1)^S \times (1 + \text{Fraction}) \times 2^{(\text{Exponent} - \text{Bias})}$

- 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

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: $(1 + 11 + 52)$

1 0111 1111 110 1000...000

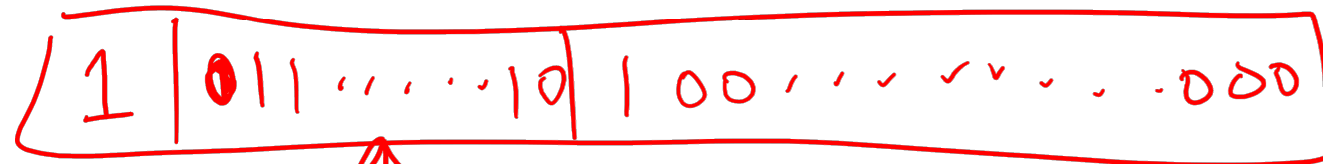
- What decimal number is represented by the following single-precision number?

1 1000 0001 01000...0000

-5.0

$$-0.75_{\text{dec}} \Rightarrow -0.11_{\text{bin}} \Rightarrow 1.1 \times 2^{-1}$$

NSN



Sign

Exp

Fraction

Remember:

+127

True exponent

Exponent in register

-127

126

+127

-1 True exp

Example 2

$$3.690625 \times 10^1$$

100100.11101000

365

NSN

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

- Represent 36.90625_{ten} in single-precision format

$$36 / 2 = 18 \text{ rem } 0$$

$$18 / 2 = 9 \text{ rem } 0$$

$$9 / 2 = 4 \text{ rem } 1$$

$$4 / 2 = 2 \text{ rem } 0$$

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

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

$$0.90625 \times 2 = 1.81250$$

$$0.8125 \times 2 = 1.6250$$

$$0.625 \times 2 = 1.250$$

$$0.25 \times 2 = 0.50$$

$$0.5 \times 2 = 1.00$$

$$0.0 \times 2 = 0.0$$

1.001001110100

$\times 2^5$

1

↓

True exp is 5

1×2^{-1}

$+ 1 \times 2^{-2}$

0.5

$+ 1 \times 2^{-7}$

132 goes in Reg

0 10000100 001 001 110100000000

36 is 100100

0.90625 is 0.1110100...0

• 11101

Example 2

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

We've calculated that $36.90625_{\text{ten}} = 100100.1110100...0$ in binary
Normalized form = $1.001001110100...0 \times 2^5$
(had to shift 5 places to get only one bit left of the point)

The sign bit is 0 (positive number)

The fraction field is 001001110100...0 (the 23 bits after the point)

The exponent field is $5 + 127$ (have to add the bias) = 132,
which in binary is 10000100

The IEEE 754 format is 0 10000100 001001110100....0
sign exponent 23 fraction bits

Remember:

True exponent

Exponent in register