

Floating Point Representation

CS520

Dept. of Computer Science
Univ. of New Hampshire

real numbers often are represented
with scientific notation

$$+ \underline{1.1066} \times 10^{\underline{22}}$$

↑ normalized form
one non-zero digit to the left of
the decimal point

$$11.066 \times 10^{21}$$

$$0.11066 \times 10^{23}$$

no reason we can't show binary numbers
in scientific notation

$$1.00 \times 2^{-1} = \frac{1}{2} = 0.5_{10}$$

$$\underline{1.11} \times 2^0 = 1 + \frac{1}{2} + \frac{1}{4} = 1.75_{10}$$

scientific notation components

sign

significant

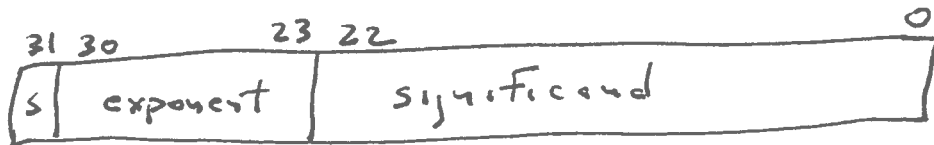
exponent - ~~possibly~~ signed

floating point number is computer support
for real numbers that uses binary numbers
in scientific notation

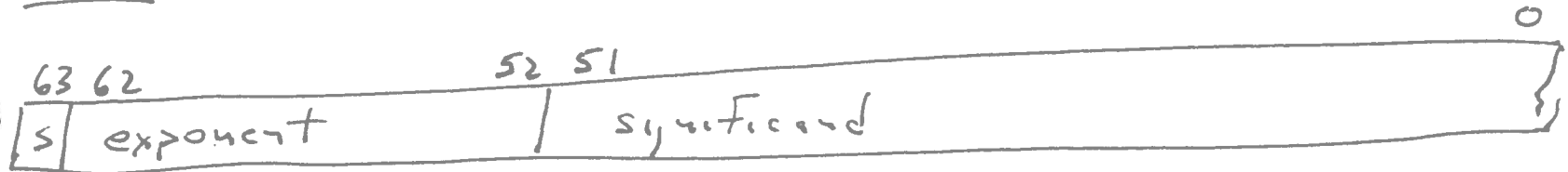
IEEE Floating Point

Standard established in 1985

single-precision (32 bits)



double-precision (64 bits)



sign bit



IEEE F.p. has normalized form

$$1.XXX\dots X \times 2^{\text{exp}}$$

↑
└ trick: since this one must be present,
don't represent it.

But how do we represent 0.0?

IEEE f.p. exponent

if exponent is all zero bits:

a) significand is all zero bits
then value is ± 0.0

b) significand is not all zero bits
then value is denormalized

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

↳ i.e. the significand bits

if exponent is all one bits

a) significand is all zero bits

then value is \pm Infinity

i.e. divide by zero
overflow

b) significand is not all zero bits

then value is NaN

↳ not a number

i.e. zero divided by zero

$\sqrt{-p}$

otherwise exponent is a biased exponent

actual exponent = stored exponent
minus bias amount

single precision: bias is 127

double precision: bias is 1023

ie treat exponent bits as unsigned integer
and subtract bias amount to
get the actual exponent

example $-0.75_{10} \rightarrow$ IEEE single-precision

$$-\left(\frac{1}{2} + \frac{1}{4}\right) = -.11 \times 2^0$$

$$= -1.10 \times 2^{-1}$$

↳ not stored!

$$\text{actual_exp} = \text{stored_exp} - 127$$

$$-1 = \text{stored_exp} - 127$$

Little Endian

$$\text{stored_exp} = +126 = \underline{01111110}_2$$

Big Endian

↓
00
00
40
BF

1011 | 1111 | 0100 | 0000 | 0000 | 0000 | 0000 | 0000
B F 4 0 0 0 0 0

BF
40
00
00

example IEEE F.p. value consisting of these bits

$20A00000_{16}$

1100 0000 1010 0000 0000 0000 0000 0000

negative

stored exp = 129
exp - 127 bias

actual 2
exp

1.0100 ————— 0 x 2²

101.0 ————— 0 x 2⁰

-5.0_{10}

example $-5.0_{10} \rightarrow$ IEEE double-precision

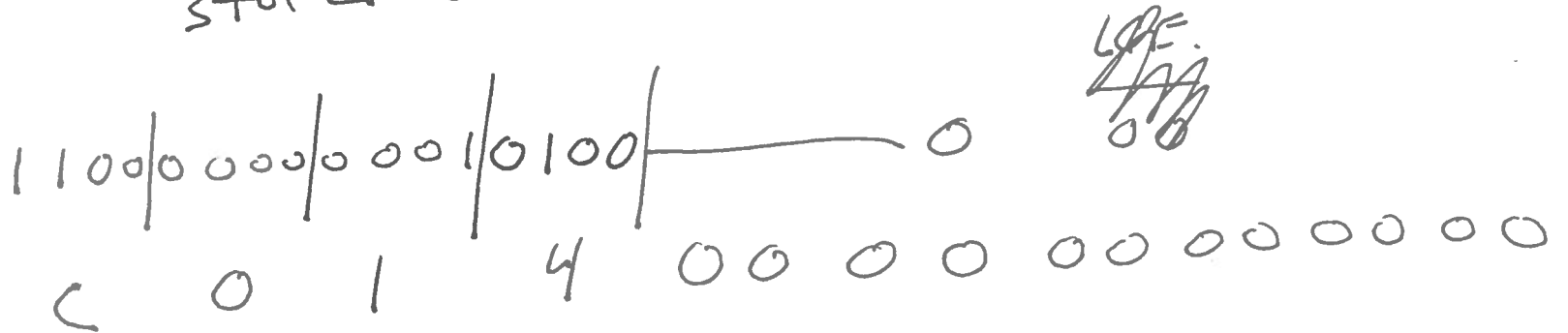
$s_{1/2} = 1$ ~~101.00~~ $\times 2^0$
~~1.0100~~ $\times 2^2$

<u>L.E.</u>	<u>B.E.</u>
00	00
00	14
00	00
00	00
00	00
00	00
00	00
14	00
00	00

actual exp = stored exp - 1023

2 = stored - 1023

stored = 1025 = 10000000000001



overflow

like with integers, value can get too big to fit in a fixed-size f.p. container

question: what component of a f.p. number actually overflows?

exponent!

Underflow

Floating-point value can get too close to zero to be seen

question: what happens to the exponent in this case?

it gets too negative