# Converting ints to floats

CS520

Dept. of Computer Science
Univ. of New Hampshire

# i2f example   0x3456789A → Float

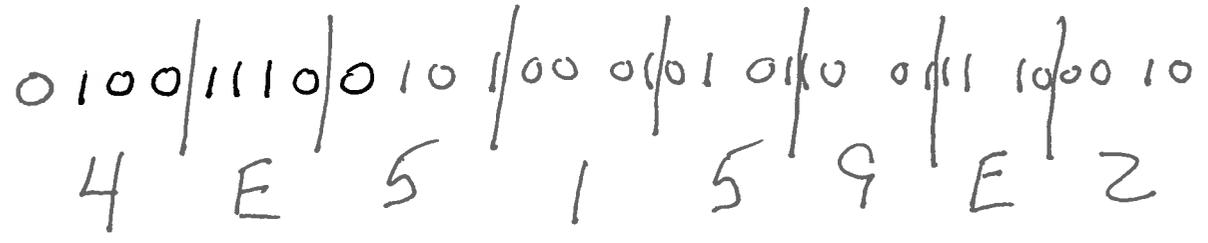$\textcircled{0}$011 0100 0101 0110 0111 1000 1001 1010. $\times 2^0$

$\underline{1.10100\ 0101\ 0110\ 0111\ 1000\ 10|01}\ \dfrac{1010}{2} \times 2^{29}$

in this case truncate extra bits    guard bits

summarized as either

$0$ or $1$

$\underset{\text{otherwise}}{\llcorner}$

all then bits are $0$

$\begin{aligned} \text{actual} \\ \text{exp} \end{aligned} = \dfrac{\text{stored}}{\text{exp}} - 127$

$29 = \dfrac{\text{stored}}{\text{exp}} - 127$

$\dfrac{\text{stored}}{\text{exp}} = 156_{10} = 10011100_2$

0100|1110|0 10 1|00 0101 0110 0111 1000 10

4    E    5    1    5    9    E    2

# IEEE Floating-point rounding

round to even

$$1. \cdots \cdots \mid GG \bigcirc$$

$$\begin{array}{c} 1 \cdots \cdots 23 \\ \uparrow \\ \text{add one} \end{array} \quad \begin{array}{c} 10 \end{array}$$

| guard bits* | sticky bit** | action |
|---|---|---|
| 00 | 0 or 1 | truncate |
| 01 | 0 or 1 | truncate |
| 10 | 0 | round to even |
| 10 | (1) | add one |
| 11 | 0 or 1 | add one |

\* two highest bits to be discarded

\*\* summarizes all other discarded bits:
  if all 0's then sticky bit is 0.
  else sticky bit is 1

truncate — just discard the bits

round to even — if low bit in significand
is 1, add 1
else do nothing

add one — add one to the significand

note: could be carry out the
top requiring re-normalization

$1.1 \longrightarrow 1$

$+1$

$10.0 \longrightarrow 0$

shift → adjust
right & the exp.

## work these cases yourself

$0x345678A0 \rightarrow 4E5159E2$

guard bits : 1 0  } round to even, but

sticky bit : 0     low bit is $\emptyset$ (i.e.

already even) so

nothing done

$0x\ 345678\ E\emptyset \rightarrow 4E5159\ E4$

guard bits: $10$ ⎫ round to even,
sticky bit: $0$ ⎭ low bit is $1$ (i.e. odd), so add $1$ to that position

$0x345678A1 \rightarrow 4E5159E3$

guard bits: 10 ⎫ add one to low bit
sticky bit: 1 ⎭ position

0x 345678B0 $\longrightarrow$ 4E5159E3

guard bits : 11 $\Big\}$ add one to low

sticky bit : ∅ bit position