## Example: Converting to IEEE 754 Form

Put 0.085 in single-precision format

1. The first step is to look at the sign of the number. Because 0.085 is positive, the sign bit =0.

 $(-1)^0 = 1.$ 

Write 0.085 in base-2 scientific notation.
 This means that we must factor it into a number in the range [1 <= n < 2] and a power of 2.</li>

```
0.085 = (-1)^{0} * (1 + fraction) * 2^{power}, or:
0.085 / 2^{power} = (1 + fraction).
```

So we can divide 0.085 by some power of 2 to get the (1 + fraction).

 $0.085 / 2^{-1} = 0.17$  $0.085 / 2^{-2} = 0.34$  $0.085 / 2^{-3} = 0.68$  $0.085 / 2^{-4} = 1.36$ 

Therefore,  $0.085 = 1.36 \times 2^{-4}$ 

## 3. Find the exponent.

The power of 2 is -4, and the bias for the single-precision format is 127. This means that the exponent =  $-4 + 127 = 123_{ten}$ . In binary, it's 01111011<sub>bin</sub>

## 4. Write the fraction in binary form

We know that the number we calculated in step 2 is in the range  $[1 \le n \le 2]$ . Therefore, we don't have to store the leading 1.

The fraction = 0.36. Unfortunately, this is not a "pretty" number, like those shown in the book. The best we can do is to approximate this value. Single-precision format allows 23 bits for the fraction.

Binary fractions look like this:

 $\begin{array}{l} 0.1 = (1/2) = 2^{-1} \\ 0.01 = (1/4) = 2^{-2} \\ 0.001 = (1/8) = 2^{-3} \end{array}$ 

To approximate 0.36, we can say:

 $\begin{array}{l} 0.36 = (0/2) + (1/4) + (0/8) + (1/16) + (1/32) + \ldots \\ 0.36 = 2^{-2} + 2^{-4} + 2^{-5} + \ldots \end{array}$ 

0.36<sub>ten</sub> ~ 0.01011100001010001111011<sub>bin</sub>.

The binary string we need is: 01011100001010001111011.

It's important to notice that you will not get 0.36 exactly. This is why floating-point numbers have error when you put them in IEEE 754 format.

5. Now put the binary strings in the correct order -1 bit for the sign, followed by 8 for the exponent, and 23 for the fraction. The answer is:

Sign Exponent Fraction Decimal 0 123 0.36 Binary 0 01111011 01011100001010001111011

Example: IEEE 754 to Float ->