

# 1 Floating Point

The IEEE 754 standard defines a binary representation for floating point values using three fields.

- The *sign* determines the sign of the number (0 for positive, 1 for negative).
- The *exponent* is in **biased notation**. For instance, the bias is 127 ( $2^{8-1} - 1$ ) for single-precision floating point numbers.
- The *significand* or *mantissa* is akin to unsigned integers, but used to store a fraction instead of an integer.

The below table shows the bit breakdown for the single precision (32-bit) representation. The leftmost bit is the MSB and the rightmost bit is the LSB.

1	8	23
Sign	Exponent	Mantissa/Significand/Fraction

For normalized floats:

$$\text{Value} = (-1)^{\text{Sign}} * 2^{\text{Exp}-\text{Bias}} * 1.\text{significand}_2$$

significand is bit  $\cdot 2^{-(n+1)}$  (for bit in  $n^{\text{th}}$  position) (from left to right)

For denormalized floats:

$$\text{Value} = (-1)^{\text{Sign}} * 2^{\text{Exp}-\text{Bias}+1} * 0.\text{significand}_2$$

Exponent	Significand	Meaning
0	Anything	Denorm
1-254	Anything	Normal
255	0	Infinity
255	Nonzero	NaN

Note that in the above table, our exponent has values from 0 to 255. When translating between binary and decimal floating point values, we must remember that there is a bias for the exponent.

1.1 How many zeroes can be represented using a float?

$$2 \text{ } +0, -0$$

1.2 What is the largest finite positive value that can be stored using a single precision float?

sign = 0    mantissa = all bits set  
Exponent = 254

$$0x7F7FFFFF = (1 + (1 - 2^{-23})) * 2^{127}$$

The mantissa for the largest value will be 23 1's. This corresponds to a value of

$$.11\dots1 = 2^{-1} + 2^{-2} + \dots + 2^{-23} = 2^{-23}(2^{22} + 2^{21} + \dots + 1)$$

Here, we apply the formula that  $\sum_{i=0}^{n-1} 2^i = 2^n - 1$ , so we have that the mantissa is

$$2^{-23}(2^{22} + 2^{21} + \dots + 1) = 2^{-23}(2^{23} - 1) = 1 - 2^{-23}$$

We have  $1 + (1 - 2^{-23})$  since we this is a normalized number and thus has a 1 to the left of the decimal point.

1.3 What is the smallest positive value that can be stored using a single precision float?

$0x00000001 = 2^{-23} * 2^{-126}$  *sign=0 Exponent=0 mantissa LSB set*

1.4 What is the smallest positive normalized value that can be stored using a single precision float?

$0x00800000 = 2^{-126}$  *sign=0 Exponent=1 mantissa=0*

1.5 Cover the following single-precision floating point numbers from binary to decimal or from decimal to binary. You may leave your answer as an expression.

•  $0x00000000$  *8.25 get mantissa  
0.25 x 2 = 0.5 0  
0.5 x 2 = 1.0 1  
0 x 2 = 0 0  
0*

•  $39.5625$  *39 = 100111.10010...*  
 $0x421E4000$  *1.001111001.2^5  
34569*

*Trick for finding mantissa for 39.5625:*  
 $0.5625 \times 2 = 1.125$  1  
 $0.125 \times 2 = 0.25$  0  
 $0.25 \times 2 = 0.5$  0  
 $0.5 \times 2 = 1.0$  1  
 $0 \times 2 = 0$  1

•  $8.25 = 1000.01_2$  *0.01<sub>2</sub> = 0.25<sub>10</sub> ✓*  
 $0x41040000$  *2<sup>3</sup> \* 1.00001  
2<sup>3</sup> \* 2<sup>-5</sup>*

•  $0xFF94BEEF$  *NaN Exponent is all ones mantissa is non zero*

•  $0x00000F00$  *(2<sup>-12</sup> + 2<sup>-13</sup> + 2<sup>-14</sup> + 2<sup>-15</sup>) \* 2<sup>-126</sup>*

•  $-\infty$  *Exponent is all ones mantissa is zero. sign is on, so negative*  
 $0xFF800000$

*Exp - Bias = 3  
Exp = 3 + 127  
Exp = 130*

## 2 More Floating Point Representation

Not every number can be represented perfectly using floating point. For example,  $\frac{1}{3}$  can only be approximated and thus must be rounded in any attempt to represent it. For this question, we will only look at positive numbers.

2.1 What is the next smallest number larger than 2 that can be represented completely?

For this question, you increment the number by the smallest amount possible. This is the same as incrementing the significand by 1 at the rightmost location.

$(1 + 2^{-23}) * 2 = 2 + 2^{-22}$

2.2 What is the next smallest number larger than 4 that can be represented completely?

For this question, you increment the number by the smallest amount possible. This is the same as incrementing the significand by 1 at the rightmost location.

$(1 + 2^{-23}) * 4 = 4 + 2^{-21}$

2.3 Define stepsize to be the distance between some value x and the smallest value larger than x that can be completely represented. What is the step size for 2? 4?

This would be the amount added in part 1. This gives  $2^{-22}$  and  $2^{-21}$ .

- 2.4 Now let's see if we can generalize the stepsize for normalized numbers (we can do so for denorms as well, but we won't in this question). If we are given a normalized number that is not the largest representable normalized number with exponent value  $x$  and with significand value  $y$ , what is the stepsize at that value? Hint: There are 23 significand bits.

Here we need to generalize the solution we got in 1 and 2. However, this is the same approach just increment the significand by the 1.

$$\begin{aligned} \text{curr\_number} &= 2^{x-127} + 2^{x-127} * y \\ \text{next\_number} &= 2^{x-127} + 2^{x-127} * y + 2^{x-127} * 2^{-23} \\ \text{stepsize} &= \text{next\_number} - \text{curr\_number} = 2^{x-150} \end{aligned}$$

- 2.5 Now let's apply this technique. What is the largest odd number that we can represent? Part 4 should be very useful in finding this answer.

To find the largest odd number we can represent, we want to find when odd numbers will stop appearing. This will be with step size of 2.

As a result, plugging into Part 4:  $2 = 2^{x-150} \rightarrow x = 151$

This means the number before  $2^{151-127}$  was a distance of 1 (it is the first value whose stepsize is 2) and no number after will be odd. Thus, the odd number is simply subtracting the previous step size of 1. This gives,  $2^{24} - 1$

### 3 RISC-V: A Rundown

RISC-V is an assembly language, which is comprised of simple instructions that each do a single task such as addition or storing a chunk of data to memory.

For example, on the left is a line of C code and on the right is a chunk of RISC-V code that accomplishes the same thing.

<pre>int x = 5, y[2]; y[0] = x; y[1] = x * x;</pre>	<pre>// x -&gt; s0, &amp;y -&gt; s1 addi s0, x0, 5 sw    s0, 0(s1) mul   t0, s0, s0 sw    t0, 4(s1)</pre>	<p><math>s0 = x0 + 5</math>  <math>s1[0] = s0</math>  <math>t0 = s0 * s0</math>  <math>s1[1] = t0</math></p>
---	---	--

- 3.1 Can you figure out what each line in the RISC-V code is doing?

`addi s0, x0, 5` evaluates to `x = 5`. `sw s0, 0(s1)` evaluates to `y[0] = x`. `mul t0, s0, s0` calculates `x * x`. `sw t0, 4(s1)` evaluates to `y[1] = x * x`.

## 4 Registers

In RISC-V, we have two methods of storing data: main memory and registers. Registers are much faster than using main memory, but are very limited in space (32 bits each). Note that you should ALWAYS use the named registers (e.g. `s0` rather than `x8`).

Register(s)	Alt.	Description
x0	zero	The zero register, always zero
x1	ra	The return address register, stores where functions should return
x2	sp	The stack pointer, where the stack ends
x5-x7, x28-x31	t0-t6	The temporary registers
x8-x9, x18-x27	s0-s11	The saved registers
x10-x17	a0-a7	The argument registers, a0-a1 are also return value

4.1 Can you convert each instruction's registers to the other form?

```
add s0, zero, a1    -->    add x8, x0, x11
or  x18, x1, x30    -->    or  s2, ra, t5
```

Look up in the table

Extra practice:

$lw\ sp, 0(a0) \Rightarrow lw\ x2, 0(a0)$

Note that you should ALWAYS use the named registers (e.g. `s0` rather than `x8`).

## 5 Basic Instructions

For your reference, here are some of the basic instructions for arithmetic operations and dealing with memory (Note: ARG1 is argument register 1, ARG2 is argument register 2, and DR is destination register):

[inst]	[destination register] [argument register 1] [argument register 2]
add	Adds the two argument registers and stores in destination register
xor	Exclusive or's the two argument registers and stores in destination register
mul	Multiplies the two argument registers and stores in destination register
sll	Logical left shifts ARG1 by ARG2 and stores in DR
srl	Logical right shifts ARG1 by ARG2 and stores in DR
sra	Arithmetic right shifts ARG1 by ARG2 and stores in DR
slt/u	If ARG1 < ARG2, stores 1 in DR, otherwise stores 0, u does unsigned comparison
[inst]	[register] [offset]([register containing base address])
sw	Stores the contents of the register to the address+offset in memory
lw	Takes the contents of address+offset in memory and stores in the register
[inst]	[argument register 1] [argument register 2] [label]
beq	If ARG1 == ARG2, moves to label
bne	If ARG1 != ARG2, moves to label

[inst]	[destination register] [label]
jal	Stores the next instruction's address into DR and moves to label

You may also see that there is an “i” at the end of certain instructions, such as addi, slli, etc. This means that ARG2 becomes an “immediate” or an integer instead of using a register. There are also immediates in some other instructions such as sw and lw. NOTE: The size of an immediate in any given instruction depends on what type of instruction it is (more on this soon!).

5.1 Assume we have an array in memory that contains `int* arr = {1,2,3,4,5,6,0}`. Let register `s0` hold the address of the zeroth element in arr. You may assume integers are four-bytes and our values are word-aligned. What do the snippets of RISC-V code do? Assume that all the instructions are run one after the other in the same context.

```

a) lw  t0, 12(s0)      -->      Sets t0 equal to arr[3]

b) slli t1, t0, 2  ←  t1 = t0 << 2 = t0 * 4
   add  t2, s0, t1  ←  t2 = s0 + t1
   lw   t3, 0(t2)  ←  t3 = t2[0] -->      Increments arr[t0] by 1
   addi t3, t3, 1  t3 += 1
   sw   t3, 0(t2)  t2[0] = t3

c) lw  t0, 0(s0)  t0 = s0[0]
   xori t0, t0, 0xFFF ← -->      Sets t0 to -1 * arr[0]
   addi t0, t0, 1  t0 += 1

```

*t0 = t0 ^ 0xFFFFFFFF gets sign extended  
12<sup>th</sup> bit is 1 so all ones  
thus xor is really  
with 0xFFFFFFFF*

## 6 C to RISC-V

## 6.1 Translate between the C and RISC-V verbatim

C	RISC-V
<pre>// s0 -&gt; a, s1 -&gt; b // s2 -&gt; c, s3 -&gt; z int a = 4, b = 5, c = 6, z; z = a + b + c + 10;</pre>	<pre>addi s0, x0, 4  a = 4 addi s1, x0, 5  b = 5 addi s2, x0, 6  c = 6 add  s3, s0, s1  z = a + b add  s3, s3, s2  z = z + c addi s3, s3, 10  z = z + 10</pre>
<pre>// s0 -&gt; int * p = intArr; // s1 -&gt; a; * p = 0; int a = 2; p[1] = p[a] = a;</pre> <p><i>Handwritten notes:</i>    put 0 to addr @ p    P[1]=a    P[a]=a</p>	<pre>sw  x0, 0(s0)  Equiv of below -&gt; *(p+a) addi s1, x0, 2 sw  s1, 4(s0)  increments it so a is now                 word addr + not byte addr. slli t0, s1, 2  ← adds the new addr add  t0, t0, s0 sw  s1, 0(t0)  ← normal set w/new addr.</pre>
<pre>// s0 -&gt; a, s1 -&gt; b int a = 5, b = 10; if(a + a == b) {     a = 0; } else {     b = a - 1; }</pre> <p><i>Handwritten notes:</i>    a = 5    b = 10    t0 = a + a    branch: if(t0 == b)    a = 0    b = a - 1</p>	<pre>addi s0, x0, 5 addi s1, x0, 10 add  t0, s0, s0 bne t0, s1, else xor  s0, x0, x0  another method of                   setting a reg to zero jal  x0, exit else: addi s1, s0, -1 exit:</pre>
<pre>// computes s1 = 2^30 s1 = 1; for(s0=0; s0&lt;30; s++) {     s1 *= 2; }</pre>	<pre>addi s0, x0, 0  s0 = 0 addi s1, x0, 1  s1 = 1 addi t0, x0, 30 t0 = 30 ← how many loops loop:     beq s0, t0, exit  while(s0 != 0)     add s1, s1, s1    s1 = s1 + s1 ← s1 * 2     addi s0, s0, 1    s0 = s0 + 1 ← inc counter     jal x0, loop exit:</pre>

```

sum=0
while (n>0){
  sum=sum+1
  n=n-1
}

```

```

// s0 -> n, s1 -> sum
// assume n > 0 to start
for(int sum = 0; n > 0; n--) {
  sum += n;
}

```

```

addi s1, x0, 0 ← s1=0
loop:
  beq s0, x0, exit  while (s0!=0){
  add s1, s1, s0    s1=s1+s0
  add s0, s0, -1    s0=s0+-1
  jal x0, loop      }
exit:

```