## CS330 Architecture and Organization Assignment Chapter 3

Your Name Here

**Problem 1** — (8 points) Convert the following binary values to decimal.

- (a) 0.101
- (b) 110.101
- (c) 1 0000 0000.0000 0000 1
- (d) 111 1111.1111 111

## Answer:

- (a)  $2^{-1} + 2^{-3} = 0.5 + 0.125 = 0.625$
- (b)  $2^2 + 2^1 + 2^{-1} + 2^{-3} = 4 + 2 + 0.5 + 0.125 = 6.625$
- (c)  $2^8 + 2^{-9} = 256.001953125$
- (d) This can be done by hand, but a clever trick is to realize that the ones left of the binary point are one less than a power of 2. And the ones to the right of the binary point are the same number but shifted.

$$(2^7 - 1) + (2^7 - 1) \cdot 2^{-7} = (2^7 - 1) + (1 - 2^{-7}) = 2^7 - 2^{-7} = 127.9921875$$

**Problem 2** — (8 points) Convert the following decimal values to binary.

- (a) 0.75
- (b) 16.0625
- (c) 0.1
- (d) 12.12

## Answer:

- (a) In base ten,  $75/100 = 3/4 = 2^{-1} + 2^{-2}$ , so in binary 0.11.
- (b) In base ten,  $16 + 625/10000 = 16 + 1/16 = 2^4 + 2^{-4}$ , so in binary 1 0000.0001.
- (c) This one is best solved by the standard algorithm, multiply by two until you get a number larger than one. Then subtract off the one and keep going.

So, the binary value is 0.000110011...

(d) The whole part is 12 which is easily seen to be 8 + 4 or in binary 1100. The fraction part is best computed via the algorithm.

```
0.12 \cdot 2 = 0.24
                          0
0.24 \cdot 2 = 0.48
                          0
0.48 \cdot 2 = 0.96
0.96 \cdot 2 = 1.92
                          1 and subtract one
0.92 \cdot 2 = 1.84
                          1 and subtract one
0.84 \cdot 2 = 1.68
                          1 and subtract one
0.68 \cdot 2 = 1.36
                          1 and subtract one
0.36 \cdot 2 = 0.72
0.72 \cdot 2 = 1.44
                          1 and subtract one
0.44 \cdot 2 = 0.88
0.88 \cdot 2 = 1.76
                          1 and subtract one
0.76 \cdot 2 = 1.52
                          1 and subtract one
0.52 \cdot 2 = 1.04
                          1 and subtract one
0.04 \cdot 2 = 0.08
                          0
```

This is rational, so we know it will eventually terminate or (more likely) repeat. I haven't calculated that far, but far enough at least to prove that I understand the method. So, 1100.00011110101110....

Answer here.

**Problem 3** — (8 points) Normalize the following fractional binary numbers. Give your answers in scientific notation, expressing the mantissa in binary and the exponents in decimal.

- (a) 100.0
- (b) 0.1110 0001
- (c) 1100 1010.01
- (d) 0.0001 01

## Answer:

- (a)  $1.00 \times 2^2$
- (b)  $1.1100\ 001\times 2^{-1}$
- (c)  $1.1001\ 0100\ 1\times 2^7$
- (d)  $1.01 \times 2^{-4}$

In class we did all of our examples in 64-bit double precision floating point format, but for the sake of brevity, I would like you to look up the specifications for **single precision** floating point format (which has different sized fields and a different exponent bias). Use single precision for the following exercises.

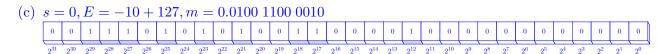
**Problem 4** — (8 points) Show how the following binary numbers would be represented in single precision floating point. In each number, the mantissa is given in binary and the exponent is given in decimal. Present your answers in binary.

- (a)  $1.01 \times 2^0$
- (b)  $-1.1000\ 1100\ 0010 \times 2^{10}$
- (c)  $1.0100\ 1100\ 0010 \times 2^{-10}$
- (d)  $1.1101\ 0001\ 1 \times 2^{-127}$

Answer: All of these use 1 sign bit, 8 bits of exponent with a bias of 127, and 23 bits of mantissa (with an understood one, except in unnormalized numbers).







(d) In this case, the value of E is 0, the non-normalized exponent value where we no longer assume a leading digit of 1. So we will have to record all of the digits of the mantissa (including the digit left of the binary point).



**Problem 5** — (8 points) Convert the following IEE754 single precision bit patterns to decimal values, using regular or scientific notation as you believe appropriate.

- (a) 0x BACO 0000
- (b) 0x CODO 0000
- (c) 0x F020 1000
- (d) 0x 7FFF FFFF

e = 224 - 127 and the mantissa is 1.01000000001. Thus, the value in decimal is  $-(1 + 2^{-2} + 2^{-11}) \times 2^{97} = -2^{97} - 2^{95} - 2^{86}$ .

(d) +NaN, although technically the sign bit is not meaningful for not-a-number, so NaN is equally (or arguably more) correct.

**Problem 6** — (5 points) The Python programming language has arbitrary integer arithmetic and double precision floating point arithmetic. Consider the following calculations:

```
>>> 2**53
9007199254740992
>>> float(2**53)
9007199254740992.0
>>> 1+2**53
9007199254740993
>>> float(1+2**53)
9007199254740992.0
>>> 3+2**53
9007199254740995
>>> float(1+2**53)
9007199254740996.0
```

Use your knowledge of floating point to explain what's going on. In your explanation, explain what happens for  $5 + 2^{53}$  and  $6 + 2^{53}$ .

Answer: The operations that are not wrapped in float() are arbitrary precision integer calculations. They are all exactly correct. When we convert them to float, they are converted to double precision floating point numbers, which we know to have 52 stored mantissa bits.

These integers would all require 53 bits of mantissa to store. So, the conversion to floating point requires rounding the mantissa. And we learned in class that a common rounding method when there is a single extra bit (as there is here) is to round toward "even mantissa" values that end with a 0 bit.

So  $1+2^{53}$  is rounded down to  $0+2^{53}$  which is the closest mantissa ending in a 0 bit, while  $3+2^{53}$  is rounded up to  $4+2^{53}$  for the same reason.

Note: The value  $2 + 2^{53}$  can be expressed exactly as a double precision floating point number, but its mantissa ends in a 1 bit, so it won't be the target of a rounding operation for the numbers above.

Similarly  $4 + 2^{53}$  and  $6 + 2^{53}$  can be expressed exactly as floating point numbers, the first ending in binary 0 and the latter ending in a 1 (so nothing will round to it). Thus, the value  $5 + 2^{53}$  will be rounded down to  $4 + 2^{53}$ , and  $6 + 2^{53}$  is represented exactly.