

Student Name:
Student ID:



Computer Systems I

Assignment 1

Problem 1

Complete the following 8-bit two's complement calculations, and represent the results by 8-bit two's complement:

1. $0101\ 0011 + 0011\ 0111$

2. $0010\ 1100 + 1010\ 0010$

3. $1111\ 1100 + 1011\ 0101$

4. $0111\ 1010 + 0011\ 1011$

5. $0110\ 0101 - 0010\ 1010$

6. $1100\ 0010 - 1111\ 1100$

7. $0110\ 1111 - 1111\ 0101$

8. $0111\ 0010 - 1000\ 0010$

Answer:

$$\begin{array}{r} 1. \ 01010011 \\ + 00110111 \\ \hline 10001010 \end{array}$$

$$\begin{array}{r} 2. \ 00101100 \\ + 10100010 \\ \hline 11001110 \end{array}$$

$$\begin{array}{r} 3. \ 11111100 \\ + 10110101 \\ \hline 10110001 \end{array}$$

$$\begin{array}{r} 4. \ 01111010 \\ + 00111011 \\ \hline 10110101 \end{array}$$

5. $\hat{A}: 11010110$

$$\begin{array}{r} 01100101 \\ + 11010110 \\ \hline 00111011 \end{array}$$

$$\begin{array}{r} 6. \ 11000010 \\ + 00001100 \\ \hline 11000110 \end{array}$$

7. $\hat{A}: 00001011$

$$\begin{array}{r} 01101111 \\ + 00001011 \\ \hline 01111010 \end{array}$$

8. $\hat{A}: 01111110$

$$\begin{array}{r} 01110010 \\ + 01111110 \\ \hline 11100000 \end{array}$$

Problem 2

We are running programs where values of type int are 32 bits. They are represented in two's complement, and they are right shifted arithmetically. Values of type unsigned are also 32 bits. We generate arbitrary values x and y, and convert them to unsigned values as follows:

```

1 /* Create some arbitrary values */
2 int x = random();
3 int y = random();
4 /* Convert to unsigned */
5 unsigned ux = (unsigned) x;
6 unsigned uy = (unsigned) y;
    
```

For each of the following C expressions, either (1) argue that it is true (evaluates to True) for all values of x and y, or (2) give values of x and y for which it is false (evaluates to False):

1. $(x < 0) == (-x > 0)$
2. $\sim x + \sim y < \sim(x + y)$
3. $((x \gg 6) \ll 6) \leq x$
4. $(ux + uy) == -(\text{unsigned})(-y - x)$

Answer:

1. $x < 0$ 对 $x \geq 0$ 对
 $x < 0 \equiv \text{True}$ $x < 0 \equiv \text{False}$
 $-x > 0 \equiv \text{True}$ $-x > 0 \equiv \text{False}$
 于是 $\forall x, 1 \equiv \text{True}$

2. $x = -\text{INT_MAX}, y = -1$
 $\sim x = \text{INT_MAX}, \sim y = 0$
 $\sim x + \sim y = \text{INT_MAX}$
 $\sim(x + y) = -\text{INT_MAX} - 1$
 False.

3. $(x \gg 6) \ll 6$ 后失去最低6位
 $x \geq 0 \quad (x \gg 6) \ll 6 < x$
 $x < 0 \quad (x \gg 6) \ll 6 < x$ 最高位为符号位
 $x = 0 \quad (x \gg 6) \ll 6 = x$

Problem 3

Consider the following two 9-bit floating-point representations based on the IEEE floating point format.

- Format A: There is 1 sign bit. There are $k = 5$ exponent bits. The exponent bias is 15. There are $n = 3$ fraction bits.
- Format B: There is 1 sign bit. There are $k = 3$ exponent bits. The exponent bias is 3. There are $n = 5$ fraction bits.

In the following table, you are given some bit patterns in format A, and your task is to convert them to the value in format B. In addition, give the values of numbers given by the format A and format B bit patterns. Give these as whole numbers (e.g., 17) or as fractions (e.g., 17/64). If the value represented by bits in format A cannot be expressed by format B, please fill **NULL** in the Value column and explain the reason briefly in **Bits** column.

| No. | Format A | | Format B | |
|-----|-------------|------------|-------------|-------------------------|
| | Bits | Value | Bits | Value |
| 1 | 0 10010 011 | 11 | 0 110 01100 | 11 |
| 2 | 1 00011 010 | $-5/16384$ | NULL | $e=0 f=0 \rightarrow 0$ |
| 3 | 0 00011 010 | $5/16384$ | NULL | $e=0 f=0 \rightarrow 0$ |
| 4 | 1 11000 000 | $-5/2$ | NULL | $+\infty$ |
| 5 | 0 10011 100 | 24 | 0 111 10000 | 24 |

Answer:

$$\begin{aligned}
 (-1)^0 \times (1.011)_2 \times 2^{18-15} &= (1011)_2 = 11 \\
 (-1)^1 \times (1.010)_2 \times 2^{3-15} &= -\frac{5}{16384} \\
 (-1)^0 \times (1.010)_2 \times 2^{3-15} &= (1.010)_2 \times 2^{-12} = \frac{5}{16384} \\
 (-1)^1 \times (1.000)_2 \times 2^{24-15} &= -5/2 \\
 (-1)^0 \times (1.100)_2 \times 2^{12-15} &= 16+8 = 24
 \end{aligned}$$

B. 最小值为 2^{-3}

最大值为 $(2-2^{-5}) \times 2^4 < 2^5 = 32$

次正规数为 $2^{-5} \times 2^{0-3} = 2^{-8}$

Problem 4

Fill in code for the following C functions, following the bit-level integer coding rules (Appendix 1). Function `srl` performs a logical right shift using an arithmetic right shift (given by value `xsra`), followed by other operations not including right shifts or division. Function `sra` performs an arithmetic right shift using a logical right shift (given by value `xsl`), followed by other operations not including right shifts or division. You may use the computation `8*sizeof(int)` to determine `w`, the number of bits in data type `int`. The shift amount `k` ranges from 0 to `w-1`.

Warning: Bit-level coding is a must!

```

1 unsigned srl(unsigned x, int k) {
2     /* Perform shift logically */
3     unsigned xsra = (int) x >> k;
4     /*
5      * Tip: Your code should be added here.
6      */
7 }
8
9 int sra(int x, int k) {
10    /* Perform shift arithmetically */
11    unsigned xsl = (unsigned) x >> k;
12    /*
13     * Tip: Your code should be added here.
14     */
15 }

```

Answer:

```

unsigned srl (unsigned x, int k) {
    /*...*/
    return x >> k;
}

int sra (int x, int k) {
    /*...*/
    unsigned ux = x;
    if (x < 0)
        ux |= ~UINT_MAX;
    unsigned shifted = ux >> k;
    if (x < 0)
        shifted |= ~UINT_MAX >> k;
    return (int) shifted;
}

```

Problem 5

Following the bit-level floating-point coding rules (Appendix 2), implement the function with the following prototype:

```
1 /* Compute -f. If f is NaN, then return f. */
2 float float_negate(float f);
```

For floating-point number f , this function computes $-f$. If f is NaN, your function should simply return f . **Warning: Bit-level coding is a must!** Testing code is shown here:

```
1 // Copyright 2023 Sycuricon Group
2 // Author: Phantom (phantom@zju.edu.cn)
3
4 #include <stdio.h>
5 typedef unsigned float_bits;
6 float float_negate(float f);
7
8 int main() {
9     printf("%f\n", float_negate(32.0));
10    printf("%f\n", float_negate(-12.8));
11    printf("%f\n", float_negate(0.0));
12    return 0;
13 }
```

Answer:

```
float float_negate(float f) {
    float_bits *p = (float_bits*)&f;
    if ((*p >> 23 & 0xFF) == 0xFF && (*p & 0x7FFFFFFF) != 0x0)
        return f;
}
*p ^= 1 << 31;
return f;
}
```

Problem 6

Explain why the following code snippets are wrong:

1. Code Snippet 1

```
1 unsigned i;
2 int cnt = 10, sum = 0;
3 for (i = cnt - 1; i >= 0; i --) {
4     sum += i
5 }
6 print(sum);
```

2. Code Snippet 2

```
1 #define Delta sizeof(int)
2 int i;
3 int cnt = 10, sum = 0;
4 for (i = cnt; i - Delta >= 0; i -= Delta) {
5     sum += i;
6 }
7 print(sum);
```

Answer:

Snippet 1: For unsigned i , $i >= 0 \equiv \text{True}$, it'll loop forever.

Snippet 2: `sizeof()` returns computes a `size_t` type (unsigned). When $i - \text{Delta}$ is performed, it occurs an implicit conversion from `int` to unsigned number, thus $i - \text{Delta} >= 0 \equiv \text{True}$ all loop forever.