

浮点数运算模拟

3180103772 张溢弛

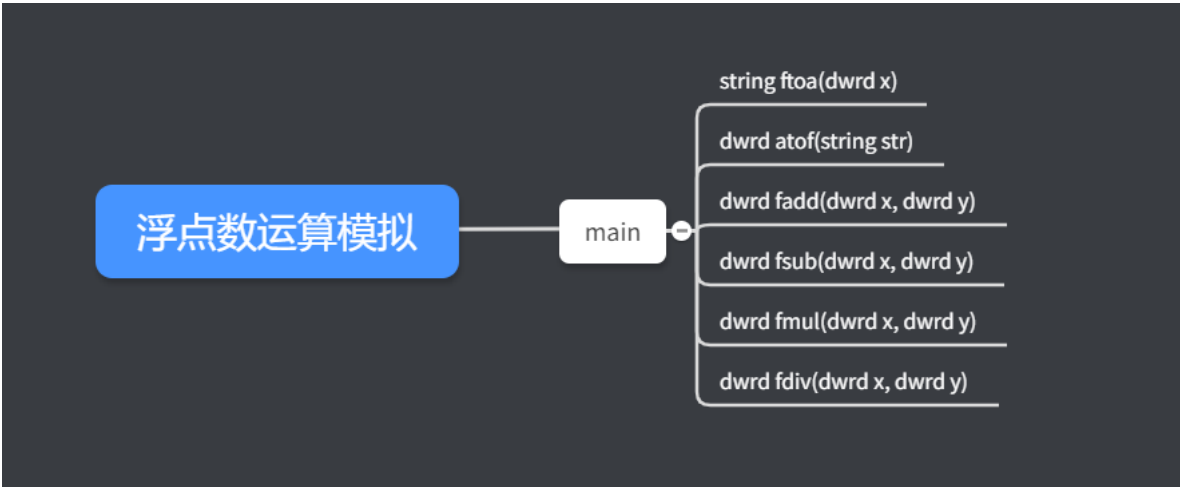
0.实验描述

浮点数的表示与算术运算算法分析，要求理论推导与程序模拟。算术运算包括字符串到浮点数、浮点数到字符串的转换，加、减、乘、除四则运算等。要有说明文档，包括算法证明、程序框图、使用方法、特殊处理(溢出、数位扩展)、实例分析等等。字符串转换可能稍难。

```
1 typedef unsigned int dwrd;
2 char* ftoa(dwrd);
3 dwrd atof(char*);
4 dwrd fadd(dwrd, dwrd);
5 dwrd fsub(dwrd, dwrd);
6 dwrd fmul(dwrd, dwrd);
7 dwrd fdiv(dwrd, dwrd);
```

1.程序结构

- 本程序的结构图如下



2.算法分析

- 加法和减法
 - 浮点数加法的运算规则为小阶向大阶看齐，根据浮点数的格式提取出ex，ey分别为浮点数x和y的阶码，而mx，my为尾数，将小阶对应的尾数右移后向大阶看齐，然后将位数相减得到新的尾数，相减后进行重新的规格化。
 - 减法可以将被减数转化为相反数相加，只要改变被减数的符号位后再进行加法就可以

```
1 dwrd fadd(dwrd x, dwrd y)
2 {
3     if(x == 0) return y;
4     else if(y == 0) return x;
5     if((x << 1) < (y << 1)) return fadd(y, x);
6     dwrd ex = (x >> 23) & 0xFF;
```

```

7     dwrdr ey = (y >> 23) & 0xFF;
8     dwrdr mx = (x & 0x7FFFFFFF) | 0x800000;
9     dwrdr my = (y & 0x7FFFFFFF) | 0x800000;
10    my >>= (ex - ey);
11    if((x ^ y) & 0x80000000) {
12        mx -= my;
13        if(mx == 0) return 0;
14        mx <<= 1;
15        ex--;
16    }
17    else {
18        mx += my;
19        if (mx == 0) return 0;
20        while (mx & 0xFF000000) {
21            mx >>= 1;
22            ex++;
23        }
24    }
25    return (x & 0x80000000) | (ex << 23) | (mx & 0x7FFFFFFF);
26 }
27
28 dwrdr fsub(dwrdr x, dwrdr y)
29 {
30     return fadd(x, y ^ 0x80000000);
31 }

```

- 乘法

- 浮点数乘法的运算规则为

- $E(xy) = x + M, E(xy) = x + y + m = E(x) + E(y) - M$
- 根据真值表，浮点数乘法运算的符号位是两个因数最高位的异或

- 乘法的尾数即为两个浮点数的尾数值相乘，但是由于一个浮点数的尾数有24位，两个浮点数相乘的时候会产生48位的尾数，容易导致溢出，所以只保留相乘的前24位

```

1     dwrdr fmul(dwrdr x, dwrdr y)
2     {
3         if(x == 0 || y == 0) return 0;
4         dwrdr ex = (x >> 23) & 0xFF;
5         dwrdr ey = (y >> 23) & 0xFF;
6         dwrdr mx = (x & 0x7FFFFFFF) | 0x800000;
7         dwrdr my = (y & 0x7FFFFFFF) | 0x800000;
8         dwrdr s = (x >> 31) ^ (y >> 31);
9         dwrdr e = ex + ey - 127;
10        dwrdr result = 0;
11        while(my != 0) {
12            result >>= 1;
13            if(my & 1) result += mx;
14            my >>= 1;
15        }
16        while(result & 0xFF000000) {
17            result >>= 1;
18            e++;
19        }
20        return (s << 31) | (e << 23) | (result & 0x7FFFFFFF);
21    }

```

- 除法

- 除法的阶码运算规则如下

- $E(x/y) = x - y + M = E(x) - E(y) + M$
 - 尾数的除法采用恢复除数的除法
 - 首先算出阶码
 - 然后对尾数做出除法
 - 为了处理最高位商0导致移位失败，我们在开始时通过不断地位移使得 $mx > my$

```
1  dwrd fdiv(dwrd x, dwrd y)
2  {
3      if(x == 0) return 0;
4      if (y == 0) return (x & 0x80000000) | 0x7F800000;
5      dwrd ex = (x >> 23) & 0xFF;
6      dwrd ey = (y >> 23) & 0xFF;
7      dwrd mx = (x & 0x7FFFFFFF) | 0x800000;
8      dwrd my = (y & 0x7FFFFFFF) | 0x800000;
9      dwrd s = (x >> 31) ^ (y >> 31);
10     dwrd e = ex - ey + 127;
11     dwrd result = 0;
12     while(mx < my) {
13         mx <<= 1;
14         e--;
15     }
16     while(!(result & 0xFF800000)) {
17         if(mx < my) result <<= 1;
18         else {
19             mx -= my;
20             result = (result << 1) | 1;
21         }
22         mx <<= 1;
23     }
24     return (s << 31) | (e << 23) | (result & 0x7FFFFFFF);
25 }
```

- ftoa的实现

- 浮点数转换成字符串的过程比较麻烦，首先要考虑特殊情况，如果是0则直接输出0，如果阶码是255则表示无穷大，用“INF”来表示
 - 然后考虑符号位，根据最高位判断符号位添加到字符串中
 - 接下来分两种情况讨论，看x是大于1的数还是小于1的数
 - 如果大于1，需要分成整数部分和小数部分进行转换
 - 如果小于1，则不需要管整数的部分

```
1  string ftoa(dwrd x)
2  {
3      string result;
4      dwrd sx = x >> 31;
5      dwrd ex = (x >> 23) & 0xFF;
6      dwrd mx = (x & 0x7FFFFFFF) | 0x800000;
7      if(ex == 0) {
8          if(mx == 0) result = "0";
9          else {
10
```

```

11     }
12     }
13     else if(ex == 0xFF) {
14         if(mx == 0) {
15             if(!sx) result = "INF";
16             else result = "-INF";
17         }
18         else result = "NaN";
19     }
20     else if(ex >= 127) {
21         dwrd e = ex - 127;
22         string integer = "1";
23         string decimal = ".";
24         int i = 0;
25         while(i < 24) {
26             if(i < e) {
27                 if((mx >> (22 - i)) & 1) integer += "1";
28                 else integer += "0";
29             }
30             else {
31                 if((mx >> (22 - i)) & 1) decimal += "1";
32                 else decimal += "0";
33             }
34             i++;
35         }
36         if(!sx) result = trans1(integer) + "." + trans2(decimal);
37         else result = "-" + trans1(integer) + "." + trans2(decimal);
38     }
39     else {
40         dwrd e = 127 - ex;
41         string decimal = ".";
42         int i = 0, j = 0;
43         while(i < e) {
44             decimal += "0";
45             i++;
46         }
47         while(j < 24) {
48             if((mx >> (22 - j)) & 1) decimal += "1";
49             else decimal += "0";
50             j++;
51         }
52         if(!sx) result = "0." + trans2(decimal);
53         else result = "-0." + trans2(decimal);
54     }
55     return result;
56 }

```

• atof的实现

- 在将字符串转换成浮点数之前，首先可以通过正则表达式来验证是否是一个合法的浮点数
- 为了能够分别计算整数部分和小数部分，需要将字符串分成符号位，整数和小数三个部分，拆分后需要考虑浮点数为0的特殊情况，此时应该直接返回0
- 整数部分的转化思路为除2取余，得到整数部分的二进制串，小数部分的转换比较麻烦，通过字符串切割得到小数部分将其当作整数转换成二进制一样的操作来转换，最后进行拼接

```

1 dwrd atof(string str)
2 {

```

```

3     string binary;
4     dwrd s, e, m = 0;
5     if(str[0] == '-') {
6         s = 1;
7         str.erase(str.begin());
8     }
9     else s = 0;
10    int i = 1, temp = (int)str.find('.');
11    if(temp == -1) binary = form1(str, 0, (int)str.length() - 1);
12    else {
13        binary = form1(str, 0, temp) + '.' + form2(str, temp + 1,
14        int(str.length()));
15    }
16    if(binary.find('.') != -1) {
17        e = (dwrd)binary.find('.') + 126;
18        binary.erase(binary.find('.'), 1);
19    }
20    else e = (dwrd)binary.find('.') + 126;
21    int len = (int)binary.length();
22    while(i < 24 && i < len) {
23        if(binary[i] == '1') m += 1 << (23 - i);
24        i++;
25    }
26    return (s << 31) | (e << 23) | m;

```

3.实验结果

- 根据提示在一行中输入两个浮点数
- 两个小数部分为0的整数

```

Please input 2 float number in a line
1.0 2.0
num1 = 1.000000000000000000
num2 = 2.000000000000000000
add = 3.000000000000000000
sub = - 1.000000000000000000
mul = 2.000000000000000000
div = 0.000000000000000000

```

- 普通的两个数

```

Please input 2 float number in a line
3.4 5.6
num1 = 3.0000000002340586730
num2 = 5.0000000000899963036
add = 8.0000000003240549766
sub = - 2.0000000003781210424
mul = 19.0000000001312362178
div = 0.0000000004110185660

```

- 包含负数

```
-1.2 2.5
num1 = - 1.0000000003781210424
num2 = 2.0000000001156841472
add = 1.0000000002597465166
sub = - 3.0000000000643084600
mul = - 2.0000000003240549766
div = -0.0000000002533156706
```

4.心得体会

- 这次实验通过编程模拟计算机中的浮点数运算，我收获颇丰，更透彻理解了计算机中的一些底层工作原理，希望在下周开始的CPU工作原理等内容中能学到更多东西。