

Floating Point Numbers

IEEE-754 representation

Single Precision - 32 bits

s	exponent	significand
1 bit	8 bits	23 bits

Double Precision - 64 bits

s	exponent	significand
1 bit	11 bits	52 bits

IEEE-754 fields

- Sign bit: 1 if -ve
- Exponent: value of exponent + 127 (single precision) -special exponent values 0 and 255 used to code for 0 and Nan and Infinity
- Significand: Leading bit is implicit - this field contains the bits *after* the binary point

Elaboration (float)

<u>Exponent</u>	<u>Significand</u>	<u>Object represented</u>
0	0	0
1-254	anything	floating-point number
255	0	infinity
255	nonzero	Nan
0	nonzero	denormalized number

C Examples

```
#include <conio.h>
#include <stdio.h>
union dami {
    float f;
    long l;
} tf[4];

void main()
{
    FILE *fp;
    fp=fopen ("xxx.txt", "w");
    clrscr();
    tf[0].f = 1.0;
    tf[1].f = -1.0;
    tf[2].f = 0.0;
    tf[3].f = 1.25;
    printf(fp, "%08lx\n", tf[0].l);
    printf(fp, "%08lx\n", tf[1].l);
    printf(fp, "%08lx\n", tf[2].l);
    printf(fp, "%08lx\n", tf[3].l);
}
```

Result:

3f800000
bf800000
00000000
3fa00000
