# Floating point representation

A scheme for representing a number very small to very large. It is

widely used in the scientific world. Consider, the floating point number

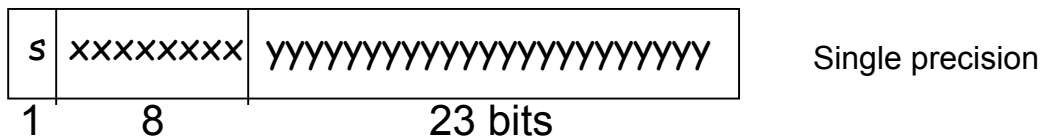Exponent E   Mantissa or Significand F

| +/- | x  x  x  x | y  y  y  y  y  y  y  y  y  y  y  y |
|-----|-----------|-----------------------------------|

In decimal it means  (+/-) **D**. yyyyyyyyyyyy x $10^{xxxx}$ (D>0)

In binary, it means   (+/-) 1. yyyyyyyyyyyy x $2^{xxxx}$

(The 1 is implied)

## IEEE 754 representation

| s | xxxxxxxx | yyyyyyyyyyyyyyyyyyyyyyyyy |
|---|----------|--------------------------|
| 1 | 8        | 23 bits                  |

Single precision

Largest = 1.111… x $2^{+127}$  ≈ 2 x $10^{+38}$

Smallest = 1.000 … x $2^{-128}$  ≈ 1 x $10^{-38}$

Sign = $(-1)^{S}$

These can be positive and negative, depending on s.

## IEEE 754 double precision (64 bits)

| S | exponent | significand |
|---|----------|-------------|

1    11 bits                                    52 bits

Largest =        $1.111... \times 2^{+1023}$

Smallest =       $1.000... \times 2^{-1024}$

## Overflow and underflow in FP

An overflow occurs when the number if <span style="color:red">too large to fit</span> in the

frame. An underflow occurs when the number <span style="color:red">is too small to</span>

<span style="color:red">fit</span> in the given frame.

## Biased Representation

Exponent = 11111111    $2^{-1}$   ⎫   awkward for sorting

Exponent = 00000000    $2^{0}$   ⎭

However, to facilitate sorting, IEEE 754 treats 00…0 as the most negative, and 1,11..1 as the most positive exponent. This amounts to using a bias of 127.

00000000 (=-127)                          11111111 (=+128)

bias +127

▲                                 ▲

smallest                               largest

⟶

So, value = $(1)^{s}$ x (1+significand) x $2^{(\text{exponent} - \text{bias})}$

-8
0000

7
1111

-7
0001

6
1110

-6
0010

5

-1   0   1

-2                    2

                      3

0011

                      4   0100

5                     5   0101

                      6   0110

-7   -8   7

1000

0111

**<u>Floating Point Addition</u>**

1. Align (Shift the smaller number to the right until the exponents are equal)

2. Add the significand

3. Normalize

4. **If** underflow or overflow **then**

   round-off to the right number of bits

   **else** flag exception

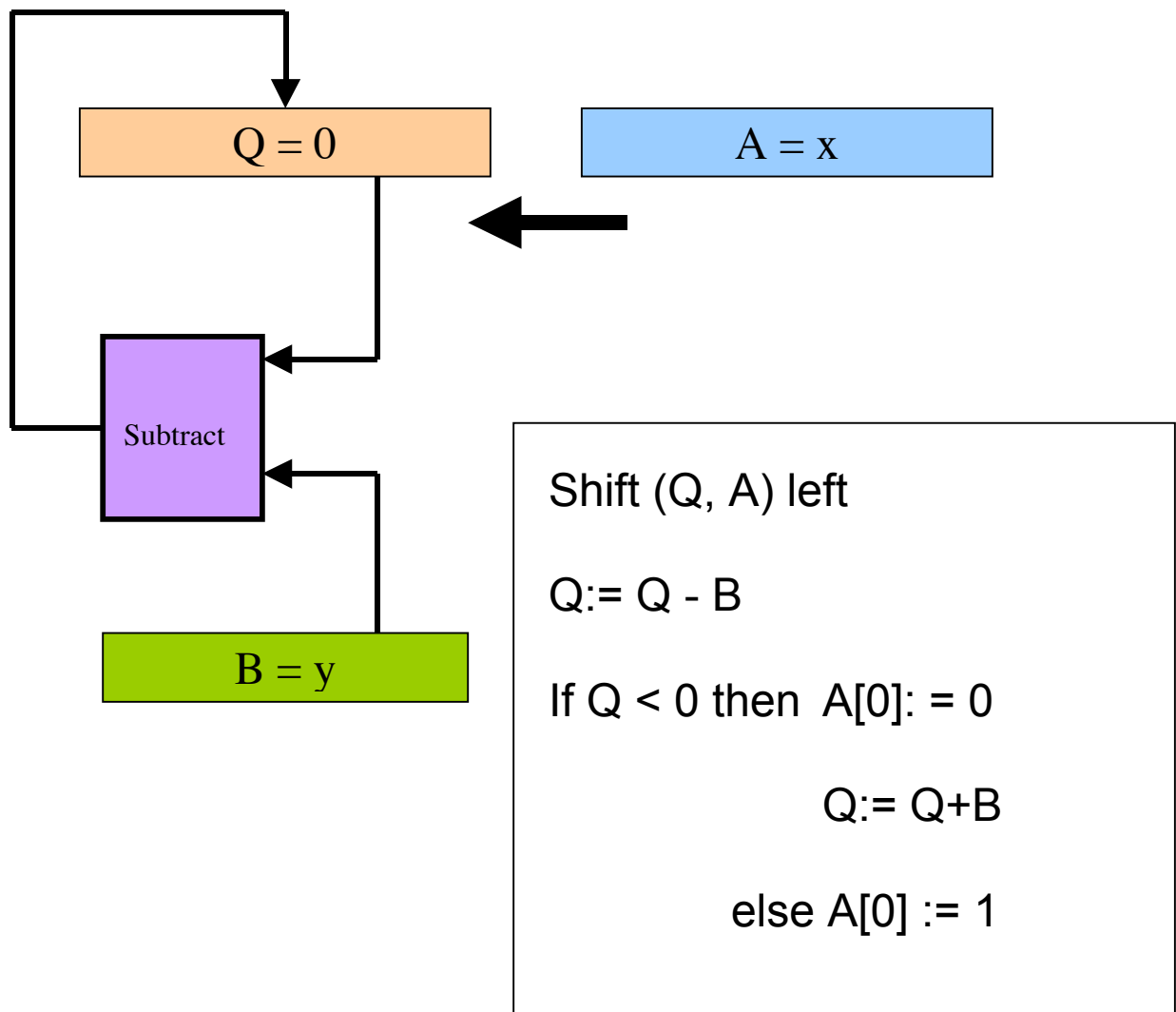Note that these stages can be pipelined.

# Floating Point Multiplication

1. Add the (biased) exponents and subtract the bias to get the new exponent

2. Multiply the significands

3. Normalize (if necessary)

4. If overflow or underflow then exception else round off the significand

5. Set the sign appropriately

These steps can be pipelined, if necessary

# Restoring division algorithm for integer operands

Divide x by y: quotient = q, remainder = r



Shift (Q, A) left

Q:= Q - B

If Q < 0 then  A[0]: = 0

          Q:= Q+B

        else A[0] := 1

After N cycles, A = quotient, Q = remainder