



3.5.3 浮点数的**乘除**运算

刘 芳 副教授

国防科学技术大学计算机学院



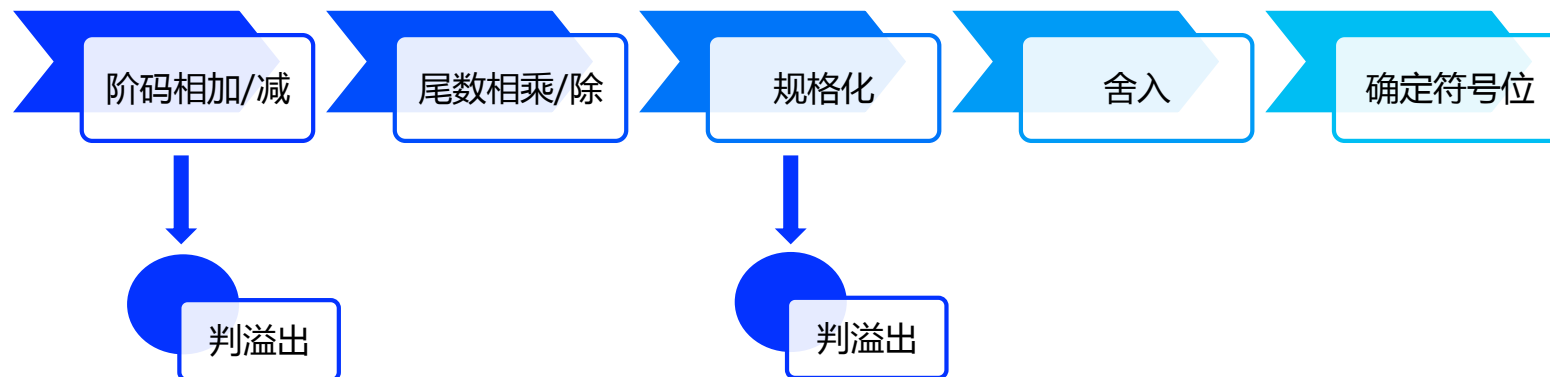
浮点数乘除

两个浮点数： $A = M_a \cdot 2^{E_a}$, $B = M_b \cdot 2^{E_b}$

$$A \times B = (M_a \times M_b) \cdot 2^{E_a + E_b}$$

$$A \div B = (M_a \div M_b) \cdot 2^{E_a - E_b}$$

计算机实现步骤：





1. 阶码运算

阶码相加（单精度浮点乘法）

$$[E_1 + E_2]_{\text{移}} = f([E_1]_{\text{移}}, [E_2]_{\text{移}}) = \boxed{([E_1]_{\text{移}} + [E_2]_{\text{移}} + 129) \bmod 2^8}$$



$$\begin{aligned} [E_1 + E_2]_{\text{移}} &= 127 + E_1 + E_2 \\ &= 127 + E_1 + 127 + E_2 - 127 \\ &= [E_1]_{\text{移}} + [E_2]_{\text{移}} - 127 \\ &= [E_1]_{\text{移}} + [E_2]_{\text{移}} + [-127]_{\text{补}} \\ &= ([E_1]_{\text{移}} + [E_2]_{\text{移}} + 10000001B) \bmod 2^8 \end{aligned}$$



1. 阶码运算

阶码相加（单精度浮点乘法）

例：若两个IEEE754操作数的阶码分别为10和 - 5，求 $10 + (-5)$ 的移码

解：

$$[E_1]_{\text{移}} = 127 + 10 = 137 = 1000\ 1001\text{B}$$

$$[E_2]_{\text{移}} = 127 + (-5) = 122 = 01111010\text{B}$$

$$\begin{aligned}[E_1 + E_2]_{\text{移}} &= [E_1]_{\text{移}} + [E_2]_{\text{移}} + 129 \\ &= 1000\ 1001 + 0111\ 1010 \\ &\quad + 1000\ 0001(\text{mod } 2^8) \\ &= 1000\ 0100\text{B} = 132\end{aligned}$$

其阶码的和为 $132 - 127 = 5$ ，正好等于 $10 + (-5) = 5$



1. 阶码运算

阶码相加（单精度浮点乘法）

$$[E_1 - E_2]_{\text{移}} = f([E_1]_{\text{移}}, [E_2]_{\text{移}}) = ([E_1]_{\text{移}} + [-[E_2]_{\text{移}}]_{\text{补}} + 127) \bmod 2^8$$



$$\begin{aligned} [E_1 - E_2]_{\text{移}} &= E_1 - E_2 + 127 \\ &= 127 + E_1 - (127 + E_2) + 127 \\ &= [E_1]_{\text{移}} - [E_2]_{\text{移}} + 127 \\ &= ([E_1]_{\text{移}} + [-[E_2]_{\text{移}}]_{\text{补}} + 01111111\text{B}) \bmod 2^8 \end{aligned}$$



1. 阶码运算

阶码相加（单精度浮点乘法）

例：若两个IEEE754操作数的阶码分别为10和 - 5，求 $10 - (-5)$ 的移码

解：

$$[E_1]_{\text{移}} = 127 + 10 = 137 = 1000\ 1001\text{B}$$

$$[E_2]_{\text{移}} = 127 + (-5) = 122 = 01111010\text{B}$$

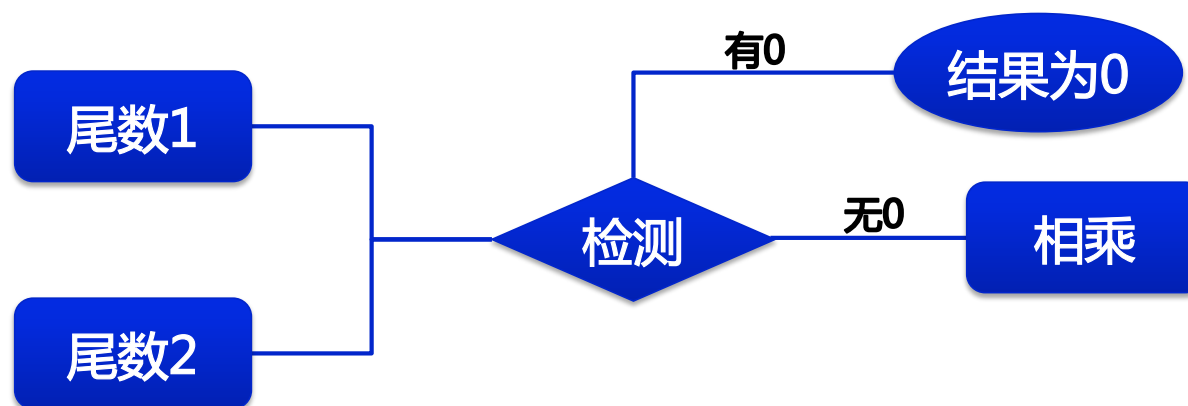
$$\begin{aligned}[E_1 + E_2]_{\text{移}} &= [E_1]_{\text{移}} + [-[E_2]_{\text{移}}]_{\text{补}} + 127 \\ &= 1000\ 1001 + 1000\ 0110 \\ &\quad + 0111\ 1111(\text{mod } 2^8) \\ &= 1000\ 1110\text{B} = 142\end{aligned}$$

其阶码的差为 $142 - 127 = 15$ ，正好等于 $10 - (-5) = 15$



2. 尾数运算

尾数相乘 预处理



相乘

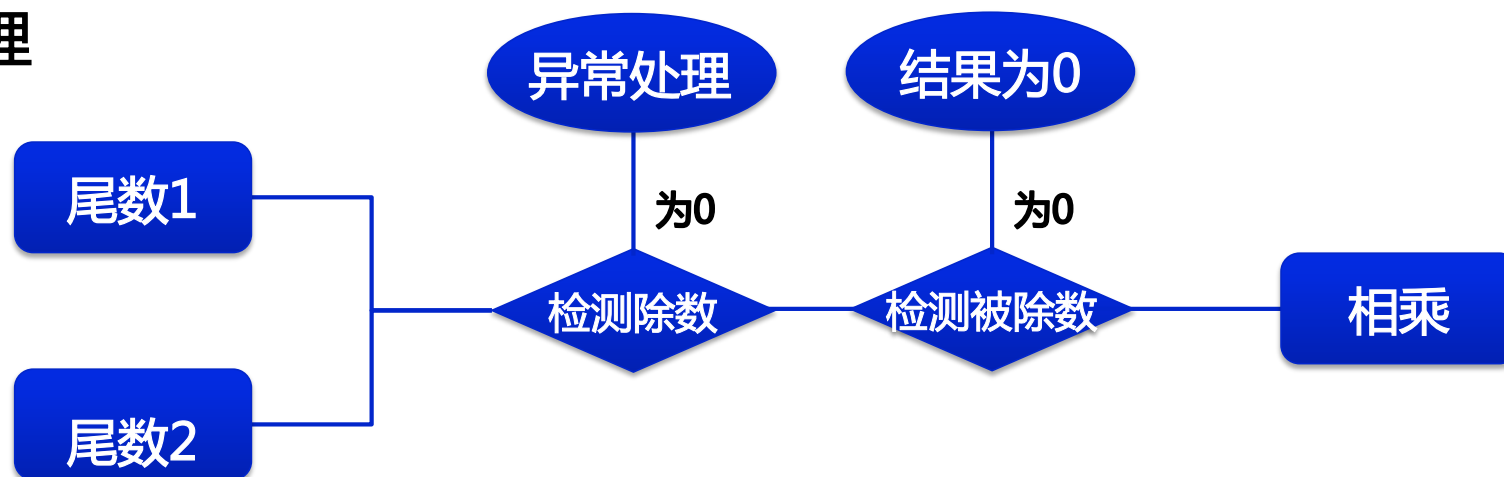
- 定点小数原码乘法运算



2. 尾数运算

尾数相乘

预处理



相乘

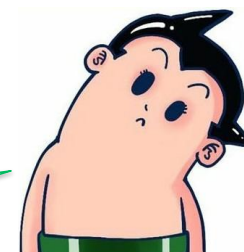
- 定点小数原码乘法运算



3. 规格化

规格化原则

- 当尾数高位为0，左规
- 当尾数产生进位，右规



乘法运算结果最多左规几次、右规几次？
除法呢？

规格化尾数： $1 \leq \text{尾数}_1 < 2$ ， $1 \leq \text{尾数}_2 < 2$ (尾数形为1.xxx)

乘法： $1 \leq \text{尾数}_1 \times \text{尾数}_2 < 4$ 不需左规、最多右规一次

除法： $0.5 < \text{尾数}_1 \div \text{尾数}_2 < 2$ 不需右规、最多左规一次



4. 舍入并确定符号

舍入

- 就近舍入
- 朝 $+\infty$ 舍入
- 朝 $-\infty$ 舍入
- 朝0舍入



若尾数是0，则需要将阶码也置0

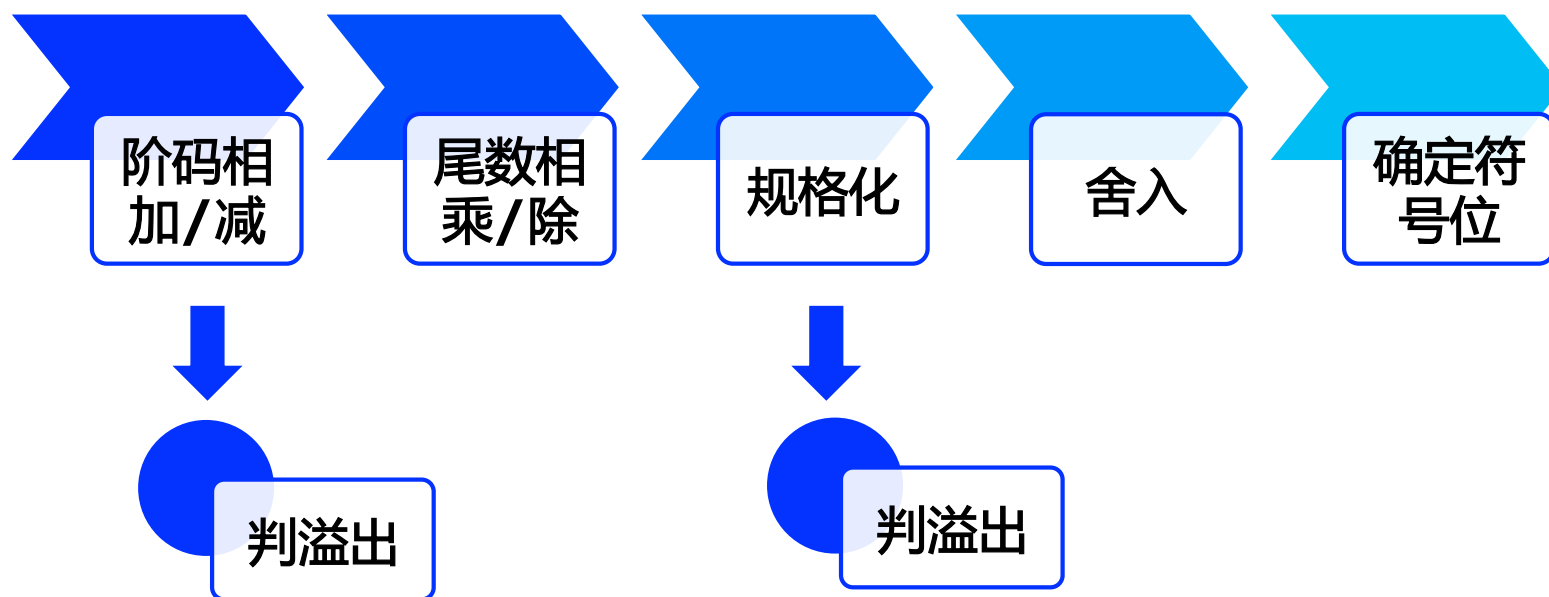
符号确定

- 同号：结果为正
- 异号：结果为负



5. 判溢出

阶码溢出情况





5. 判溢出

阶码溢出情况

阶码求和

$$E_x + E_y + 129 = E_b$$

阶码求差

| | | | |
|-------|-------|-------|----|
| 1xxxx | 1xxxx | 0xxxx | 上溢 |
| 0xxxx | 0xxxx | 1xxxx | 下溢 |

规格化

| | |
|-------|----|
| 11111 | 上溢 |
| 00000 | 下溢 |



5. 判溢出

阶码溢出情况

阶码求和

$$E_x - E_y + 127 = E_b$$

阶码求差

| | | | |
|-------|-------|-------|----|
| 1xxxx | 0xxxx | 0xxxx | 上溢 |
| 0xxxx | 1xxxx | 1xxxx | 下溢 |

规格化

| | |
|-------|----|
| 11111 | 上溢 |
| 00000 | 下溢 |



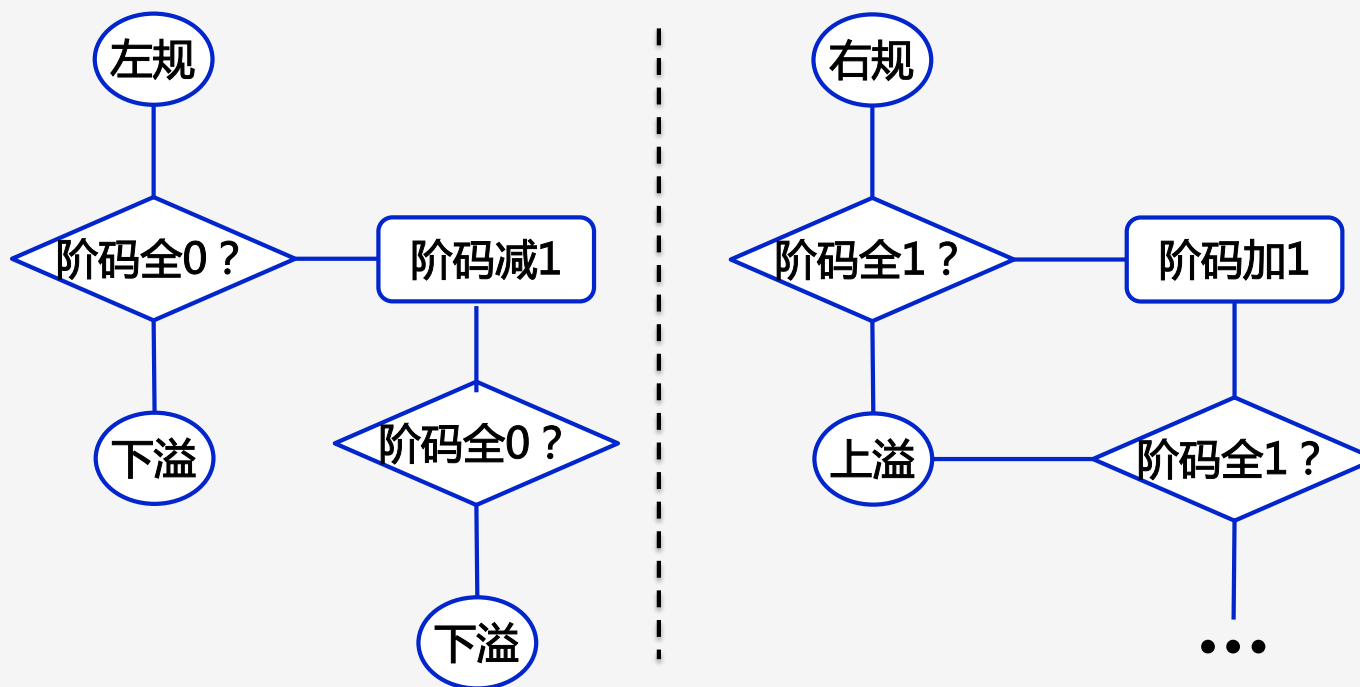
5. 判溢出

阶码溢出情况

阶码求和

阶码求差

规格化





5. 判溢出

溢出判断举例

例1：若 $E_b=0000\ 0001$ ，则尾数左规一次后，结果的阶码 $E_b=?$

解： $E_b = E_b + [-1]_{\text{补}} = 0000\ 0001 + 1111\ 1111$
 $= 0000\ 0000$ 阶码下溢！

例2：若 $E_x=1111\ 1110$ ， $E_y=1000\ 0000$ ，则乘法运算时，结果的阶码 $E_b=?$

解： $E_b = E_x + E_y + 129 = 1111\ 1110 + 1000\ 0000 + 1000\ 0001$
 $= 1111\ 1111$ 阶码上溢！



浮点数乘法实例

例：用二进制的形式求出浮点数 0.5_{10} 与 -0.4375_{10} 之积(假设保留4位有效数位的精度)。

解： $0.5_{10} = 1/2_{10} = 0.1_2 = 1.000_2 \times 2^{-1}$
 $-0.4375_{10} = -7/16_{10} = -0.0111_2 =$
 $-1.110_2 \times 2^{-2}$

■阶码相加： $-1 + (-2) = -3$

■尾数相乘：

$$1.000_2 \times 1.110_2 = 1.110000_2$$

乘积为： $1.110000_2 \times 2^{-3}$

$$\begin{array}{r} 1.000_{\text{two}} \\ \times 1.110_{\text{two}} \\ \hline 0000 \\ 1000 \\ 1000 \\ 1000 \\ \hline 1110000_{\text{two}} \end{array}$$



浮点数乘法实例

例：用二进制的形式求出浮点数 0.5_{10} 与 -0.4375_{10} 之积(假设保留4位有效数位的精度)。

解： $0.5_{10} = 1/2_{10} = 0.1_2 = 1.000_2 \times 2^{-1}$
 $-0.4375_{10} = -7/16_{10} = -0.0111_2 = -1.110_2 \times 2^{-2}$

■阶码相加： $-1 + (-2) = -3$

■尾数相乘：

$$1.000_2 \times 1.110_2 = 1.110000_2$$

乘积为： $1.110000_2 \times 2^{-3}$

■规格化并判溢出：

尾数 1.110000 已是规格化数，且 $127 \geq -3 \geq -126$ ，没有溢出！

■舍入： $1.110_2 \times 2^{-3}$

■确定符号位：

乘积为负数， $-1.110_2 \times 2^{-3}$

所以浮点乘法结果：

$$\begin{aligned} -1.110_2 \times 2^{-3} &= -0.00111_2 \\ &= -7/32_{10} \\ &= -0.21875_{10} \end{aligned}$$



浮点数除法实例

例：用二进制的形式求出浮点数 0.5_{10} 与 -0.4375_{10} 的商(假设保留4位有效数位的精度)。

解： $0.5_{10} = 1/2_{10} = 0.1_2 = 1.000_2 \times 2^{-1}$
 $-0.4375_{10} = -7/16_{10} = -0.0111_2 = -1.110_2 \times 2^{-2}$

■阶码相减： $-1 - (-2) = 1$

■尾数相除

$: 1.000_2 \div 1.110_2 = 0.1001001_2$

商为： $0.1001001_2 \times 2^1$

■规格化并判溢出：

规格化： $1.0010010_2 \times 2^0$

且 $127 \geq 0 \geq -126$ ，没有溢出

■舍入： $1.001_2 \times 2^0$

■确定符号位：

商为负数， $-1.001_2 \times 2^0$

所以浮点除法结果：

$$\begin{aligned} -1.001_2 \times 2^0 &= -1.001_2 \\ &= -9/8_{10} \\ &= -1.125_{10} \end{aligned}$$



浮点数运算小结

浮点数运算：由多个ALU + 移位器实现 加/减运算

- 对阶、尾数相加减、规格化处理、舍入、判断溢出

乘/除运算

- 尾数用定点原码乘/除运算实现，阶码用定点数加/减运算实现

溢出判断

- 当结果发生阶码上溢时，结果发生溢出；发生阶码下溢时，结果为0

精确表示运算结果

- 中间结果增设保护位、舍入位，等
- 最终结果四种舍入方式：就近舍入 / 正向舍入 / 负向舍入 / 截去