



3.5.2 浮点数的加法运算

刘 芳 副教授

国防科学技术大学计算机学院



3.5.2 浮点数加法

十进制科学计数法的加法实例：

$$A=1.23 \times 10^5; B=4.56 \times 10^2; \text{求 } A+B$$

其计算过程为：

$$\begin{aligned} & 1.23 \times 10^5 + 4.56 \times 10^2 \\ = & 1.23 \times 10^5 + 0.00456 \times 10^5 \\ = & (1.23 + 0.00456) \times 10^5 = 1.23456 \times 10^5 \end{aligned}$$

进行尾数加法运算前，必须“对阶”！

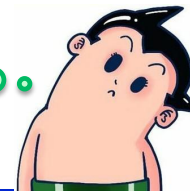


3.5.2 浮点数加法

对阶：目的是使两个操作数的阶码相等(对齐小数点)

- **规则：**小阶向大阶看齐，阶小的数的尾数右移，右移位数等于两个阶码差的绝对值。

为什么？



IEEE754尾数右移时，需要注意的是什么？

- 1、要将**隐含的“1”**移到小数部分，空出位补**0**；
- 2、移出的低位保留到特定的**“附加位”**上



3.5.2 浮点数加法

两个规格化浮点数分别为A和B

- $A = M_a \cdot 2^{E_a}$, $B = M_b \cdot 2^{E_b}$
- 假设 $E_a \geq E_b$

浮点数加法步骤

- 求阶差 $E_a - E_b$
- 对阶 $M_b \cdot 2^{-(E_a - E_b)}$
- 尾数相加 $M_a + M_b \cdot 2^{-(E_a - E_b)}$
- 结果规格化 $A + B = (M_a + M_b \cdot 2^{-(E_a - E_b)}) \cdot 2^{E_a}$



浮点数运算及结果

两个规格化浮点数分别为A和B

- $A = M_a \cdot 2^{E_a}$, $B = M_b \cdot 2^{E_b}$
- 假设 $E_a \geq E_b$

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

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

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



3.5.2 浮点数加法



自然世界中的浮点数



无穷位



计算机世界中的浮点数

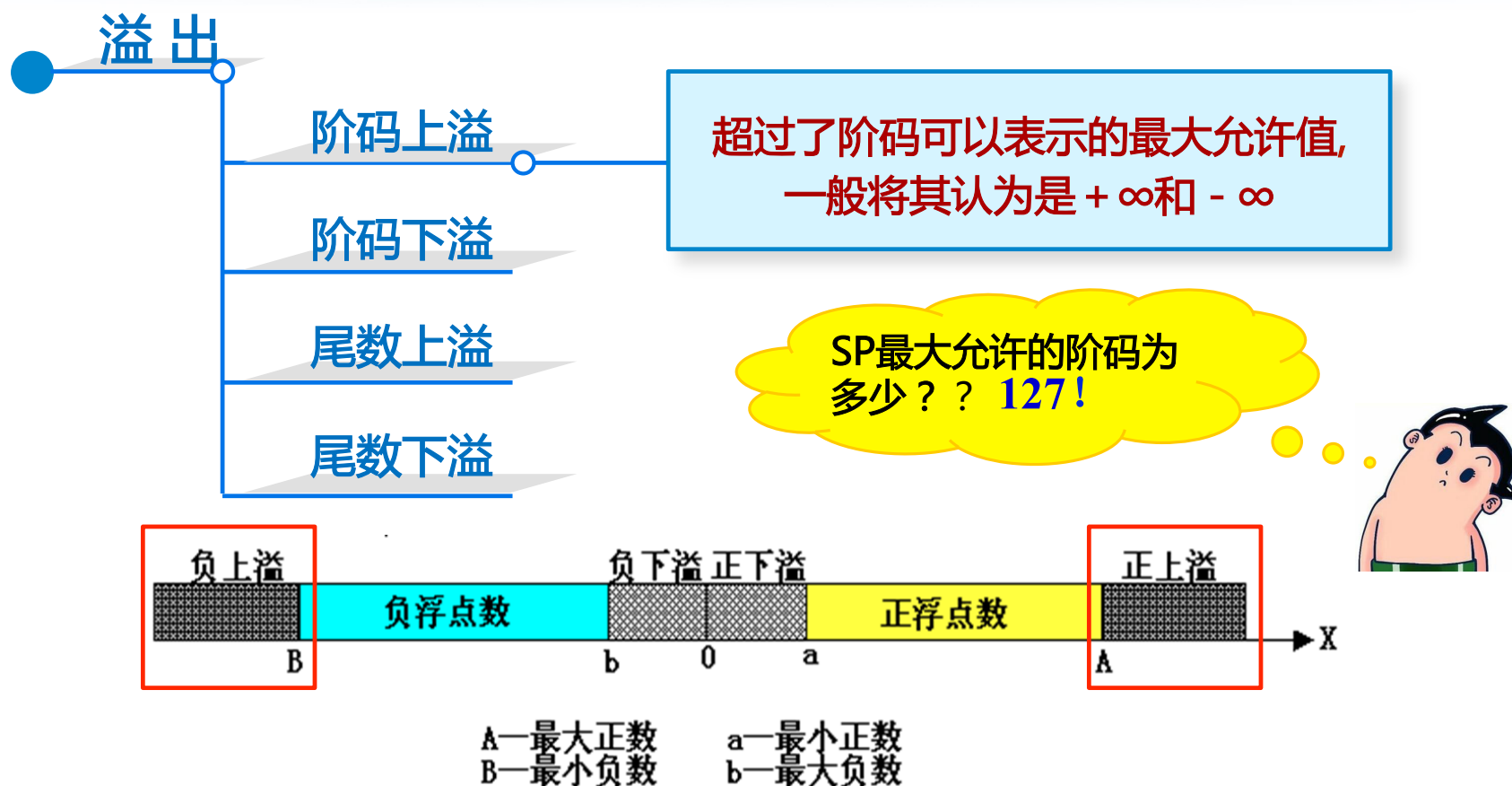


有限位

超出表示范围：
溢出！

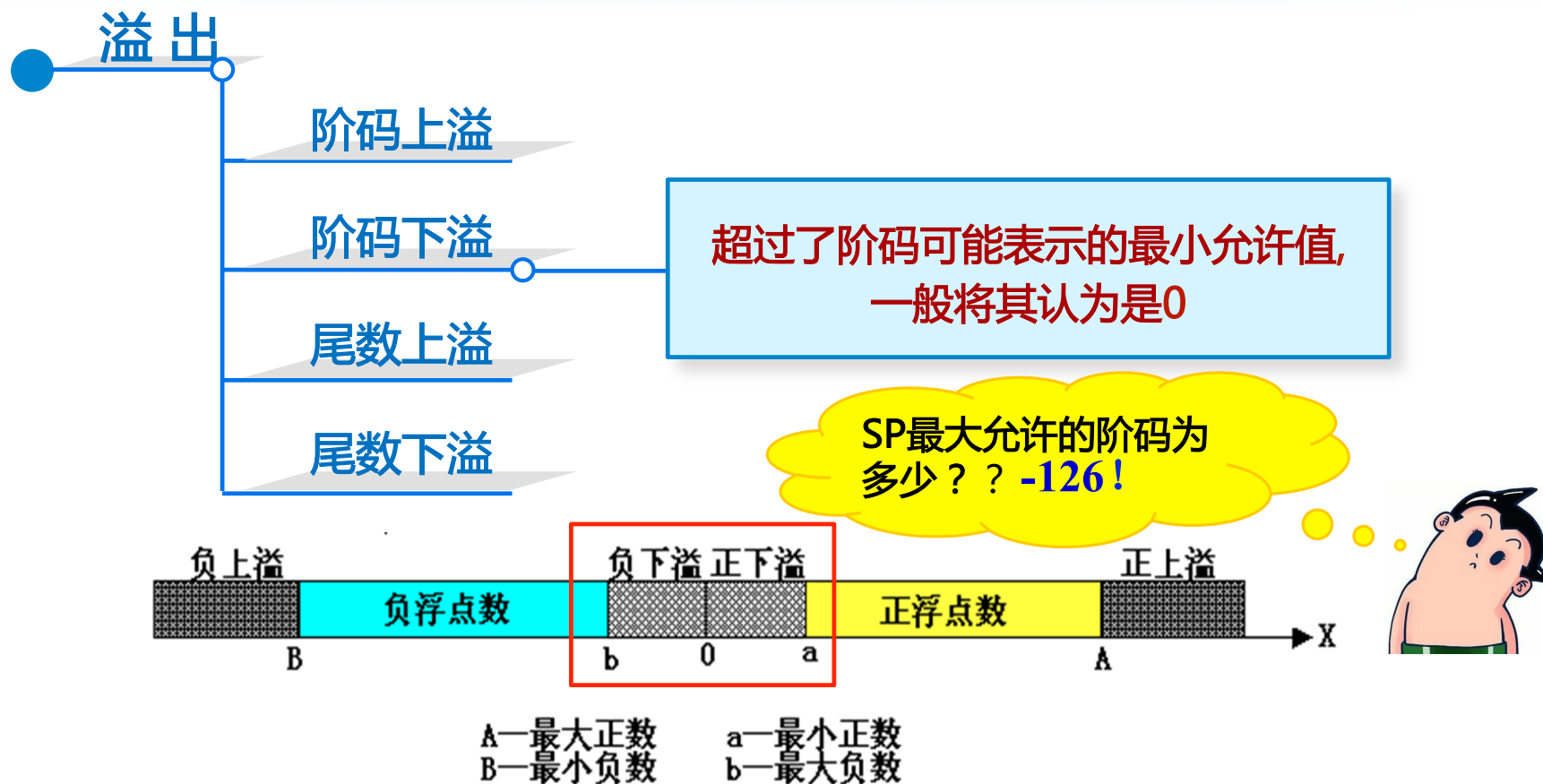


3.5.2 浮点数加法





3.5.2 浮点数加法





3.5.2 浮点数加法



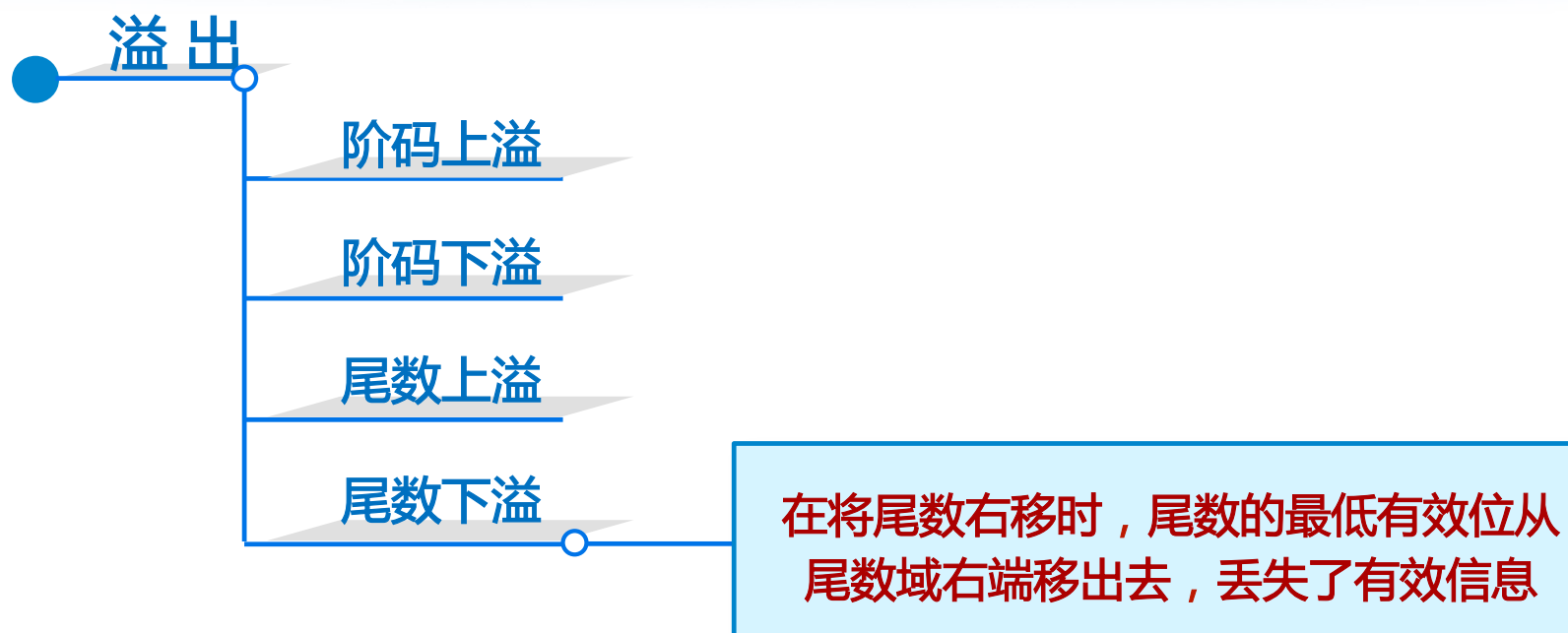
尾数溢出：

- 不一定浮点数溢出，即不一定会发生“异常”
- 右规：将尾数右移，阶码增1来重新对齐

$$\begin{aligned} & 1.010...0_2 \times 2^2 \\ & + 1.100...0_2 \times 2^2 \\ \hline & = 10.110...0_2 \times 2^2 \\ & => 1.0110...0_2 \times 2^3 \end{aligned}$$



3.5.2 浮点数加法



- 1、进行舍入处理
- 2、在运算过程中，添加保护位



3.5.2 浮点数加法

浮点数加法步骤

- 求阶差
- 对阶
- 尾数相加
- 规格化并判溢出
- 舍入
- 置0



3.5.2 浮点数加法

浮点数加法步骤

- 求阶差
- 对阶
- 尾数相加
- 规格化并判溢出
- 舍入
- 置0

如果尾数比规定位数长，需考虑舍入



3.5.2 浮点数加法

浮点数加法步骤

- 求阶差
- 对阶
- 尾数相加
- 规格化并判溢出
- 舍入
- **置0**

尾数为0说明结果也为0，根据IEEE754，阶码和尾数全为0



3.5.2 浮点数加法

IEEE 754标准的四种舍入方式

就近舍入

- 舍入为最近可表示的数

例：1.1101**11** ~ 1.1110; 1.1101**01** ~ 1.1101;
 1.1101**10** ~ 1.1110; 1.1111**00** ~ 1.1111

附加位为：

11：入

01：舍

10（强制结果为偶数）

00：保持结果不变



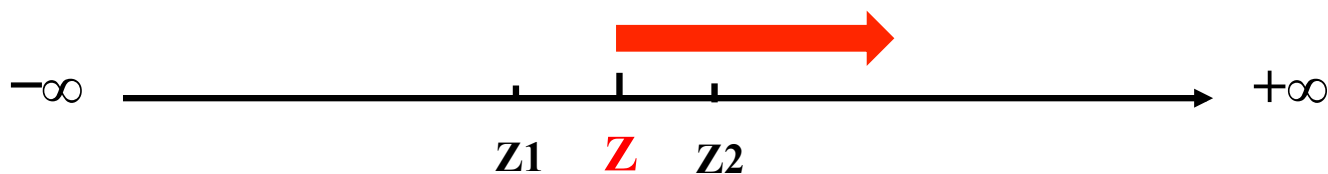
3.5.2 浮点数加法

IEEE 754标准的四种舍入方式

就近舍入

朝 $+\infty$ 方向舍入

- 舍入为Z2(正向舍入)



Z1和Z2分别是结果Z的最近可表示的左、右数



3.5.2 浮点数加法

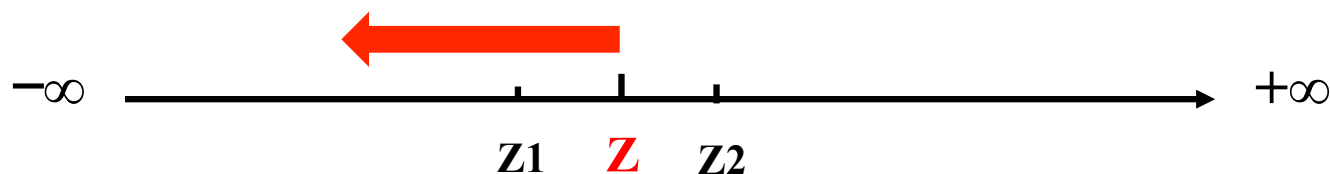
IEEE 754标准的四种舍入方式

就近舍入

朝 $+\infty$ 方向舍入

朝 $-\infty$ 方向舍入

- 舍入为Z1(负向舍入)



Z1和Z2分别是结果Z的最近可表示的左、右数



3.5.2 浮点数加法

IEEE 754标准的四种舍入方式

就近舍入

朝 $+\infty$ 方向舍入

朝 $-\infty$ 方向舍入

朝0方向舍入

- 总是舍入成Z1与Z2中绝对值较小的数 (更接近于0)



Z1和Z2分别是结果Z的最近可表示的左、右数



浮点数加法举例

用IEEE 754单精度形式，求出浮点数 $X=0.5_{10}$ 与 $Y=-0.4375_{10}$ 之和

解： $0.5_{10} = 1/2_{10} = 0.1_2 = (1.00...0_2) \times 2^{-1}$

$-0.4375_{10} = -7/16_{10} = -0.0111_2 = -1.110_2 \times 2^{-2}$

$[0.5]_{\text{浮}} = 0 \text{ 01111110 000...0}, [-0.4375]_{\text{浮}} = 1 \text{ 01111101 110...0}$

↓
符号位

↓
阶码

↓
尾数

↓
符号位

↓
阶码

↓
尾数



浮点数加法举例

用IEEE 754单精度形式，求出浮点数 $X=0.5_{10}$ 与 $Y=-0.4375_{10}$ 之和

解： $0.5_{10} = 1/2_{10} = 0.1_2 = (1.00...0_2) \times 2^{-1}$

$$-0.4375_{10} = -7/16_{10} = -0.0111_2 = -1.110_2 \times 2^{-2}$$

$$[0.5]_{\text{浮}} = 0\ 01111110\ 000...0, [-0.4375]_{\text{浮}} = 1\ 01111101\ 110...0$$

对阶（求阶差）：

$$[\Delta E]_{\text{补}} = [Ex]_{\text{移}} + [-[Ey]_{\text{移}}]_{\text{补}} \pmod{256}$$

$$[\Delta E]_{\text{补}} = 0111\ 1110 + 1000\ 0011 = 0000\ 0001, \Delta E = 1$$

所以：对Y进行对阶， $[Y]_{\text{浮}} = 1\ \underline{0111\ 1110}\ \underline{1110...0}$



浮点数加法举例

用IEEE 754单精度形式，求出浮点数 $X=0.5_{10}$ 与 $Y=-0.4375_{10}$ 之和

解： $0.5_{10} = 1/2_{10} = 0.1_2 = (1.00...0_2) \times 2^{-1}$

$-0.4375_{10} = -7/16_{10} = -0.0111_2 = -1.110_2 \times 2^{-2}$

$[0.5]_{\text{浮}} = 0\ 01111110\ 000...0$, $[-0.4375]_{\text{浮}} = 1\ 01111101\ 110...0$

对阶（求阶差）： $[Y]_{\text{浮}} = 1\ \underline{0111\ 1110}\ \underline{1110...0}$

尾数相加： $01.0000...0 + (10.1110...0) = 00.00100...0$

两个异号的原码数加法
(尾数为原码表示，最左边一位为符号位)



浮点数加法举例

用IEEE 754单精度形式，求出浮点数 $X=0.5_{10}$ 与 $Y=-0.4375_{10}$ 之和

解： $0.5_{10} = 1/2_{10} = 0.1_2 = (1.00...0_2) \times 2^{-1}$

$$-0.4375_{10} = -7/16_{10} = -0.0111_2 = -1.110_2 \times 2^{-2}$$

$$[0.5]_{\text{浮}} = 0\ 01111110\ 000...0, [-0.4375]_{\text{浮}} = 1\ 01111101\ 110...0$$

对阶（求阶差）： $[Y]_{\text{浮}} = 1\ \underline{0111\ 1110}\ \underline{1110...0}$

尾数相加： $01.0000...0 + (10.1110...0) = 00.00100...0$

规格化和判溢出： $+(0.00100...0)_2 \times 2^{-1} = +(1.00...0)_2 \times 2^{-4}$ (阶码减3)

因为 $127 \geq -4 \geq -126$ ，没有溢出



浮点数加法举例

用IEEE 754单精度形式，求出浮点数 $X=0.5_{10}$ 与 $Y=-0.4375_{10}$ 之和

解： $0.5_{10} = 1/2_{10} = 0.1_2 = (1.00...0_2) \times 2^{-1}$

$-0.4375_{10} = -7/16_{10} = -0.0111_2 = -1.110_2 \times 2^{-2}$

$[0.5]_{\text{浮}} = 0\ 01111110\ 000...0$, $[-0.4375]_{\text{浮}} = 1\ 01111101\ 110...0$

对阶（求阶差）： $[Y]_{\text{浮}} = 1\ \underline{0111\ 1110}\ \underline{1110...0}$

尾数相加： $01.0000...0 + (10.1110...0) = 00.00100...0$

规格化和判溢出： $+(0.00100...0)_2 \times 2^{-1} = +(1.00...0)_2 \times 2^{-4}$ (阶码减3)

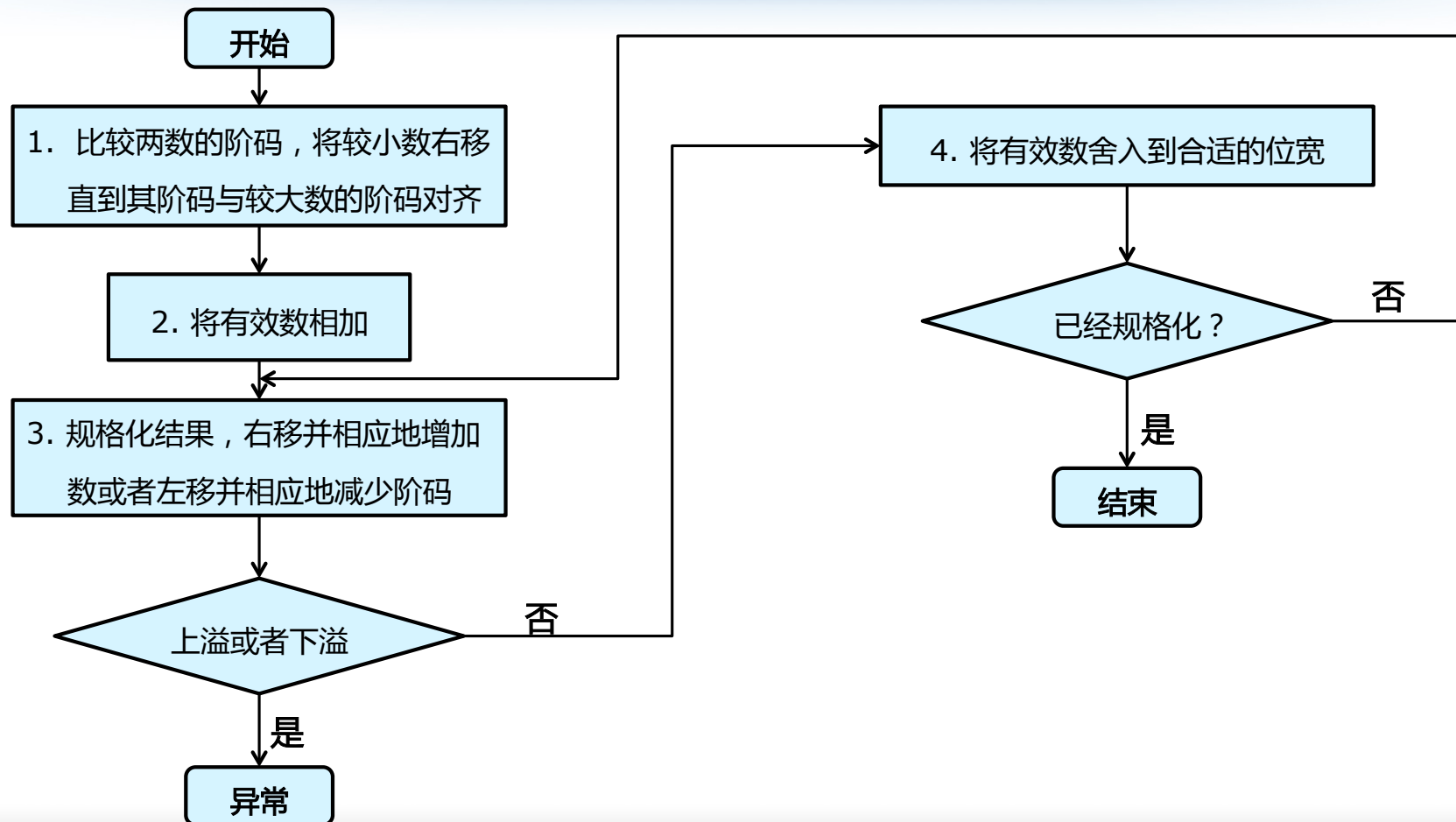
因为 $127 \geq -4 \geq -126$ ，没有溢出

舍入：无需舍入

结果为： $(1.00...0)_2 \times 2^{-4}$ 即 $1/16 = 0.0625$

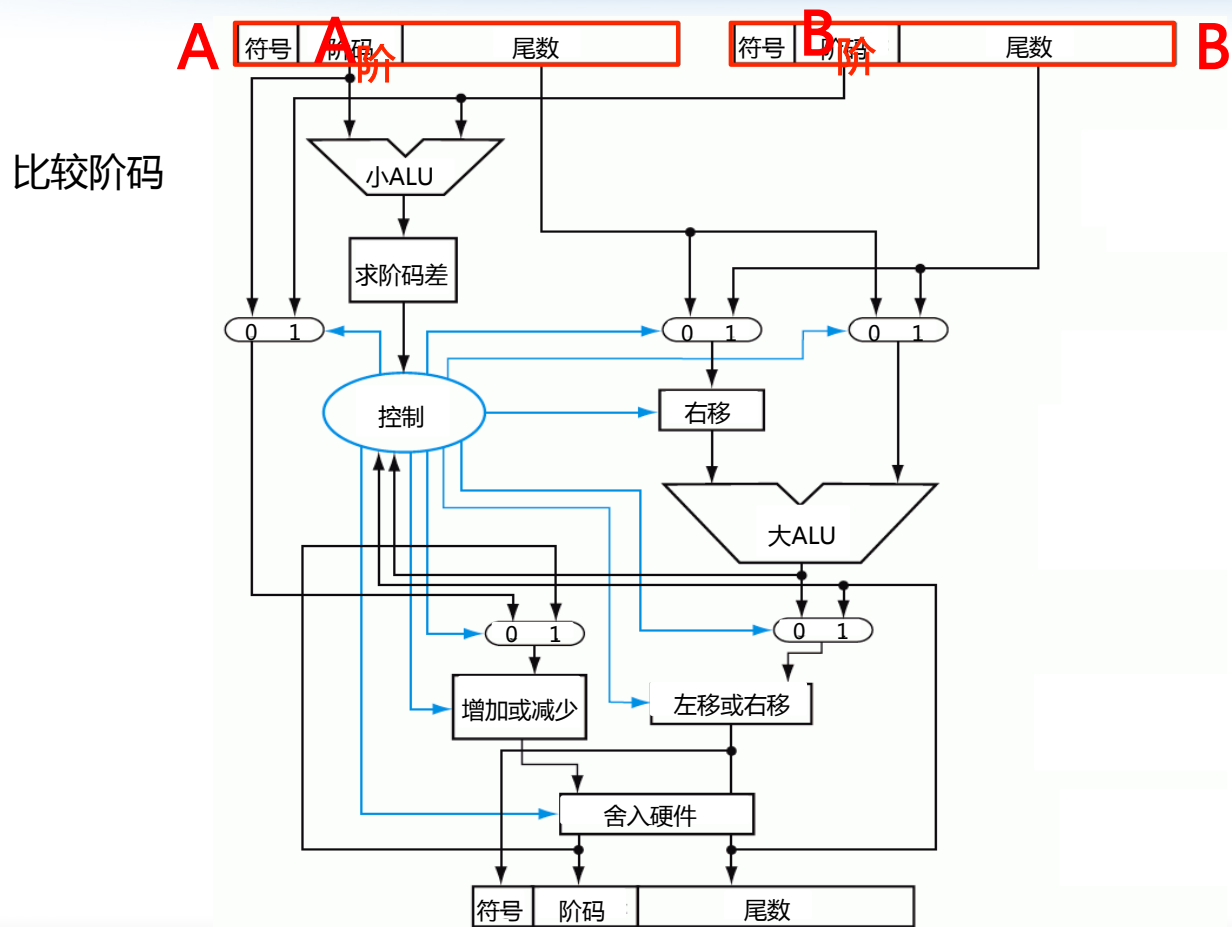


浮点数加法算法



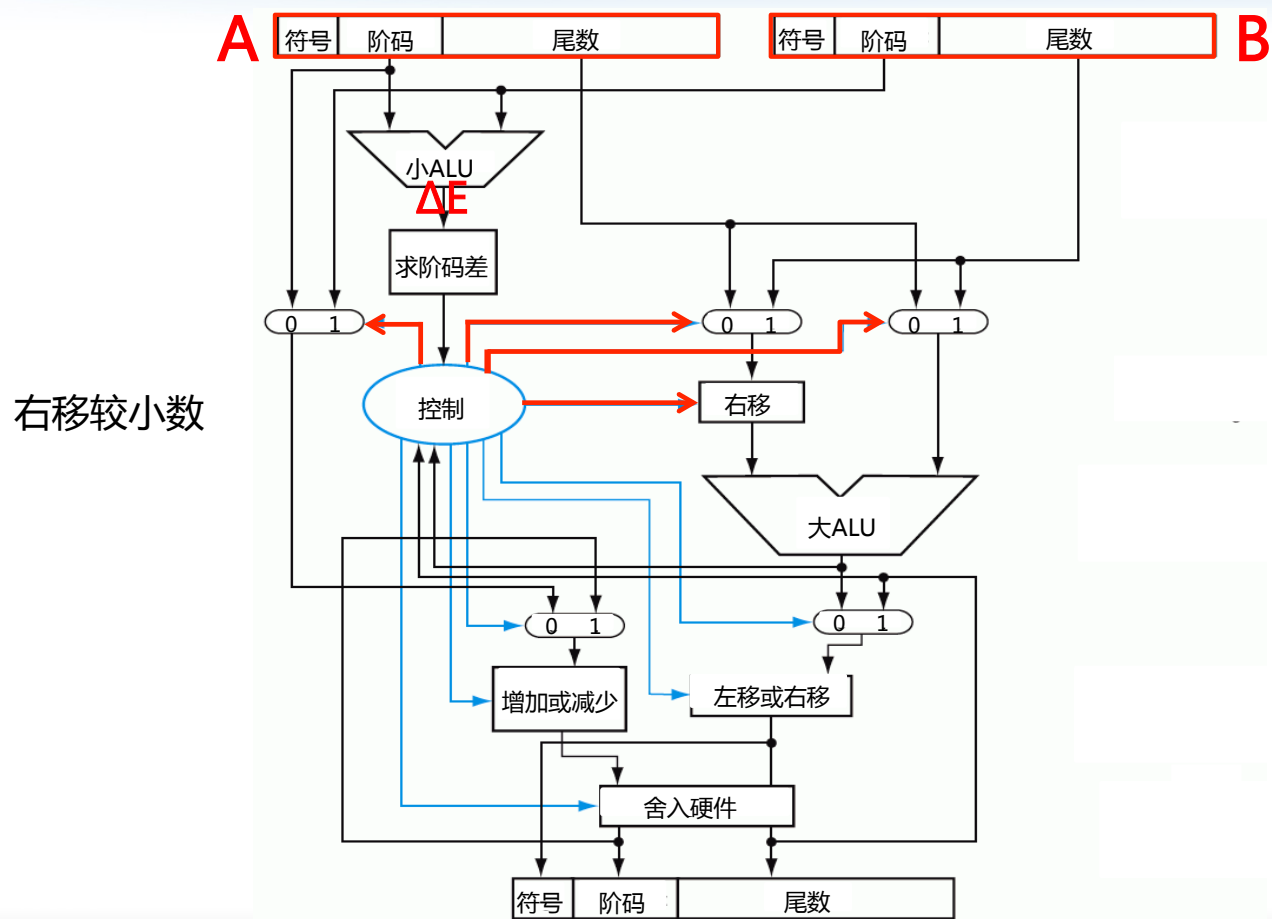


浮点数加法运算的硬件逻辑结构图



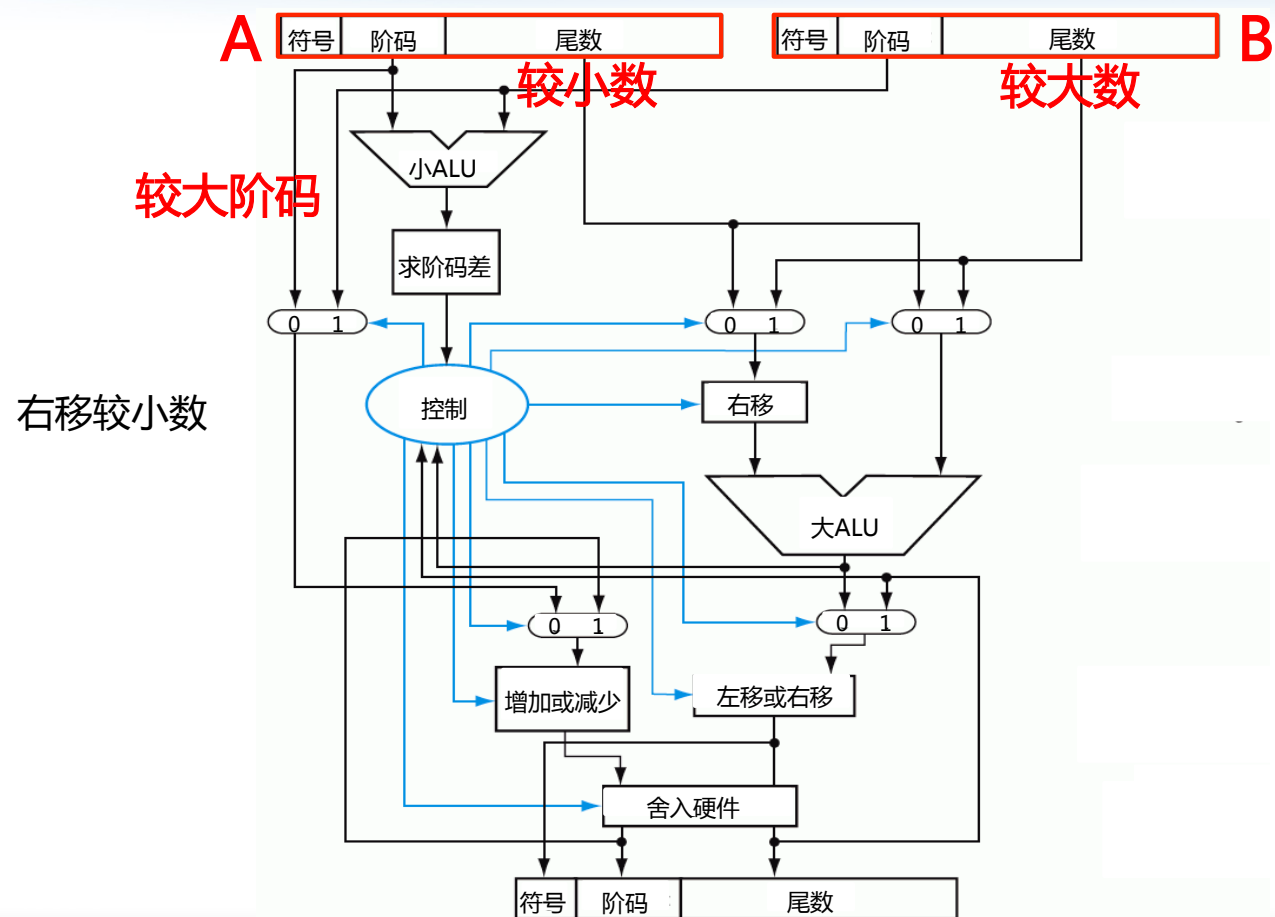


浮点数加法运算的硬件逻辑结构图



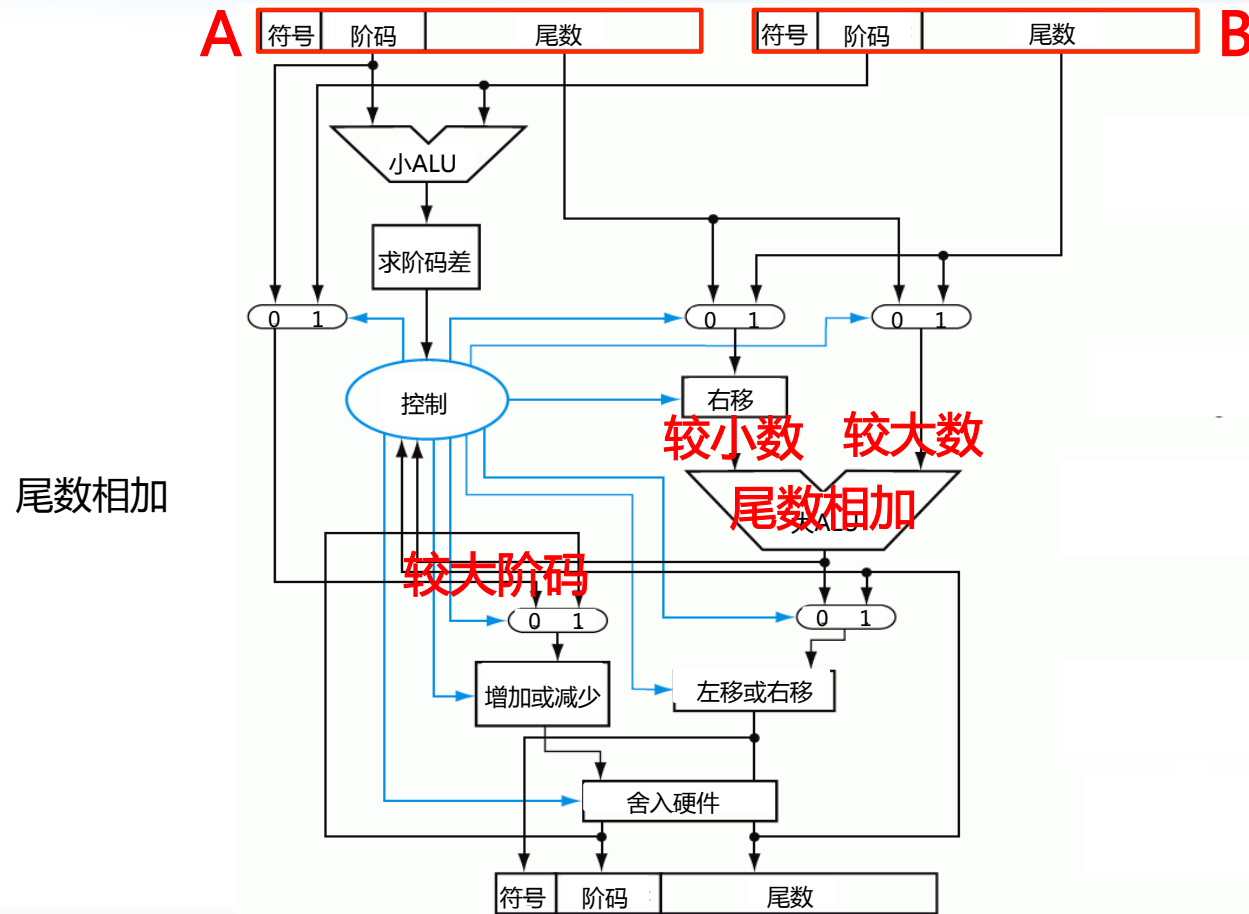


浮点数加法运算的硬件逻辑结构图



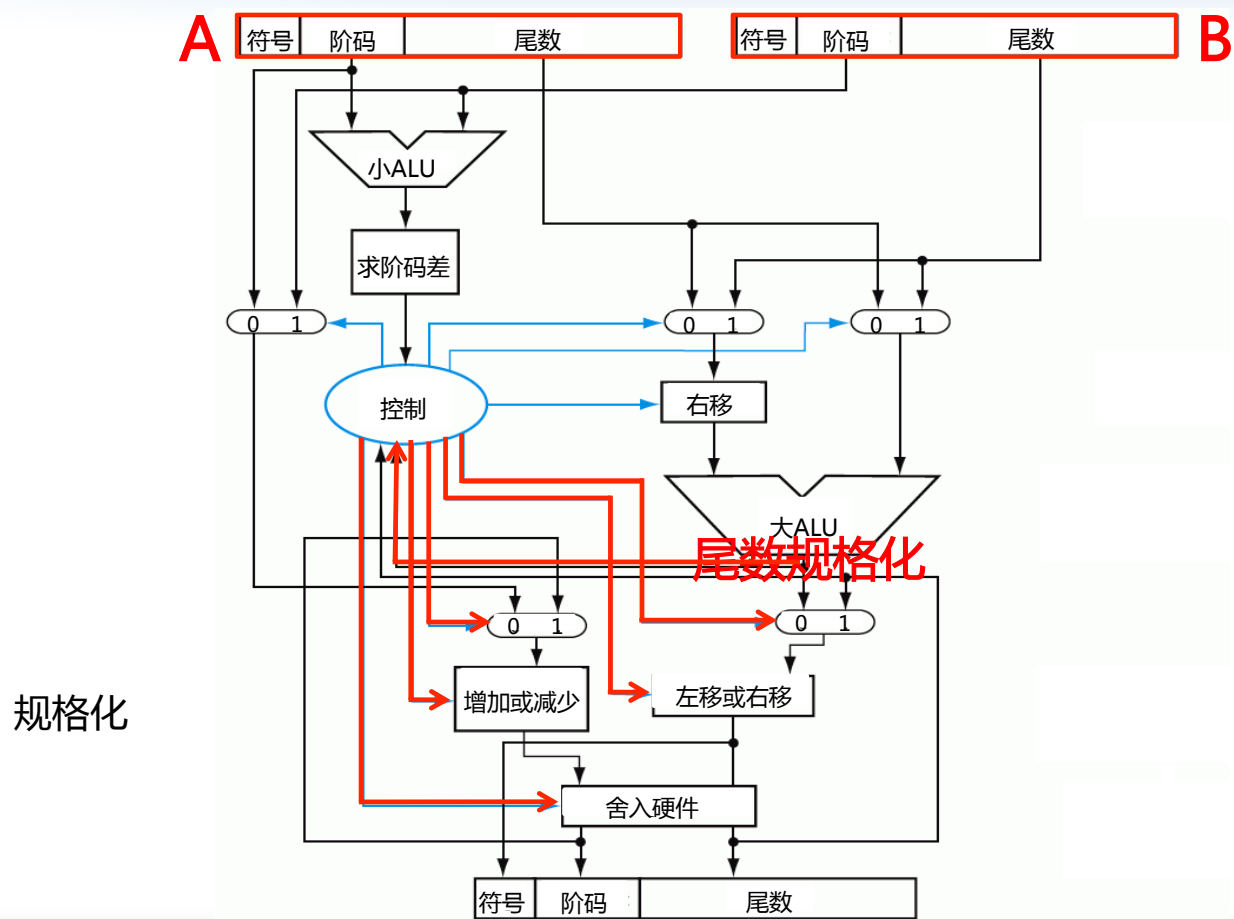


浮点数加法运算的硬件逻辑结构图



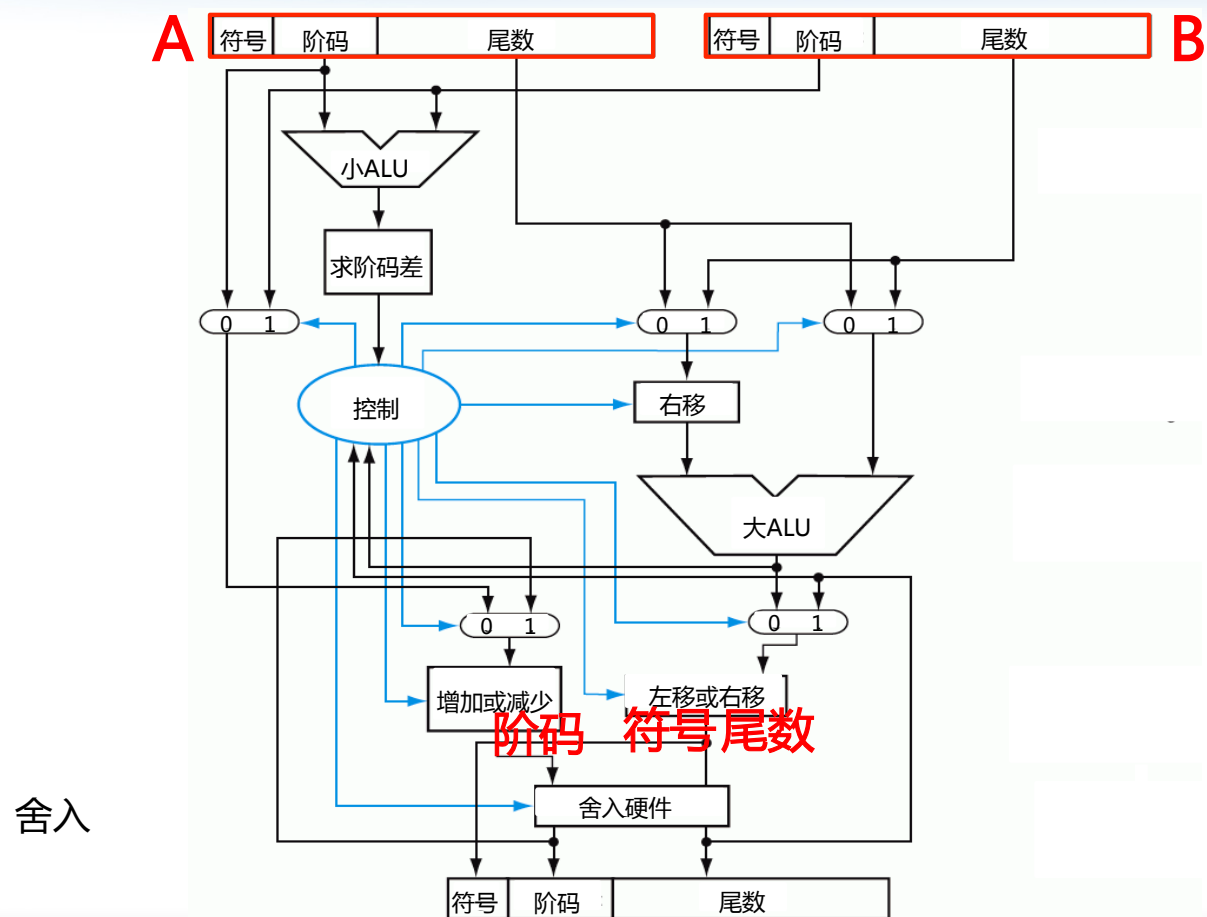


浮点数加法运算的硬件逻辑结构图





浮点数加法运算的硬件逻辑结构图





浮点数加法运算的硬件逻辑结构图

