

第 1 章 概述

1.1 记法

本书既不是普通的数学算式教程，也不是单纯的电脑程序算法手册，它讲的是“计算机算术”（computer arithmetic），参与运算的数是长度固定的位串与位元数组^①。计算机算术中的表达式与普通的数学表达式近似，不同之处在于其中的变量指代的是 CPU 寄存器中的内容，而且计算机算术表达式的值是一串不具备特定符号性的位元。在此类表达式中，操作符可能会以不同方式解读其操作数，例如“比较操作符”（comparison operator）有时会将其操作数当成有正负号的二进制整数，有时却把它视为无符号的二进制整数。为免混淆，本书所列计算机算式以不同的符号来区分上述两种情况。

计算机算式与数学算式的主要差别在于：无论是加减法还是乘法，计算机算式的结果总是要跟 2 的 n 次方取模， n 是指当前机器的字长^②。此外，计算机算式的运算种类也远多于数学算式：除了基本的四则运算外，还有逻辑与、异或、比较、左移，等等。

如未特别指明，则默认字长为 32 位，且带符号的整数均以“2 补码”^③形式表示。

-
- ① bit vector，即 bit array，也写作 bitmap、bitset、bit string 等，是由若干位元排列而成的数组，又叫“位向量”、“位矢量”等，详情参见：https://en.wikipedia.org/wiki/Bit_array。译文在不致混淆时，均以“位元数组”称之。——译者注
- ② word size，就是字组中所含的位元个数，也称 word length、word width。在不致混淆的情况下，“字长”、“字宽”均指这一概念。——译者注
- ③ two's complement，也叫“2 的补码”、“二补数”，是一种用二进制表示带符号数字的方式。整数和零的补码表示法与其二进制写法相同，只是左方要补足 0，而表示负数，则需先将其绝对值按位取反，也就是求绝对值的一补数（也称反码），然后再加 1。例如用 8 位二进制表示数字时，5 可以表示为 0000 0101，而表示 -5，则需先求其绝对值 5，写成二进制形式 0000 0101，然后对其按位取反得到 1111 1010，再加 1。最终的结果 1111 1011 就是 -5 的补码表示。详情参见：<https://zh.wikipedia.org/wiki/补码>。带符号整数与无符号整数的英文分别为 signed integer 与 unsigned integer，其中“带符号”与“无符号”是计算机处理二进制数的两种不同方式。前者会根据最高有效位来判定数字的正负，0 为非负，1 为负，而后者则一律将其视为非负数。以上面的 1111 1011 来说，若视为带符号整数，则其值为 -5（因最高位为 1，故是负数。先将其按位取反得 0000 0100，再加 1 得 0000 0101，即十进制的 5，再添上负号得 -5），若视为无符号整数，则其值为 251（最高位的 1 不再表示负数，而当做 128 来算）。详情参见：<https://zh.wikipedia.org/wiki/有符号数处理>。——译者注

计算机算式与数学算式的书写方式相同，只是其中指代 CPU 寄存器中内容的变量会以粗体标出。为了和位元数组运算的通则一致，我们把电脑中的字组也视为由一串位元构成的数组。若某常量表示 CPU 寄存器中的值，那么它也会以粗体字出现。（这种情况在位元数组运算中找不到对应物，如果要在位元数组算式里书写常量，只能逐个列出其中的每个位元。）然而常量如果作为 shift 等指令的操作数，则不加粗。

对于“+”这样的操作符，若其操作数为粗体，则表明执行的是计算机加法，也就是“向量加法”，否则意味着执行的是纯数字加法。如果操作数未加粗，则其对应的操作符就是纯数学运算中的含义。要是我们想拿原来做计算机运算的粗体 x 值做数学运算的话，那就会把它写成不加粗的 x ，其符号性应该能够从上下文中推出。假如 $x = 0x8000\ 0000$ ， $y = 0x8000\ 0000$ ，那么在做带符号的整数运算时， $x = y = -2^{31}$ ， $x + y = -2^{32}$ ，而 $x + y = 0^{\ominus}$ 。此处的 $0x8000\ 0000$ 是用十六进制表示的位串，它最左边的位元是“1”，后面跟着 31 个“0”。

位元的序号从右侧算起，最右方的位元（也就是最低有效位）叫做 0 号位元。术语“位”、“半字节”、“字节”、“半字组”、“字组”、“双字”^㉑所对应的位元数量分别是 1、4、8、16、32、64。

简短的代码段用的都是计算机算式，并以左箭头表示赋值操作，偶尔还会用 if 语句。在这种情况下，它只是以一种与电脑平台无关的方式编写汇编语言代码罢了。

过长或过于复杂的计算机算式则用 C 语言来写，其代码遵循 ISO 1999 标准^㉒。

完整描述 C 语言不是本书该做的事，不过书中用到的大部分 C 语言基本表达式 [H & S] 都总结到表 1.1 中了，用过程序语言编程但是不熟悉 C 语言的读者应该看看。表中也列出了笔者在计算机算式中用到的对应操作符，它们按照优先级从高到低的顺序排列（高者先算），在优先级这一列中，L 表示左结合，比如乘号就是左结合的运算符： $a \cdot b \cdot c = (a \cdot b) \cdot c$ ，而 R 则表示右结合^㉓。本书计算机算式里的运算符，其优先级与结合性同 C 语言一致。

表 1.1 C 语言与计算机算术表达式对照表

优先级	C	计算机算式	含义
	0x...	0x..., 0b...	十六进制常数，二进制常数

⊖ 不加粗的 $x + y$ 是普通的数学运算，故结果为 -2^{32} ，而粗体的 $x + y$ 则是按照计算机算术规则做二进制加法，其结果为 1 00000000 00000000 00000000 00000000，也就是 1 后面 32 个零。宽度为 32 的字组无法容纳这个 33 位元的数，故而最高有效位的“1”会被舍去，留下 32 个 0，因此最终的运算结果就是 0。——译者注

㉑ 对应英文写法分别为：bit、nibble、byte、halfword、word、doubleword。——译者注

㉒ 该标准的正式名称是 ISO/IEC 9899: 1999，俗称 C99，详情参见：<https://zh.wikipedia.org/wiki/C语言#C99>。C 语言的最新标准是 2011 年发布的 ISO/IEC 9899: 2011，俗称 C11。——译者注

㉓ 左结合与右结合的英文分别是 left-associative 与 right-associative，结合性指的是遇到两个相邻的同优先级运算符时，要从左先算还是从右先算。详情参见：https://en.wikipedia.org/wiki/Operator_associativity。——译者注

(续)

优先级	C	计算机算式	含义
16	a [k]		数组 a 中索引为 k 的元素
16		x_0, x_1, \dots	若干个变量或位元 (具体含义会在文中说明)
16	f (x, ...)	$f(x, \dots)$	求函数值
16		abs (x)	求绝对值 (例外: $\text{abs}(-2^{31}) = -2^{31}$)
16		nabs (x)	绝对值对应的负数
15	x++, x--		后置自增与自减 ^①
14	++x, --x		前置自增与自减 ^②
14	(数据类型名) x		类型转换
14R		x^k	x 的 k 次方
14	~x	$\neg x, \bar{x}$	按位取反 (也就是求 x 的反码) ^③
14	! x		逻辑非 (若 x 是 0 则结果为 1, 若 x 非 0 则结果为 0)
14	- x	-x	取相反数
13L	x * y	$x * y$	乘法, 根据字组长度裁剪其运算结果 ^④
13L	x/y	$x \div y$	带符号整数的除法
13L	x/y	$x \overset{u}{\div} y$	无符号整数的除法
13L	x % y	rem (x, y)	已知 x 与 y 为带符号的数, 求 $x \div y$ 的余数。结果有可能是负数 ^⑤
13L	x % y	remu (x, y)	已知 x 与 y 为无符号的数, 求 $x \div y$ 的余数
		mod (x, y)	已知 x 与 y 为带符号的数, 求 x 除以 y 的余数, 并将结果调整到 $[0, \text{abs}(y)-1]$ 之间
12L	x + y, x - y	$x + y, x - y$	加法与减法
11L	x << y, x >> y	$x \ll y, x \gg y$	左移位, 右移位 (以 0 填补空缺位元, 又名逻辑移位, logical shift)
11L	x >> y	$x \overset{s}{\gg} y$	右移位 (根据 x 的符号来填补空缺位元, 又名算术移位或数学移位)
11L		$x \overset{ml}{\ll} y, x \overset{mr}{\gg} y$	循环左移, 循环右移
10L	$x < y, x < = y,$ $x > y, x > = y$	$x < y, x \leq y,$ $x > y, x \geq y$	带符号数的关系比较表达式
10L	$x < y, x < = y,$ $x > y, x > = y$	$x \overset{u}{<} y, x \overset{u}{\leq} y,$ $x \overset{u}{>} y, x \overset{u}{\geq} y$	无符号数的关系比较表达式
9L	$x = y, x ! = y$	$x = y, x \neq y$	比较两数是否相等, 比较两数是否不等
8L	x & y	$x \& y$	按位与
7L	x ^ y	$x \oplus y$	按位异或
7L		$x \equiv y$	按位等值, 结果与 $\neg(x \oplus y)$ 相同
6L	x y	$x y$	按位或
5L	x && y	$x \bar{\&} y$	条件与 (x 与 y 均不为 0 时, 结果是 1, 否则是 0)
4L	x y	$x \bar{ } y$	条件或 (x 与 y 均为 0 时, 结果是 0, 否则是 1)

(续)

优先级	C	计算机算式	含义
3L		$x \parallel y$	连接
2R	$x = y$	$x \leftarrow y$	赋值

- ① postincrement 与 postdecrement，意为先求表达式的值，然后再对其增减。假设 x 是 5，那么 $x++$ 的值还是 5，求完值之后， x 才变成 6。——译者注
- ② preincrement 与 predecrement，意为先对其增减，然后再求表达式的值。假设 x 是 5，那么先将 x 加 1，然后再判定 $++x$ 的值为 6。——译者注
- ③ one's-complement，又叫一补数，是将二进制数中的每个位元反转之后得到的数。例如 10010 的一补数为 01101。详情参见：<https://zh.wikipedia.org/wiki/一补数>。——译者注
- ④ modulo word size，意思是如果运算结果位数太多，超过了字组长度，那么会丢弃超出的部分。此概念直译为“与字长取模”或“模字长”，为了不使人误认为是和 32（字组长数值）取模，故译文用了意译，下文皆同。——译者注
- ⑤ 在 C99 标准中，模除的结果与被除数的正负号相同，详情参见：https://en.wikipedia.org/wiki/Modulo_operation。——译者注

除了表 1.1 列出的写法之外，书中也将借用一些布尔代数运算符与标准数学符号，在用到它们时会给出解释。

计算机算术中除了“abs”、“rem”之外，还会用到其他一些函数，等讲到那些函数时再给出其定义。

C 语言中， $x < y < z$ 这个表达式的意思是：先判断 $x < y$ 是否成立，若成立，则这个子表达式的值就是 1，否则为 0，然后，再将刚才的结果与 z 比较[⊖]。而在本书所讲的计算机算术中，该表达式的意思就是想判断 $x < y$ 与 $y < z$ 是否同时成立而已。

C 语言中有三种循环控制语句：while、do、for。while 语句的语法[⊖]是：

$$\text{while}(\text{expression})\text{statement}$$

首先评估 expression 的值，若为真（也就是非 0），那么就执行 statement 对应的语句，然后再次判断 expression。一旦 expression 为假（也就是 0），while 循环就终止了。

do 语句与之相似，只是判断放在了循环尾部。它的语法是：

$$\text{do statement while}(\text{expression})$$

先执行 statement 中的语句，然后再判断 expression。若为真，则继续执行循环，若为假，则循环终止。

for 语句的格式为：

$$\text{for}(e_1; e_2; e_3)\text{statement}$$

首先执行 e_1 ，这部分通常是赋值表达式。然后判断 e_2 ，它一般是个关系表达式，如果

⊖ 例如 C 语言中表达式 $-5 < -4 < -3$ 就不成立，因为必须先求 $-5 < -4$ 这个子表达式的值，它成立，故被视为 1，然后再判断 1 是否小于 -3，此判断不成立，故整个表达式不成立。——译者注

⊖ 为了遵循业内惯例，语法格式中的英文单词在译文中保留原样。其中的 expression 意为“表达式”，statement 意为“语句”。——译者注

为假，则终止 for 循环，如果为真，那就执行 statement。最后再执行 e_3 ，这部分通常也是个赋值表达式。执行完毕后，又跳回 e_2 判断。因此，“do $i = 1$ to n ”用 for 语句写出来就是：

```
for (i = 1; i <= n; i++)
```

(书中用到后置自增运算符的场合不多，这是一例。)

ISO C 标准并未规定右移 (“ \gg ”操作符)带符号的数时，左边多出来的空位是用 0 填充还是用符号填充^①。在本书的 C 语言代码中，我们假定如果左操作数是带符号的，那么右移操作就使用符号填充，若无符号，则按 ISO 规范以 0 填充。大多数 C 编译器也都这么做。

左移操作符都采用“逻辑”填充。(某些电脑还有“算术”左移操作，也就是左移之后保持符号不变。这通常出现在老式计算机上。)

移位操作还会出现一个问题，那就是 ISO C 标准规定：假如移位的数量等于或大于左操作数的位宽，则结果未定义。虽说如此，但是几乎所有 32 位机在遇到这种情况时，都会把移位数量和 32 或 64 取模^②。书中代码在对待此问题时，上述各种处理方式都会用到。若不同方式之间的差异关乎运算结果，则会给出解释。

1.2 指令集与执行时间模型

我们假定测试书中算法时所用的这台电脑，其 CPU 使用当前通行的 RISC 指令集，这样就好粗略估量其效率了。IBM RS/6000 系列电脑^③的 CPU、Oracle SPARC[®]及 ARM 架构[®]的处理器都用的是 RISC。这台虚拟电脑的 CPU 使用“三段式地址”表示操作数，其寄存器相当多，至少 16 个。除非另有说明，否则寄存器都是 32 位的。0 号通用寄存器的值永远是 0，其他寄存器作用不限。

有些 CPU 中具备“特殊用途”寄存器^④，专门保存条件比较的结果或“溢出”等状

-
- ① 用 0 填充 (0-propagating) 就是表 1.1 中所说的“逻辑”移位，而用符号填充 (sign-propagating) 的意思是，如果待移位的数是负值，那么就用 1 来填充，这样移位后的结果还是负的；而如果待移位的数非负，则用 0 填充，移位后的结果仍然非负。这种保持正负性的填充方式也就是上表所说的“算术”移位或“数学”移位。——译者注
- ② 举例来说，若要将一个 32 位元的数左移 80 位，那么由于 80 大于等于 32，所以必须先求它除以 32 的余数，也就是 16，然后将该数左移 16 位，得出运算结果。——译者注
- ③ RS6000 是 IBM 公司使用其 RISC 架构的 Power 处理器设计生产的小型计算机。详情参见：<https://zh.wikipedia.org/wiki/RS/6000>。——译者注
- ④ SPARC，全称为“可扩充处理器架构”(Scalable Processor ARChitecture)，是 RISC 微处理器架构之一。详情参见：<https://zh.wikipedia.org/wiki/SPARC>。——译者注
- ⑤ ARM 架构是一种 32 位元精简指令集 (RISC) 处理器架构，因其具备低成本、高效能、低功耗的特性，故广泛用于嵌入式系统。详情参见：https://zh.wikipedia.org/wiki/ARM_架构。——译者注
- ⑥ special purpose register，又叫 special function register 或 special register，详情参见：https://en.wikipedia.org/wiki/Special_function_register。——译者注

态码之类的数据，而我们为了简洁起见，假设虚拟电脑的 CPU 中没有这种寄存器。此外，它还没有浮点指令。浮点数在书中只占少量篇幅，基本上局限于第 17 章。

现在规定两套 RISC 指令集：表 1.2 是“基本 RISC 指令集”，这些指令再加上表 1.3 中的那些，就构成“完整 RISC 指令集”。

表 1.2 基本 RISC 指令集

操作码助记符	操作数	含义
add, sub, mul, div, divu, rem, remu	RT, RA, RB	$RT \leftarrow RA \text{ op } RB$, 其中 op 可以是加法、减法、乘法、带符号数的除法、无符号数的除法、带符号数的求余、无符号数的求余
addi, muli	RT, RA, I	$RT \leftarrow RA \text{ op } I$, 其中 op 可以是加法或乘法, I 是 16 位元的带符号常量 ^①
addis	RT, RA, I	$RT \leftarrow RA + (I \ll 16)$
and, or, xor	RT, RA, RB	$RT \leftarrow RA \text{ op } RB$, 其中 op 可以是按位与、按位或、按位异或
andi, ori, xori	RT, RA, Iu	同上, 但最后一个操作数是 16 位元的无符号常量
beq, bne, blt, ble, bgt, bge	RT, target	在条件成立时转入 target 分支。这 6 个操作符的判断条件分别是: $RT=0, RT \neq 0, RT < 0, RT \leq 0, RT > 0, RT \geq 0$
bt, bf	RT, target	当 RT 是 true/false 时转入 target, 与 bne/beq 等效
cmpeq, cmpne, cmplt, cmple, cmpgt, cmpge, cmpltu, cmpleu, cmpgtu, cmpgeu	RT, RA, RB	将 RA 与 RB 比较的结果存入 RT。若判断不成立, 则 RT 为 0, 若成立, 为 1。cmpeq 表示 compare for equality (判断二者是否相等), cmpne 表示 compare for inequality (判断两者是否不等), cmplt 表示 compare for less than (判断前者是否小于后者), 以此类推。cmp 后面的字母与分支指令中的含义相同。后缀“u”表示无符号数的比较
cmpieq, cmpine, cmpilt, cmpile, cmpigt, cmpige	RT, RA, I	与 cmpeq 等 6 个操作符含义相同, 只是 I 为 16 位元的带符号常量
cmpiequ, cmpineu, cmpiltu, cmpileu, cmpigtu, cmpigeu	RT, RA, Iu	与 cmpltu 等 6 个操作符含义相同, 只是 I 为 16 位元的无符号常量
ldbu, ldh, ldhu, ldw	RT, d (RA)	分别将地址 $RA + d$ 中的无符号字节、半字组、无符号半字组、字组载入 RT。d 为 16 位元的带符号常量
mulhs, mulhu	RT, RA, RB	将 RA 与 RB 之积的高 32 位存入 RT。两指令分别适用于带符号数及无符号数
not	RT, RA	将 RA 按位取反后的值存入 RT
shl, shr, shrs	RT, RA, RB	将 RA 左移位或右移位后的值存入 RT。RB 最右方的 6 个位元代表移位数。shl 与 shr 用 0 填充, shrs 用 RA 的符号位填充 (移位数只取最右方 6 个位元的原因是, 它需要和 64 取模, 以便调整到 0 至 63 之间)
shli, shri, shrsi	RT, RA, Iu	将 RA 左移位或右移位后的值存入 RT。移位数是 Iu 所代表的 5 位元常数
stb, sth, stw	RS, d (RA)	将 RS 里的字节、半字、字组存入 $RA + d$ 所表示的内存地址中。d 是 16 位元的带符号常量

^① 本书多次出现 immediate value 一词, 字面意思为“立即值”、“立即数”, 它是谈到计算机指令时的一种术语, 在不影响文意的情况下, 译文均以“常量”、“常数”称呼之。——译者注

表 1.2、表 1.3 与表 1.4 中的 RA 与 RB，如果作为源操作数使用，就表示这些寄存器本身的内容。^①

在实际使用的 CPU 中，尚有分支链接跳转指令（branch and link）以及返回到寄存器中所含地址的分支返回指令，前者可用来调用子程序，后者则用于从子程序中返回或实现“switch”功能。此外，可能存在一些处理专用寄存器的指令，特权指令与调用管理服务的指令当然也不少，或许还会有浮点数指令。

表 1.3 列出了 RISC 指令集中可能会用到的其他一些计算类指令，后面的章节中将会讲到它们。

汇编语言中提供了一些“扩展助记符”（Extended Mnemonic），开发人员用起来很方便，它们有点像宏，展开之后通常是一条指令。表 1.4 列举了一些可能会用到的扩展助记符。

常量加载操作会根据常量 I 的值而展开成一条或两条指令。如果 $0 \leq I < 2^{16}$ ，那么就展开成 R0 与 I 的常量按位或操作（ori 指令），而当 $-2^{15} \leq I < 0$ 时，则会展开成 R0 与 I 的常量加法操作（addi 指令）。若 I 最右方的 16 个位元都是 0，那么就视为常量移位加法操作（addis 指令），而在不属于上述三种情况时，则会展开成一条 addis 与一条 ori 指令。^②（在最后一情况下，也可以直接使用从内存中加载数值的指令来做。但是为了便于估量算法所需的执行时间和空间，我们假定此种常量加载操作总是相当于两个算术指令。）

表 1.3 “完整 RISC 指令集”中的其余附加指令

操作码助记符	操作数	含 义
abs, nabs	RT, RA	将 RA 的绝对值或“负绝对值” ^① 存入 RT
andc, eqv, nand, nor, orc	RT, RA, RB	将 RB 按位取反后与 RA 按位与，按位比较，按位与后按位取反，按位或后按位取反，将 RB 按位取反后与 RA 按位或
extr	RT, RA, I, L	将 RA 中序号为 I 至 I+L-1 之间的位元提取出来，靠右存放在 RT 中，余位填 0
extrs	RT, RA, I, L	与 extr 相同，但是用符号填充
ins	RT, RA, I, L	将 RA 中序号为 0 至 L-1 的位元插入 RT 中序号为 I 至 I+L-1 的位置上
nlz	RT, RA	取得 RA 中前导 0 的个数（该操作的结果在 0 至 32 之间）
pop	RT, RA	RT 计算 RA 中有多少个值为 1 的位元（该操作的结果在 0 至 32 之间）
ldb	RT, d (RA)	从 RA+d 所在的内存地址处加载一个带符号的字节到 RT 中。d 为 16 位元的带符号常量
moveq, movne, movlt, movle, movgt, movge	RT, RA, RB	如果 RA 与 0 的关系满足判断条件，就将 RB 赋值给 RT，否则 RT 的值不变。eq 与 ne 分别表示 RA=0、RA≠0，其余依此类推
shlr, shrr	RT, RA, RB	将 RA 循环左移位或循环右移位。RB 最右方的 5 个位元表示移位的数量

① 若 RA、RB 不以带括号形式出现，即表示其本身值，否则视为定位所用的内存地址。——译者注

② 常量加载操作要分成三种情况展开的原因，请参考：<http://www.engr.uconn.edu/~jeffm/Classes/CSE240-Spring-2000/Lectures/lecture6/node4.html>。——译者注

(续)

操作码助记符	操作数	含义
shlri, shrrri	RT, RA, Iu	将 RA 循环左移位或者循环右移位，移位的数量在 5 位立即数字段中给出
trpeq, trpne, trplt, trple, trpgt, trpge, trpltu, trpleu, trpgtu, trpgeu	RA, RB	当两个操作数满足 $RA = RB$ 、 $RA \neq RB$ 等条件时，令处理器中断 (trap)
trpieq, trpine, trpilt, trpile, trpigt, trpige	RA, I	与 trpeq 等指令含义相同，只是第 2 个操作数为 16 位元的带符号常量
trpiequ, trpineu, trpiltu, trpileu, trpigt, trpigeu	RA, Iu	与 trpltu 等指令含义相同，只是第 2 个操作数为 16 位元的无符号常量

① negative of the absolute value 是一种特有的绝对值解读方式，详情参见：http://pic.dhe.ibm.com/infocenter/aix/v7r1/index.jsp?topic=%2Fcom.ibm.aix.aixassem%2Fdoc%2Falangref%2Fidalangref_nabs_negabs_instrs.htm。——译者注

表 1.4 扩展助记符

扩展助记符	展开式	含义
b target	beq R0, target	无条件执行分支
li RT, I	详见书中解释	将常量加载至 RT 中。 $-2^{31} \leq I < 2^{32}$
mov RT, RA	ori RT, RA, 0	将寄存器 RA 中的值赋给 RT
neg RT, RA	sub RT, R0, RA	将 RA 的相反数（也就是二补码）放入 RT
subi RT, RA, I	addi RT, RA, -I	将 RA 减 I 的差放入 RT ($I \neq -2^{15}$) ^①

① 因为 -2^{15} 的相反数是 2^{15} ，该值无法以 16 个位元容纳，所以在这种特殊情况下，subi 与 addi 展开式不等效。——译者注

一个指令到底应该属于基本 RISC 指令集还是完整 RISC 指令集，其实取决于个人的解读。有人可能会觉得无符号数除法与求余操作应该放在完整 RISC 指令集中，而与此相反，带符号的字节加载 (load byte signed, ldb) 操作反而应该归入基本 RISC 指令集。ldb 操作之所以会被放在完整指令集中，是因为其使用频率相当低，而且该指令需要将符号位扩展很多次^②，这样做比较耗费 CPU 周期。

划分基本 RISC 指令集与完整 RISC 指令集时，还有许多可以商榷的地方，此处不再赘言。

书中所用指令最多只会用到两个源寄存器与一个目标寄存器，这样做也简化了计算机的 CPU 架构（比方说，寄存器堆^③中最多只需两个读端口与一个写端口即可）。此外也缓解了优化编译^④的工作量，使它不用再处理那些用到多个目标寄存器的指令了。这样做

② 如果待加载字节的最高有效位是 1，则表明其为负数，将它加载到一个 32 位宽的寄存器时，必须把多出来的 24 个空位元全部扩展为 1，这样才能保证寄存器中的值和正负性与原字节相符。——译者注

③ Register File，是 CPU 中多个寄存器组成的阵列，通常由快速的静态随机读写存储器 (SRAM) 实现。这种 RAM 具有专门的读端口与写端口，可以多路并发访问不同的寄存器。详情参见：<https://zh.wikipedia.org/wiki/寄存器堆>。——译者注

④ Optimizing Compiler，是优化程序所用的编译器，可以加快程序执行速度或减低内存占用量等。详情参见：https://en.wikipedia.org/wiki/Optimizing_compiler。——译者注

的代价是，有些本来执行一条指令就能实现的操作现在得分解成两条。比如某些程序想同时知道两个数的商和余数（这种情况很罕见），那么现在就必须用除法和求余两个操作才能完成。在现实环境下的 CPU 中，余数是除法操作的副产品（by-product），很多电脑在执行除法指令时顺便计算好了余数。要获取保存两个字组乘积的双字时，也会遇到这个问题。

条件移动指令（conditional move，例如 `moveq`）表面上看只有两个源操作数，但换个角度看，也可以说是 3 个。因为该指令的结果取决于 RT、RA 与 RB 的值，所以乱序执行的 CPU 必须把这种指令里的 RT 标注为 use 和 set。也就是说，如果前面一条指令修改了 RT，而后面紧跟一个将要再度修改 RT 的条件移动指令，那么就只能按照这个顺序执行，同时不能丢弃前一条指令的运算结果。有鉴于此，此类处理器的设计者可以选择略过条件移动操作以避免处理这种（从逻辑上讲）需要 3 个源操作数的指令。反之，使用条件移动指令确实能简化分支的实现。

指令格式与本书内容无关，不过上面列出的所有 RISC 指令集，再加上浮点指令和一些管理用的指令，都可以在有 32 个通用寄存器的电脑上用 32 个位元（其中用于指示寄存器序号的字段占 5 个位元^①）加以实现。如果将 `compare`、`load`、`store`、`trap` 指令所支持的常数字段缩减为 14 位，那么这些指令在具备 64 个通用寄存器的电脑上，也照样能用 32 个位元做出来（其中用于指示寄存器序号的字段占 6 个位元^②）。

执行时间

本书假设所有指令都只需 1 个 CPU 周期即可执行完毕，但是乘法、除法及求余指令例外，我们不去预估这 3 个指令的执行时间。而分支指令不论其执行结果是转入分支还是继续沿主线执行，都只花 1 个周期。

常量加载指令会花费一至两个周期，因为要把常量存入寄存器中，所需的算术指令数可能是 1 个，也可能是两个。

书里用到加载与存储指令的场合不多，我们也假设其只需 1 个周期，而且忽略加载延迟（也就是从算术逻辑单元^③执行完加载指令算起，到待加载的数据真正可以用于后续指令为止，这两者的时间间隔）。

然而，只晓得算术与逻辑指令的执行周期，通常无法估算出程序执行时间，因为数

① 在 CPU 中，除了操作数之外，指令本身也是用二进制来表示的。在由 32 个位元构成的指令中，一般会有指示操作类型的操作码字段、标识寄存器序号的字段，以及容纳常量的字段等。5 个位元能表达值为 0~31 的 32 个数，所以刚好可以对应 32 个通用寄存器的序号。详情参见：<https://zh.wikipedia.org/wiki/指令集架构>。——译者注

② 一般常量用 16 个位元表示，而指令码本身通常含有两个表示寄存器的字段，现在要将其由 5 位扩展至 6 位，所以原来表示常量用的字段就会缩减为 14 位。——译者注

③ 原文为 arithmetic unit，应是 arithmetic logic unit（ALU）的简称，它是中央处理器的执行单元，也是所有 CPU 的核心组成部分，由“`And Gate`”和“`Or Gate`”构成，主要用于进行加、减、乘等二进制算术运算。详情参见：<https://zh.wikipedia.org/wiki/算术逻辑单元>。——译者注

据加载延迟及获取指令延迟可能会大大减缓程序执行速度。尽管近些年人们越来越重视这个问题了，但是本书不打算讨论。另外还有一个因素会改善执行效率，那就是所谓的“指令级并发”（instruction-level parallelism）。很多时下流行的 RISC 芯片都支持该技术，那些“高端”（high-end）电脑尤其如此。

此类机器的 CPU 都有多个执行单元，并具备足够的指令派发能力，以便并发执行互不依赖的一组指令（也就是说，这些指令的执行结果都不取决于其他指令，而且不会同时修改某个寄存器或状态位）。由于此技术已非常普遍，所以本书会标明那些相互独立的指令。如此一来，我们就可以说：某某公式用 8 条指令即可实现，并且在一台具备无限指令级并发能力的电脑上运行，只需 5 个周期。这也意味着，假使 CPU 能够合理安排（schedule）各条指令的执行顺序，而且其中加法器[⊖]、移位器（shifter）、逻辑单元及寄存器的数量足够多，那么从理论上讲，只要 5 个周期就能执行完这段代码了。

然而也不能过分强调这一点，因为不同电脑的指令级并行能力差得很远。比如一台 1992 年加州产的 IBM RS/6000 系列电脑[⊖]，其 CPU 的加法器就有 3 个输入端，甚至能够并发执行两条首位衔接的加法类指令（比如加法指令的目标寄存器当做比较指令或加载指令的源寄存器用[⊗]）。与之相对，那种在低端嵌入式设备上运行应用程序所用的 CPU 的构造非常简单，其寄存器堆栈可能只有一个读端口。一般说来，这种机器执行需要两个源寄存器的指令时，会多花一个周期来读取第 2 个操作数。然而这种 CPU 也许会有一个旁路（bypass），当某条指令的目标寄存器恰好是下面那条指令的源操作数时，它不用再通过寄存器堆的读端口就可以把寄存器拿来用。在这种设备上以非并行方式执行首位衔接的代码反倒发挥了它的优势。

1.3 习题

1. 将如下循环用 while 方式改写：

$$\text{for } (e_1; e_2; e_3) \text{ statement}$$

此循环能用 do 表示吗？

2. 用 C 语言写一个循环，控制名为 *i* 的无符号整数变量，使其值从 0 开始，递增到 32 位机所能表示的最大无符号数 0xFFFF FFFF（循环范围包含这两个极值）。
3. 此题留给知识丰富的读者：本书所列的基本 RISC 指令集和完整 RISC 指令集中，每条指令最多只需读取两次寄存器，写入一次寄存器。请问在常见的 RISC 指令中，有没有那种貌似简单但实际上却需要操作两个以上源数据，或多次写入寄存器的指令呢？

⊖ 在电子学中，加法器（adder）是一种用于执行加法运算的数字电路部件，是计算机核心微处理器中算术逻辑单元的基础。在这些电子系统中，它主要负责计算地址、索引等数据。此外，它还是二进制数乘法器等其他硬件的重要组成部分。详情参见：<https://zh.wikipedia.org/wiki/加法器>。——译者注

⊖ IBM RS/6000 系列电脑的发展历程可参阅其官方文档：<http://www-03.ibm.com/ibm/history/documents/pdf/rs6000.pdf>。——译者注

⊗ 原书的说法是 one feeds the other，即前一条指令将操作结果作为数据源，“喂”给后一条指令。——译者注