

## 11-二进制编码：“手持两把锒斤拷，口中疾呼烫烫烫”？

上算法和数据结构课的时候，老师们都会和你说，程序 = 算法 + 数据结构。如果对应到组成原理或者说硬件层面，算法就是我们前面讲的各种计算机指令，数据结构就对应我们接下来要讲的二进制数据。

众所周知，现代计算机都是用0和1组成的二进制，来表示所有的信息。前面几讲的程序指令用到的机器码，也是使用二进制表示的；我们存储在内存里面的字符串、整数、浮点数也都是用二进制表示的。万事万物在计算机里都是0和1，所以呢，搞清楚各种数据在二进制层面是怎么表示的，是我们必备的一课。

大部分教科书都会详细地从整数的二进制表示讲起，相信你在各种地方都能看到对应的材料，所以我就不再啰啰嗦嗦地讲这个了，只会快速地浏览一遍整数的二进制表示。

然后呢，我们重点来看一看，大家在实际应用中最常遇到的问题，也就是文本字符串是怎么表示成二进制的，特别是我们会遇到的乱码究竟是怎么回事儿。我们平时在开发的时候，所说的Unicode和UTF-8之间有什么关系。理解了这些，相信以后遇到任何乱码问题，你都能手到擒来了。

### 理解二进制的“逢二进一”

二进制和我们平时用的十进制，其实并没有什么本质区别，只是平时我们是“逢十进一”，这里变成了“逢二进一”而已。每一位，相比于十进制下的0~9这十个数字，我们只能用0和1这两个数字。

任何一个十进制的整数，都能通过二进制表示出来。把一个二进制数，对应到十进制，非常简单，就是把从右到左的第N位，乘上一个2的N次方，然后加起来，就变成了一个十进制数。当然，既然二进制是一个面向程序员的“语言”，这个从右到左的位置，自然是从0开始的。

比如0011这个二进制数，对应的十进制表示，就是 $0 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0$   
 $= 3$ ，代表十进制的3。

对应地，如果我们想要把一个十进制的数，转化成二进制，使用**短除法**就可以了。也就是，把十进制数除以2的余数，作为最右边的一位。然后用商继续除以2，把对应的余数紧靠着刚才余数的右侧，这样递归迭代，直到商为0就可以了。

比如，我们想把13这个十进制数，用短除法转化成二进制，需要经历以下几个步骤：

	商	余数	二进制位
13 / 2	6	1	1
6 / 2	3	0	0
3 / 2	1	1	1
1 / 2	0	1	1

因此，对应的二进制数，就是1101。

刚才我们举的例子都是正数，对于负数来说，情况也是一样的吗？我们可以把一个数最左侧的一位，当成是

对应的正负号，比如0为正数，1为负数，这样来进行标记。

这样，一个4位的二进制数，0011就表示为+3。而1011最左侧的第一位是1，所以它就表示-3。这个其实就是整数的**原码表示法**。原码表示法有一个很直观的缺点就是，0可以用两个不同的编码来表示，1000代表0，0000也代表0。习惯万事一对应的程序员看到这种情况，必然会被“逼死”。

于是，我们就有了另一种表示方法。我们仍然通过最左侧第一位的0和1，来判断这个数的正负。但是，我们不再把这一位当成单独的符号位，在剩下几位计算出的十进制前加上正负号，而是在计算整个二进制值的时候，在左侧最高位前面加个负号。

比如，一个4位的二进制补码数值1011，转换成十进制，就是 $-1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0$   
 $= -5$ 。如果最高位是1，这个数必然是负数；最高位是0，必然是正数。并且，只有0000表示0，1000在这样的情况下表示-8。一个4位的二进制数，可以表示从-8到7这16个整数，不会白白浪费一位。

当然更重要的一点是，用补码来表示负数，使得我们的整数相加变得很容易，不需要做任何特殊处理，只是把它当成普通的二进制相加，就能得到正确的结果。

我们简单一点，拿一个4位的整数来算一下，比如  $-5 + 1 = -4$ ， $-5 + 6 = 1$ 。我们各自把它们转换成二进制来看一看。如果它们和无符号的二进制整数的加法用的是同样的计算方式，这也就意味着它们是同样的电路。

	十进制数		二进制表示			
	-5		1	0	1	1
+	1		0	0	0	1
进位			-	1	1	-
加法结果	-4		1	1	1	0

	十进制数		二进制表示			
	-5		1	0	1	1
+	6		0	1	1	0
进位		1	1	1	-	-
加法结果	1	溢出丢弃	0	0	0	1

### 字符串的表示，从编码到数字

不仅数值可以用二进制表示，字符乃至更多的信息都能用二进制表示。最典型的例子就是**字符串**（Character String）。最早计算机只需要使用英文字符，加上数字和一些特殊符号，然后用8位的二进制，就能表示我们日常需要的所有字符了，这个就是我们常常说的**ASCII码**（American Standard Code for Information Interchange，美国信息交换标准代码）。

ASCII (1977/1986)

	_0	_1	_2	_3	_4	_5	_6	_7	_8	_9	_A	_B	_C	_D	_E	_F
0_0	NUL 0000	SOH 0001	STX 0002	ETX 0003	EOT 0004	ENQ 0005	ACK 0006	BEL 0007	BS 0008	HT 0009	LF 000A	VT 000B	FF 000C	CR 000D	SO 000E	SI 000F
1_16	DLE 0010	DC1 0011	DC2 0012	DC3 0013	DC4 0014	NAK 0015	SYN 0016	ETB 0017	CAN 0018	EM 0019	SUB 001A	ESC 001B	FS 001C	GS 001D	RS 001E	US 001F
2_32	SP 0020	! 0021	" 0022	# 0023	\$ 0024	% 0025	& 0026	' 0027	( 0028	) 0029	* 002A	+ 002B	, 002C	- 002D	. 002E	/ 002F
3_48	0 0030	1 0031	2 0032	3 0033	4 0034	5 0035	6 0036	7 0037	8 0038	9 0039	: 003A	; 003B	< 003C	= 003D	> 003E	? 003F
4_64	@ 0040	A 0041	B 0042	C 0043	D 0044	E 0045	F 0046	G 0047	H 0048	I 0049	J 004A	K 004B	L 004C	M 004D	N 004E	O 004F
5_80	P 0050	Q 0051	R 0052	S 0053	T 0054	U 0055	V 0056	W 0057	X 0058	Y 0059	Z 005A	[ 005B	\ 005C	] 005D	^ 005E	_ 005F
6_96	` 0060	a 0061	b 0062	c 0063	d 0064	e 0065	f 0066	g 0067	h 0068	i 0069	j 006A	k 006B	l 006C	m 006D	n 006E	o 006F
7_112	p 0070	q 0071	r 0072	s 0073	t 0074	u 0075	v 0076	w 0077	x 0078	y 0079	z 007A	{ 007B	 007C	} 007D	~ 007E	DEL 007F

☐ Letter ☐ Number ☐ Punctuation ☐ Symbol ☐ Other ☐ undefined ☐ Changed from 1963 version

[图片来源](#)

ASCII码就好比一个字典，用8位二进制中的128个不同的数，映射到128个不同的字符里。比如，小写字母a在ASCII里面，就是第97个，也就是二进制的0110 0001，对应的十六进制表示就是 61。而大写字母 A，就是第65个，也就是二进制的0100 0001，对应的十六进制表示就是41。

在ASCII码里面，数字9不再像整数表示法里一样，用0000 1001来表示，而是用0011 1001 来表示。字符串15也不是用0000 1111 这8位来表示，而是变成两个字符1和5连续放在一起，也就是 0011 0001 和 0011 0101，需要用两个8位来表示。

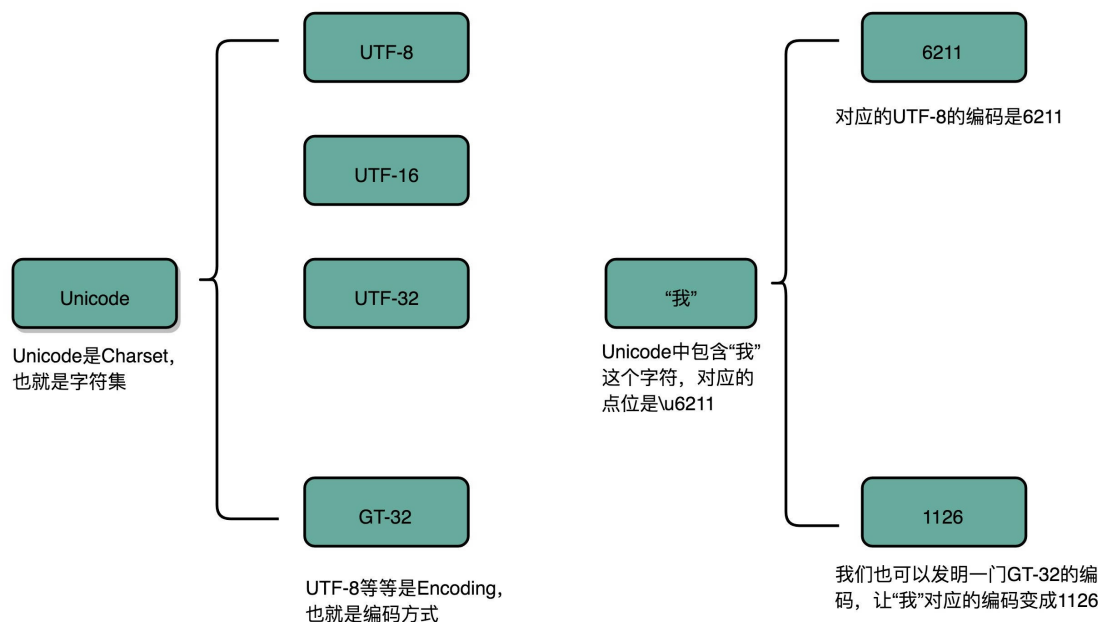
我们可以看到，最大的32位整数，就是2147483647。如果用整数表示法，只需要32位就能表示了。但是如果用字符串来表示，一共有10个字符，每个字符用8位的话，需要整整80位。比起整数表示法，要多占很多空间。

这也是为什么，很多时候我们在存储数据的时候，要采用二进制序列化这样的方式，而不是简单地把数据通过CSV或者JSON，这样的文本格式存储来进行序列化。**不管是整数也好，浮点数也好，采用二进制序列化会比存储文本省下不少空间。**

ASCII码只表示了128个字符，一开始倒也可用，毕竟计算机是在美国发明的。然而随着越来越多的不同国家的人都用上了计算机，想要表示譬如中文这样的文字，128个字符显然是不太够用的。于是，计算机工程师们开始各显神通，给自己国家的语言创建了对应的**字符集**（Charset）和**字符编码**（Character Encoding）。

字符集，表示的可以是字符的一个集合。比如“中文”就是一个字符集，不过这样描述一个字符集并不准确。想要更精确一点，我们可以说，“第一版《新华字典》里面出现的所有汉字”，这是一个字符集。这样，我们才能明确知道，一个字符在不在这个集合里面。比如，我们日常说的Unicode，其实就是一个字符集，包含了150种语言的14万个不同的字符。

而字符编码则是对于字符集里的这些字符，怎么一一用二进制表示出来的一个字典。我们上面说的Unicode，就可以用UTF-8、UTF-16，乃至UTF-32来进行编码，存储成二进制。所以，有了Unicode，其实我们可以用不止UTF-8一种编码形式，我们也可以自己发明一套 GT-32 编码，比如就叫作Geek Time 32好了。只要别人知道这套编码规则，就可以正常传输、显示这段代码。



同样的文本，采用不同的编码存储下来。如果另外一个程序，用一种不同的编码方式来进行解码和展示，就会出现乱码。这就好像两个军队用密语通信，如果用错了密码本，那看到的消息就会不知所云。在中文世界里，最典型的就是“手持两把锏斤拷，口中疾呼烫烫烫”的典故。

我曾经听说过这么一个笑话，没有经验的同学，在看到程序输出“烫烫烫”的时候，以为是程序让CPU过热发出报警，于是尝试给CPU降频来解决问题。

既然今天要彻底搞清楚编码知识，我们就来弄清楚“锏斤拷”和“烫烫烫”的来龙去脉。

锏斤拷锏斤拷PP么

On 2013 锏斤拷 12 锏斤拷 13 锏斤拷 10:28, akirya(锏斤拷[锏斤拷实锏揭谗拷锏斤拷什么锏斤拷谓锏侥伙拷锏斤拷]) wrote:

锏斤拷锏斤拷锏斤拷锏街。拷锏杰久讹拷没锏斤拷么锏斤拷锏街癸拷锏斤拷

锏斤拷 2013 锏斤拷 6 锏斤拷 2 锏斤拷锏斤拷锏斤拷锏斤拷 UTC+8 锏斤拷锏斤拷 10 时 31 锏斤拷 34 锏斤拷, carrie 写 锏斤拷锏斤拷

锏揭爸帮拷说锏斤拷锏斤拷要女 锏斤拷锏斤拷

锏斤拷女锏斤拷说锏斤拷锏斤拷要锏斤拷锏绞懃拷锏斤拷锏斤拷炼锏斤拷锏斤拷汁

平时锏斤拷识锏斤拷锏剿谗拷锏劫。拷锏叫碉拷也锏杰多, 锏斤拷只锏斤拷锏斤拷识锏斤拷锏斤拷知锏斤拷锏斤拷么锏

醇筹拷锏斤拷剩女锏剿。拷锏斤拷

锏斤拷锏斤拷锏剿筹拷蠮糠锏斤拷羌锏斤拷锏结。拷锏斤拷锏斤拷锏斤拷锏皆懃拷锏斤拷潜由峡锏斤拷椎摹锏。wbr>锏

斤拷锏斤拷锏斤拷锏斤拷坛锏斤拷锏斤拷锏斤拷锏斤拷锏斤拷揭搦拷注锏斤拷锏剿。拷锏斤拷锏斤拷锏斤拷还锏角比斤

拷锏阶的★拷锏斤拷锏斤拷 锏斤拷锏斤拷锏斤拷锏斤拷锏揭的达拷锏斤拷

...

搜索了一下我自己的个人邮件历史记录，不出意外,里面出现了各种“锏斤拷”

首先，“锏斤拷”的来源是这样的。如果我们想要用Unicode编码记录一些文本，特别是一些遗留的老字符集内的文本，但是这些字符在Unicode中可能并不存在。于是，Unicode会统一把这些字符记录为U+FFFD这个编码。如果用UTF-8的格式存储下来，就是\xef\xbf\xbd。如果连续两个这样的字符放在一起，

\xef\xbf\xbd\xef\xbf\xbd，这个时候，如果程序把这个字符，用GB2312的方式进行decode，就会变成“锕斤拷”。这就好比我们用GB2312这本密码本，去解密别人用UTF-8加密的信息，自然没办法读出有用的信息。

而“烫烫烫”，则是因为如果你用了Visual Studio的调试器，默认使用MBCS字符集。“烫”在里面是由0xCCCC来表示的，而0xCC又恰好是未初始化的内存的赋值。于是，在读到没有赋值的内存地址或者变量的时候，电脑就开始大叫“烫烫烫”了。

了解了这些原理，相信你未来在遇到中文的编码问题的时候，可以做到“手中有粮，心中不慌”了。

## 总结延伸

到这里，相信你发现，我们可以用二进制编码的方式，表示任意的信息。只要建立起字符集和字符编码，并且得到大家的认同，我们就可以在计算机里面表示这样的信息了。所以说，如果你有心，要发明一门自己的克林贡语并不是什么难事。

不过，光是明白怎么把数值和字符在逻辑层面用二进制表示是不够的。我们在计算机组成里面，关心的不只是数值和字符的逻辑表示，更要弄明白，在硬件层面，这些数值和我们一直提的晶体管和电路有什么关系。下一讲，我就会为你揭开神秘的面纱。我会从时钟和D触发器讲起，最终让你明白，计算机里的加法，是如何通过电路来实现的。

## 推荐阅读

关于二进制和编码，我推荐你读一读《编码：隐匿在计算机软硬件背后的语言》。从电报机到计算机，这本书讲述了很多计算设备的历史故事，当然，也包含了二进制及其背后对应的电路原理。

## 课后思考

你肯定会计算十进制整数的加减法，二进制的加减法也是一样的。如果二进制的加法中，有数是负数的时候该怎么处理呢？我们今天讲了补码的表示形式，如果这个负数是原码表示的，又应该如何处理？如果是补码表示的呢？请你用二进制加法试着算一算， $-5+4=-1$ ，通过原码和补码是如何进行的？

欢迎你在留言区写下你的思考和疑问，和大家一起探讨。你也可以把今天的文章分享给你朋友，和他一起学习和进步。

# 深入浅出计算机组成原理

带你掌握计算机体系全貌

徐文浩 bothub 创始人



新版升级：点击「👤请朋友读」，20位好友免费读，邀请订阅更有**现金**奖励。

## 精选留言：

- 豫樟 2019-05-20 09:14:48  
这是因为9的ASCII码是0039，换算成二进制，就是0011 1001了 [1赞]
- 小海海 2019-05-20 13:42:17  
1、-5 + 1那个表格的二进制结果好像写错了，应该是1100。  
2、补充一点对原码、反码、补码的理解：计算机中所有的加减法运算都可以转换为加法，所以一般ALU中也只实现了加法运算器，原码表示时正正、负负相加都没问题，而唯独正负相加的时候无法直接处理，如：0001(1)+1001(-1) = 1010(-2)但是结果本应该是0的，所以这时候聪明的计算机先驱们想到用反码来计算（正数的反码同原码，负数反码符号位不变，其余位取反）这样就解决了正负相加的问题，但是负负相加又有问题了，不过问题不大，可以转换为正正相加再把符号位置1即可，然后更聪明的先驱们又发明了补码（模的思想）..... 这里我放个链接：[https://www.imooc.com/article/16813?block\\_id=tuijian\\_wz](https://www.imooc.com/article/16813?block_id=tuijian_wz)  
总结：大学学组成原理时对补码一直是懵的，今天借机回顾一把又有新收获，感谢老师。还有“烫烫烫”是当年的噩梦啊，现在终于搞懂原因了，哈哈哈。 [3赞]
- 铁皮 2019-05-20 10:22:11  
老师，-5 + 1 = -4 的二进制的结果是不是错了？文中是“1110”，正确结果应该是“1100”吧？ [3赞]
- 猫头鹰爱拿铁 2019-05-20 15:32:25  
[-5+4]补=[-5]补+[4]补=[1011+0100]补=[1111]补 原码1001 [2赞]
- lzhaoh 2019-05-20 08:20:51  
在 ASCII 码里面，数字 9 不再像整数表示法里一样，用 0000 1001 来表示，而是用 0011 1001 来表示。  
  
这里不明白 [2赞]
- matter 2019-05-21 07:01:49  
老师，我明白了补码和原码表示的区别。但是对这两种方法的使用场景很模糊。什么情况下计算机使用补码，什么情况下使用原码。使用的时候会标识出来吗？



- Yyyyyy 2019-05-20 22:24:39

手持两把锒斤拷，口中疾呼烫烫烫  
脚踏千朵屯屯屯，笑看万物锒锒锒

- 阿乔 2019-05-20 21:05:34

“就是把从右到左的第  $N$  位，乘上一个  $2$  的  $N$  次方”，应该是乘以  $2^{N-1}$  次方吧？

- 一步 2019-05-20 20:45:41

文章中的这个写错了：“对应的二进制数，就是 1101” 应该是 1011

- 活的潇洒 2019-05-20 14:35:09

打卡day11

刚完成笔记如下：<https://www.cnblogs.com/luoahong/p/10893437.html>

去做踮着脚尖才能够到的事情，你的储备就会发挥出最大效能

- 龙猫 2019-05-20 14:35:08

终于要开始讲门电路了吗 哇 期待

- 姜暖 2019-05-20 11:22:46

终于知道烫烫烫咋回事了……

- 庄小P 2019-05-20 09:45:20

首先，“锒斤拷”的来源是这样的。如果我们想要用 Unicod...果我们想要用 Unicode 编码记录一些文本，特别是一些遗留的老字符集内的文本，但是这些字符在 Unicode 中可能并不存在。于是，Unicode 会统一把这些字符记录为 FFFD 这个编码。如果用 UTF-8 的格式存储下来，就是...

这里的意思是说在文本中输入不在Unicode字符集的字符，那这字符会长什么样子呢？？老师，能不能举个例子呢。

- 胖胖胖 2019-05-20 09:27:59

二进制的加减法，溢出之后有两位判断位，由此判断正负。低位正常加法进位就行

- 业余爱好者 2019-05-20 09:12:20

字符集编码也是一种加密，只不过对非计算机从业者来说是透明的，一般只需要关心屏幕上的字符图形，而无需关心他们是如何存储，传输的。想起了网络分层，用户只需关心应用层，低层协议大可不必关心。

- 古夜 2019-05-20 09:10:49

加了笑话更有料了