

link: null  
title: 珠峰架构师成长计划  
description: 最开始计算机只在美国用，八位的字节可以组合出256种不同状态。0-32种状态规定了特殊用途，一旦终端、打印机遇上约定好的这些字节被传过来时，就要做一些约定的动作，如：  
keywords: null  
author: null  
date: null  
publisher: 珠峰架构师成长计划  
stats: paragraph=63 sentences=48, words=497

1. 字符发展历史 #

1.1 字节 #

- 计算机内部，所有信息最终都是一个二进制值
- 每一个二进制位（bit）有0和1两种状态，因此八个二进制位就可以组合出256种状态，这被称为一个字节(byte)

1.2 单位 #

- 8位 = 1字节
- 1024字节 = 1K
- 1024K = 1M
- 1024M = 1G
- 1024G = 1T

1.3 JavaScript中的进制 #

1.3.1 进制表示 #

```
let a = 0b10100; // 二进制
let b = 0o24; // 八进制
let c = 20; // 十进制
let d = 0x14; // 十六进制
console.log(a == b);
console.log(b == c);
console.log(c == d);
```

1.3.2 进制转换 #

- 10进制转任意进制 10进制数.toString(目标进制)

```
console.log(c.toString(2));
```

- 任意进制转十进制 parseInt('任意进制字符串', 原始进制);

```
console.log(parseInt('10100', 2));
```

1.4 ASCII #

最开始计算机只在美国用，八位的字节可以组合出256种不同状态。0-32种状态规定了特殊用途，一旦终端、打印机遇上约定好的这些字节被传过来时，就要做一些约定的动作，如：

- 遇上0×10, 终端就换行；
- 遇上0×07, 终端就问人们干嘛；

又把所有的空格、标点符号、数字、大小写字母分别用连续的字节状态表示，一直编到了第 127 号，这样计算机就可以用不同字节来存储英语的文字了

这128个符号（包括32个不能打印出来的控制符号），只占了一个字节的后面7位，最前面的一位统一规定为0

ASCII表																									
( American Standard Code for Information Interchange 美国标准信息交换代码 )																									
高四位	ASCII控制字符													ASCII打印字符											
	0000						0001							0010	0011	0100	0101	0110	0111						
	0						1							2	3	4	5	6	7						
低四位	十进制	字符	Ctrl	代码	转义	字符解释	十进制	字符	Ctrl	代码	转义	字符解释	十进制	十进制	十进制	十进制	十进制	十进制	十进制	十进制	十进制	十进制	十进制	十进制	Ctrl
0000	0	0		^@	NUL	空字符	16	▶	^P	DLE	数据链路转义	32		48	0	64	@	80	P	96	`	112	p		
0001	1	1	☺	^A	SOH	标题开始	17	◀	^Q	DC1	设备控制 1	33	!	49	1	65	A	81	Q	97	a	113	q		
0010	2	2	☹	^B	STX	正文开始	18	↑	^R	DC2	设备控制 2	34	"	50	2	66	B	82	R	98	b	114	r		
0011	3	3	♥	^C	ETX	正文结束	19	!!	^S	DC3	设备控制 3	35	#	51	3	67	C	83	S	99	c	115	s		
0100	4	4	♦	^D	EOF	传输结束	20	◀	^T	DC4	设备控制 4	36	\$	52	4	68	D	84	T	100	d	116	t		
0101	5	5	♣	^E	ENQ	查询	21	§	^U	NAK	否定应答	37	%	53	5	69	E	85	U	101	e	117	u		
0110	6	6	♠	^F	ACK	肯定应答	22	—	^V	SYN	同步空闲	38	&	54	6	70	F	86	V	102	f	118	v		
0111	7	7	•	^G	BEL	响铃	23	↑↓	^W	ETB	传输块结束	39	'	55	7	71	G	87	W	103	g	119	w		
1000	8	8	☐	^H	BS	退格	24	↑	^X	CAN	取消	40	(	56	8	72	H	88	X	104	h	120	x		
1001	9	9	○	^I	HT	横向制表	25	↓	^Y	EM	介质结束	41	)	57	9	73	I	89	Y	105	i	121	y		
1010	A	10	☑	^J	LF	换行	26	→	^Z	SUB	替代	42	*	58	:	74	J	90	Z	106	j	122	z		
1011	B	11	♂	^K	VT	纵向制表	27	←	^[	ESC	溢出	43	+	59	;	75	K	91	[	107	k	123	{		
1100	C	12	♀	^L	FF	换页	28	└	^_	FS	文件分隔符	44	,	60	<	76	L	92	\	108	l	124			
1101	D	13	♪	^M	CR	回车	29	↔	^]	GS	组分隔符	45	-	61	=	77	M	93	]	109	m	125	}		
1110	E	14	♫	^N	SO	移出	30	▲	^^	RS	记录分隔符	46	.	62	>	78	N	94	^	110	n	126	~		
1111	F	15	♫	^O	移入	移入	31	▼	^_	US	单元分隔符	47	/	63	?	79	O	95	_	111	o	127	␣	^Backspace 代码: DEL	
注: 表 1-1 注: 11-15 符可以用“Alt + 小键盘上的数字键”方法输入。																									
2013/08/08																									

这个方案叫做 ASCII 编码

American Standard Code for Information Interchange: 美国信息互换标准代码

1.5 GB2312 #

后来西欧一些国家用的不是英文，它们的字母在ASCII里没有为了可以保存他们的文字，他们使用127号这后的空位来保存新的字母，一直编到了最后一位255。比如法语中的é的编码为130。当然了不同国家表示的符号也不一样，比如，130在法语编码中代表了é，在希伯来语编码中却代表了字母Gimel (א)。



Bytes	Encoding
FE FF	UTF16BE
FF FE	UTF16LE
EF BB BF	UTF8

因此，我们可以根据文本文件头几个字节等来判断文件是否包含BOM，以及使用哪种Unicode编码。但是，BOM字符虽然起到了标记文件编码的作用，其本身却不属于文件内容的一部分，如果读取文本文件时不去掉BOM，在某些使用场景下就会有问题。例如我们把几个JS文件合并成一个文件后，如果文件中间含有BOM字符，就会导致浏览器JS语法错误。因此，使用NodeJS读取文本文件时，一般需要去掉BOM

```
function readText(pathname) {
  var bin = fs.readFileSync(pathname);
  if (bin[0] === 0xEF && bin[1] === 0xBB && bin[2] === 0xBF) {
    bin = bin.slice(3);
  }
  return bin.toString('utf-8');
}
```

1.12.2 GBK转UTF8 #

NodeJS支持在读取文本文件时，或者在Buffer转换为字符串时指定文本编码，但遗憾的是，GBK编码不在NodeJS自身支持范围内。因此，一般我们借助iconv-lite这个三方包来转换编码。使用NPM下载该包后，我们可以按下边方式编写一个读取GBK文本文件的函数。

```
var iconv = require('iconv-lite');
function readGBKText(pathname) {
  var bin = fs.readFileSync(pathname);
  return iconv.decode(bin, 'gbk');
}
```

1.12 扩展阅读 #

- utf-8 ([http://www.ruanyifeng.com/blog/2007/10/ascii\\_unicode\\_and\\_utf-8.html](http://www.ruanyifeng.com/blog/2007/10/ascii_unicode_and_utf-8.html))
- 字符编码的故事 (<https://tianziyao.github.io/2017/07/03/%E5%AD%A6%E7%BC%96%E7%A0%81%E7%9A%84%E6%95%85%E4%BA%8B/>)