

# 字符编码

# 二进制

山东肥城

字节

字符

a123212 哈 😄

ASCII

Oct	Char	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr
000	<b>NUL</b> (null)	32	20	040	&#32;	Space	64	40	100	&#64;	@	96	60	140	&#96;	
001	<b>SOH</b> (start of heading)	33	21	041	&#33;	!	65	41	101	&#65;	A	97	61	141	&#97;	
002	<b>STX</b> (start of text)	34	22	042	&#34;	"	66	42	102	&#66;	B	98	62	142	&#98;	
003	<b>ETX</b> (end of text)	35	23	043	&#35;	#	67	43	103	&#67;	C	99	63	143	&#99;	
004	<b>EOT</b> (end of transmission)	36	24	044	&#36;	\$	68	44	104	&#68;	D	100	64	144	&#100;	
005	<b>ENQ</b> (enquiry)	37	25	045	&#37;	%	69	45	105	&#69;	E	101	65	145	&#101;	
006	<b>ACK</b> (acknowledge)	38	26	046	&#38;	&	70	46	106	&#70;	F	102	66	146	&#102;	
007	<b>BEL</b> (bell)	39	27	047	&#39;	'	71	47	107	&#71;	G	103	67	147	&#103;	
010	<b>BS</b> (backspace)	40	28	050	&#40;	(	72	48	110	&#72;	H	104	68	150	&#104;	
011	<b>TAB</b> (horizontal tab)	41	29	051	&#41;	)	73	49	111	&#73;	I	105	69	151	&#105;	
012	<b>LF</b> (NL line feed, new line)	42	2A	052	&#42;	*	74	4A	112	&#74;	J	106	6A	152	&#106;	
013	<b>VT</b> (vertical tab)	43	2B	053	&#43;	+	75	4B	113	&#75;	K	107	6B	153	&#107;	
014	<b>FF</b> (NP form feed, new page)	44	2C	054	&#44;	,	76	4C	114	&#76;	L	108	6C	154	&#108;	
015	<b>CR</b> (carriage return)	45	2D	055	&#45;	-	77	4D	115	&#77;	M	109	6D	155	&#109;	
016	<b>SO</b> (shift out)	46	2E	056	&#46;	.	78	4E	116	&#78;	N	110	6E	156	&#110;	
017	<b>SI</b> (shift in)	47	2F	057	&#47;	/	79	4F	117	&#79;	O	111	6F	157	&#111;	
020	<b>DLE</b> (data link escape)	48	30	060	&#48;	0	80	50	120	&#80;	P	112	70	160	&#112;	
021	<b>DC1</b> (device control 1)	49	31	061	&#49;	1	81	51	121	&#81;	Q	113	71	161	&#113;	
022	<b>DC2</b> (device control 2)	50	32	062	&#50;	2	82	52	122	&#82;	R	114	72	162	&#114;	
023	<b>DC3</b> (device control 3)	51	33	063	&#51;	3	83	53	123	&#83;	S	115	73	163	&#115;	
024	<b>DC4</b> (device control 4)	52	34	064	&#52;	4	84	54	124	&#84;	T	116	74	164	&#116;	
025	<b>NAK</b> (negative acknowledge)	53	35	065	&#53;	5	85	55	125	&#85;	U	117	75	165	&#117;	
026	<b>SYN</b> (synchronous idle)	54	36	066	&#54;	6	86	56	126	&#86;	V	118	76	166	&#118;	
027	<b>ETB</b> (end of trans. block)	55	37	067	&#55;	7	87	57	127	&#87;	W	119	77	167	&#119;	
030	<b>CAN</b> (cancel)	56	38	070	&#56;	8	88	58	130	&#88;	X	120	78	170	&#120;	
031	<b>EM</b> (end of medium)	57	39	071	&#57;	9	89	59	131	&#89;	Y	121	79	171	&#121;	
032	<b>SUB</b> (substitute)	58	3A	072	&#58;	:	90	5A	132	&#90;	Z	122	7A	172	&#122;	
033	<b>ESC</b> (escape)	59	3B	073	&#59;	;	91	5B	133	&#91;	[	123	7B	173	&#123;	
034	<b>FS</b> (file separator)	60	3C	074	&#60;	<	92	5C	134	&#92;	\	124	7C	174	&#124;	
035	<b>GS</b> (group separator)	61	3D	075	&#61;	=	93	5D	135	&#93;	]	125	7D	175	&#125;	
036	<b>RS</b> (record separator)	62	3E	076	&#62;	>	94	5E	136	&#94;	^	126	7E	176	&#126;	
037	<b>US</b> (unit separator)	63	3F	077	&#63;	?	95	5F	137	&#95;	_	127	7F	177	&#127;	



144	É	160	á	176	☒	192	⊥	208	⊥	224	α	240
145	æ	161	í	177	☒	193	⊥	209	⊥	225	β	241
146	Æ	162	ó	178	☒	194	⊥	210	⊥	226	Γ	242
147	ô	163	ú	179		195	⊥	211	⊥	227	π	243
148	ö	164	ñ	180	⊥	196	—	212	⊥	228	Σ	244
149	ò	165	Ñ	181	⊥	197	⊥	213	⊥	229	σ	245
150	û	166	ª	182	⊥	198	⊥	214	⊥	230	μ	246
151	ù	167	º	183	⊥	199	⊥	215	⊥	231	τ	247
152	ÿ	168	¿	184	⊥	200	⊥	216	⊥	232	Φ	248
153	Ö	169	⌈	185	⊥	201	⊥	217	⊥	233	⊕	249
154	Ü	170	⌋	186	⊥	202	⊥	218	⊥	234	Ω	250
155	¢	171	½	187	⊥	203	⊥	219	■	235	δ	251
156	£	172	¼	188	⊥	204	⊥	220	■	236	∞	252
157	⌘	173	¡	189	⊥	205	=	221	■	237	φ	253
158	⌚	174	«	190	⊥	206	⊥	222	■	238	ε	254
159	ƒ	175	»	191	⊥	207	⊥	223	■	239	∧	255

Source : [www.LookupTable.com](http://www.LookupTable.com)

混乱时代

你会如何设计？

- 字符表 (Character repertoire)
- 给字符表里的抽象字符编上一个数字，也就是字符集合到一个整数集合的映射。这种映射称为编码字符集 (CCS:Coded Character Set)
- 将CCS里字符对应的整数转换成有限长度的比特值，便于以后计算机使用一定长度的二进制形式表示该整数。这个对应关系被称为字符编码表 (CEF:Character Encoding Form)

Unicode

UTF-8 UTF-16 GBK

# Unicode Transformation Format

Unicode符号范围 | UTF-8编码方式  
(十六进制) | (二进制)

-----+-----

0000	0000-0000	007F		0xxxxxxx
0000	0080-0000	07FF		110xxxxx 10xxxxxx
0000	0800-0000	FFFF		1110xxxx 10xxxxxx 10xxxxxx
0001	0000-0010	FFFF		11110xxx 10xxxxxx 10xxxxxx 10xxxxxx



“𡗗” unicode

4E25 (100111000100101)

从“严”的最后一个二进制位开始，依次从后向前填入格式中的x，多出的位补0。这样就得到了，“严”的UTF-8编码>是“11100100 10111000 10100101”，转换成十六进制就是0xE4B8A5。

Byte order mark

0010 0000 1010 1100

字体

1. 字体是用来显示字符的，将不可见的内存中的字符显示为可见的图形。
2. 日常使用的字符是由特定字符集（比如 Unicode）定义的。
3. 并非所有字体都遵循你用的字符集。

```
iconv -f UTF8 -t  
GB18030 a.csv >b.csv
```

QA