

Character Frequency by Unicode Version

Last updated: May 2, 2009 11:58 AM

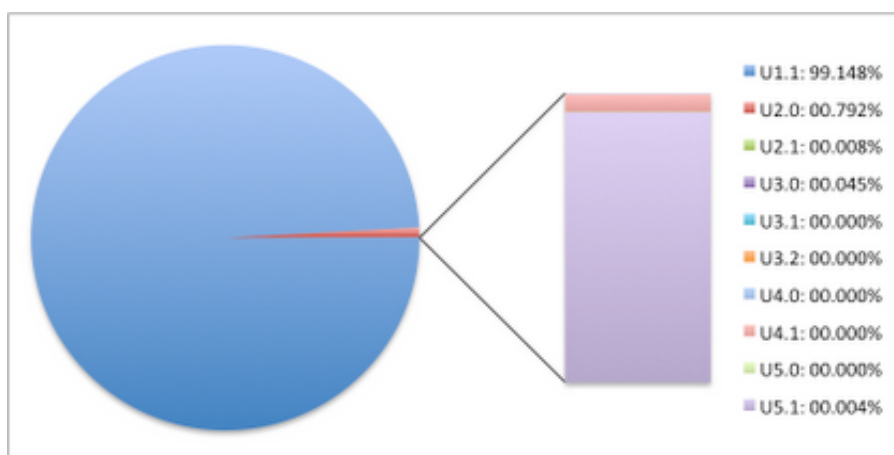


This page: <http://www.macchiato.com/unicode/character-frequency-by-unicode-version>

The image below shows frequency of usage of Unicode characters on the web based on the Unicode version where each character first appeared. One must not read too much into this data; some factors to keep in mind are:

- certain characters occur with much higher frequency in specific languages; if character is low in frequency it doesn't mean it isn't important for some modern orthography.
- there is a certain amount of noise in the data. Below a certain threshold the counts are unreliable because they are more likely to be from garbaged files.
- data is from a sample of the web (a question for the reader: how many pages are there on the web?)
- the characters are not by usage (views); that is, an occurrence of 'é' on [Le Monde.fr](http://www.lemonde.fr) counts the same as one on Jean Bleau's Facebook page, even though the former is viewed vastly more often.
- it only include public web pages; in particular it does not include the vast amount of character data in emails and other sources.

For dates of Unicode versions, see <http://www.unicode.org/history/publicationdates.html>



Below are the numbers, including the ten most frequent characters.

Value	Count	Relative Frequency	Most Frequent Chars
1.1.0.0	25,244	99.148%	e, a, i, o, t, n, r, s, l, d, ...
2.0.0.0	13,400	0.792%	이, 기, 다, 지, 의, 가, 는, 하, 사, 예, ...
2.1.0.0	2	0.008%	€,
3.0.0.0	10,083	0.045%	⌘, ⌘, ⌘, ⌘, ⌘, ⌘, ⌘, ⌘, ⌘, ⌘, ...
3.1.0.0	43,398	0%	ı, ş, ð, é, ñ, ã, ß, ð, 々, ...
3.2.0.0	904	0%	ø, U+2062, ₣, ™, , , , , *, /, ...
4.0.0.0	1,145	0%	~, 3, 4, 5, 6, 7, 8, 9, M, %, ...
4.1.0.0	1,108	0%	o, s, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, ...
5.0.0.0	1,367	0%	c, 1, 2, 3, 4, 5, 6, 7, 8, 9, 0, ...
5.1.0.0	1,621	0.004%	1, 2, 3, 4, 5, 6, 7, 8, 9, 0, ...

By Age | [By Script/Category](#)

On the second sheet are numbers by script. If there is no explicit script, the general category is given.

Character Frequency by Script/Category			
Value	Count	Relative Frequency	Most Frequent Chars
Latn - Latin	1016	73.771%	e, a, i, o, t, n, r, s, l, d, ...
Nd - Decimal_Number	20	4.883%	0, 1, 2, 3, 5, 9, 4, 8, 6, 7, ...
Cyrl - Cyrillic	403	4.592%	o, a, и, e, н, т, р, с, в, л, ...
Po - Other_Punctuation	103	3.791%	., ,, ;, /, U+0027, U+0022, !, ?, , , ., ...
Hani - Han	70362	3.701%	的, 人, 一, 中, 日, 年, 大, 有, 用, 新, ...
Arab - Arabic	264	2.442%	ب, د, ن, ت, و, ر, ي, م, ل, ا, ...
Hira - Hiragana	88	0.92%	の, い, に, し, で, て, す, る, な, た, ...
Hang - Hangul	11412	0.789%	이, 기, 다, 지, 의, 가, 는, 하, 사, 에, ...
Pd - Dash_Punctuation	11	0.768%	-, -, —, —, —, ~, -, , -, ...
Kana - Katakana	108	0.641%	ン, ト, ス, イ, ツ, ル, ク, リ, ラ, ア, ...
Thai - Thai	85	0.611%	า, น, ะ, อ, , ี, ุ, ู, ุ, ...
Pe - Close_Punctuation	47	0.395%),],) , ,] , } , } , } , } , ...
Ps - Open_Punctuation	48	0.393%	(, [, (, (, [, { , { , { , { , ...
Sm - Math_Symbol	899	0.392%	, >, =, +, <, ~, ~, , →, ×, ...
GreK - Greek	452	0.306%	α, ο, ι, τ, ε, ν, σ, ρ, κ, η, ...
Hebr - Hebrew	88	0.297%	ש, א, ב, ג, ד, ה, ו, ז, ח, ט, י, ...
Deva - Devanagari	99	0.168%	र, क, न, स, त, म, ...
Armn - Armenian	84	0.125%	ա, ւ, ր, ռ, է, Լ, Ի, մ, Կ, ւ, ...
Taml - Tamil	72	0.118%	, ட, க, த, ி, ப, ம, ர, ல, ற, ...
Lm - Modifier_Letter	41	0.114%	—, ~, ~, ~, ~, ~, ~, ~, ~, ...
Pc - Connector_Punctuation	4	0.112%	_, _ , _ , _ , _ , , { , ^ , ^
Geor - Georgian	119	0.086%	ა, ო, ე, რ, ზ, ც, ძ, ზ, ლ, ბ, ...
So - Other_Symbol	1811	0.073%	©, °, ♦, ►, ®, ★, ■, ●, ™, †, ...
Pf - Final_Punctuation	10	0.072%	», ’, ”,), ’, ’, ’, ’, ’, †
Beng - Bengali	88	0.064%	ই, র, ট, ি, ন, ক, ব, ম, ত, ...
Telu - Telugu	93	0.047%	ఁ, ం, ఱ, ల, న, ం, క, త, ...
Mlym - Malayalam	95	0.042%	ു, ി, ക, ന, റ, ര, ത, ു, യ, ല, ...
Lao - Lao	62	0.041%	າ, ບ, ັ, ື, ຸ, ົ, ຼ, ູ, ົ, ...
Knda - Kannada	84	0.039%	ಁ, ಂ, ಱ, ಲ, ನ, ದ, ಕ, ಗ, ತ, ...
Sc - Currency_Symbol	27	0.039%	\$, €, £, ¢, ¥, ¥, ¢, ¢, ¢, \$, ...
Mymr - Myanmar	156	0.03%	န, ဝ, ဂ, င, က, ဝ, ဝ, တ, င, ...
Pi - Initial_Punctuation	12	0.029%	“, «, ‘, ‘, ‘, ‘, ‘, ‘, ‘, ...
Gujr - Gujarati	83	0.026%	લ, ર, વ, ળ, ળ, ળ, ળ, ળ, ...
Mn - Nonspacing_Mark	483	0.017%	̣, ̤, ̥, ̦, ̧, ̨, ̩, ̪, ...
Guru - Gurmukhi	73	0.016%	ੳ, ਵ, ਿ, ਿ, ਿ, ਿ, ਿ, ਿ, ...
Sk - Modifier_Symbol	66	0.015%	^, ^, ^, ^, ^, ^, ^, ^, ...
Thaa - Thaana	50	0.005%	ޅ, ކ, އ, ވ, މ, ފ, ދ, ...
Khmr - Khmer	146	0.004%	្ក, ្ខ, ្គ, ្ឃ, ្ង, ្ច, ្ឆ, ្ជ, ...
Tibt - Tibetan	181	0.004%	ཨ, ཉ, ཐ, ད, དྷ, ན, པ, ...
Sinh - Sinhala	80	0.004%	ආ, භ, ඈ, ඉ, ඊ, උ, ඌ, ...
Cf - Format	136	0.004%	U+200E, U+200B, U+FEFF, U+200F, U+202C, U+202A, U+202B, U+202E, ☐, U+202D, ...
Ethi - Ethiopic	461	0.004%	ሀ, ት, ር, ሞ, የ, ክ, ስ, ዎ, በ, ል, ...
No - Other_Number	114	0.003%	², ½, ³, ¾, ¼, 1, ①, ②, ③, ④, ...
Orya - Oriya	82	0.001%	ି, ଋ, ି, କ, ଟ, ନ, ବ, ଘ, ଙ, ...
Cans - Canadian_Aboriginal	630	0%	ᑦ, ᑦ, ᑦ, ᑦ, ᑦ, ᑦ, ᑦ, ...
Bopo - Bopomofo	65	0%	ㄉ, ㄌ, ㄎ, ㄍ, ㄎ, ㄎ, ㄎ, ...
Syrc - Syriac	77	0%	ܐ, ܒ, ܓ, ܔ, ܕ, ܕ, ܕ, ܕ, ...
Cher - Cherokee	85	0%	Ⴀ, Ⴁ, Ⴂ, Ⴃ, Ⴄ, Ⴅ, Ⴆ, Ⴇ, ...
...			

