

Source: <http://www.columbia.edu/kermit/utf8.html>



characters, depending on one's personal preference or etymological theory. In this sentence, for example, "-tāng", "chiāh", "mā" and "bē" are problematic using Chinese characters. "Góá" (I/me) and "po-lé" (glass) are as written in other Sinitic languages (e.g. Mandarin, Hakka)." Wagner Amaral of Pinse & Amaral Associados notes that the Brazilian Portuguese sentence for "I can eat glass" should be identical to the Portuguese one, as the word "machueca" means "inflict pain", or rather "injuries". The words "faz mal" would more correctly translate as "cause harm". Burmese: In English the first person pronoun "I" stands for both genders, male and female. In Burmese (except in the central part of Burma) kyundaw (ကုဏ်တာ) for male and kyanma (ကုဏ်မာ) for female. Using here a fully-compliant Unicode Burmese font -- sadly one and only Padauk Graphite font exists -- rendering using graphite engine. CLICK HERE to test Burmese characters. The Quick Brown Fox The "I can eat glass" sentences do not necessarily show off the orthography of each language to best advantage. In many alphabetic written languages it is possible to include all (or most) letters (or "special" characters) in a single (often nonsense) pangram. These were traditionally used in typewriter instruction; now they are useful for stress-testing computer fonts and keyboard input methods. Here are a few examples (SEND MORE): English: The quick brown fox jumps over the lazy dog. Jamaican: Chruu, a kwik di kwik brong fox a jomp huova di liezi daag de, yu no siit? Irish: "An bfuil do érói ag bualad ó fáitíos an grá a meall lena póg éada ó slí do leasa tú?" "D'fúasail íosa Úrmac na hÓige Beannaíte pór Éava agus Ádaim." Dutch: Pa's wijze lynx bezag vroom het fikse aqueduct. German: Falsches Üben von Xylophonmusik quält jeden größeren Zwerg. (1) German: Im finsternen Jagdschloß am offenen Felsquellwasser patzte der affig-flatterhafte kauzig-höfliche Bäcker über seinem versifften kniffligen C-Xylophon. (2) Swedish: Flygande bäckasiner söka strax hwila på mjuka tuvor. Icelandic: Sævör grét áðan því úlpan var ónýt. Polish: Pchnąć w tę łódź jeżę lub ośm skrzyni fig. Czech: Příliš žluťoučký kůň úpěl dábelské kdy. Slovak: Starý kôň na hríbe knížuje tíska povädnuté ruže, na stípe sa dátel učí kvákať novú ódu o živote. Greek (monotonic): ξεσκεπάζω την ψυχοφόρα βδελυγμία Greek (polytonic): ξεσκεπάζω τὴν ψυχοφόρα βδελυγμία Russian: Съешь же ещё этих мягких французских булок да выпей чаю. Russian: В часах юга жил-был цитрус? Да, но фальшивый экземпляр! єъ. Bulgarian: Жълтата дюля беше щастлива, че пухът, който цъфна, замръзна като гъон. Sami (Northern): Vuol Ruota gedggiid leat mánge luosa ja čuovžza. Hungarian: Árvíztűrő tükörfúrógép. Spanish: El pingüino Wenceslao hizo kilómetros bajo exhaustiva lluvia y frío, añoraba a su querido cachorro. Portuguese: O próximo vôo à noite sobre o Atlântico, põe freqüentemente o único médico. (3) French: Les naïfs aegithales hâtifs pondant à Noël où il gèle sont sûrs d'être déçus et de voir leurs drôles d'œufs abîmés. Esperanto: Ehošanĝo ĉiujaude. Hebrew: זג בוט עכ דפרק חצנת דיא נומשל מותס זיך זה. Japanese (Hiragana): いろはにはくへど ちりぬるを わがよたれぞ つねならむ うゐのおくやま けふこえて あさきゆめみじ えひもせざ (4) Notes: Other phrases commonly used in Germany include: "Ein wackerer Bayer vertilgt ja bequem zwo Pfund Kalbshaxe" and, more recently, "Franz jagt im komplett verwahrlosten Taxi quer durch Bayern", but both lack umlauts and esszet. Previously, going for the shortest sentence that has all the umlauts and special characters, I had "Grüße aus Bärenhöfe (und Oechtringen)" Acute accents are not used in native German words, so I was surprised to discover "Oechtringen" in the Deutsche Bundespost Postleitzahlenbuch: It's a small village in eastern Lower Saxony. The "oe" in this case turns out to be the Lower Saxon "lengthening e" (Dehnungs-e), which makes the previous vowel long (used in a number of Lower Saxon place names such as Soest and Itzehoe), not the "e" that indicates umlaut of the preceding vowel. Many thanks to the Oechtringen-Namenschreibungsuntersuchungskomitee (Alex Bochanek, Manfred Erren, Asmus Freytag, Christoph Páper, plus Werner Lemberg who serves as Oechtringen-Namenschreibungsuntersuchungskomiteerechtschreibungsprüfer) for their relentless pursuit of the facts in this case. Conclusion: the accent almost certainly does not belong on this (or any other native German) word, but neither can it be dismissed as dirt on the page. To add to the mystery, it has been reported that other copies of the same edition of the PLZB do not show the accent! UPDATE (March 2006): David Krings was intrigued enough by this report to contact the mayor of Ebstorf, of which Oechtringen is a borough, who responded: Sehr geehrter Mr. Krings, wenn Oechtringen irgendwo mit einem Akzent auf dem O geschrieben wurde, dann kann das nur ein Fehldruck sein. Die offizielle Schreibweise lautet jedenfalls „Oechtringen“. Mit freundlichen Grüßen Der Samtgemeindebürgermeister i.A. Lothar Jessel From Karl Pentzlin (Kochel am See, Bavaria, Germany): "This German phrase is suited for display by a Fraktur (broken letter) font. It contains: all common three-letter ligatures: ffi ffl fft and all two-letter ligatures required by the Duden for Fraktur typesetting: ch ck ff fi fl ft ll fh fi ff ft tz (all in a manner such they are not part of a three-letter ligature), one example of f- where German typesetting rules prohibit ligating (marked by a ZWNJ), and all German letters a...z, ä, ö, ü, ß, f [long s] (all in a manner such that they are not part of a two-letter Fraktur ligature)." Otto Stoltz notes that "Schloß" is now spelled 'Schloss', in contrast to 'größer' (example 4) which has kept its 'ß'. Fraktur has been banned from general use, in 1942, and long-s (f) has ceased to be used with Antiqua (Roman) even earlier (the latest Antiqua-f I have seen is from 1913, but then I am no expert, so there may well be a later instance." Later Otto confirms the latter theory, "Now I've run across a book "Deutsche Rechtschreibung" (edited by Lutz Mackensen) from 1954 (my reprint is from 1956) that has kept the Antiqua-f in its dictionary part (but neither in the preface nor in the appendix)." Diaeresis is not used in Iberian Portuguese. From Yurio Miyazawa: "This poetry contains all the sounds in the Japanese language and used to be the first thing for children to learn in their Japanese class. The Hiragana version is particularly neat because it covers every character in the phonetic Hiragana character set." Yurio also sent the Kanji version: 色は匂へど 散りぬるを 我が世誰ぞ 常ならむ 有為の奥山 今日越えて 浅き夢見じ 酔ひもせざ Accented Cyrillic: (This section contributed by Vladimir Marinov.) In Bulgarian it is desirable, customary, or in some cases required to write accents over vowels. Unfortunately, no computer character sets contain the full repertoire of accented Cyrillic letters. With Unicode, however, it is possible to combine any Cyrillic letter with any combining accent. The appearance of the result depends on the font and the rendering engine. Here are two examples. Той видя бялата коса по главата ѝ и коса на рамото ѝ, и рече да ѝ рече: "Парата по пъти от пърата, не ща пари!", но си помисли: "Хей, помисли си! А ѝ река, ѝ е скочила в тази река, която щеше да тече, а не тече." По пъти пъти ват кърди и югославяни. HTML Features Here is the Russian alphabet (uppercase only) coded in three different ways, which should look identical: АБВГДЕЖЗИЙКЛМНОПРСТУФХЦЧШЩЬЫЭЮЯ (Literal UTF-8) АБВГДЕЖЗИЙКЛМНОПРСТУФХЦЧШЩЬЫЭЮЯ (Decimal numeric character reference) АБВГДЕЖЗИЙКЛМНОПРСТУФХЦЧШЩЬЫЭЮЯ (Hexadecimal numeric character reference) In another test, we use HTML language tags to distinguish Bulgarian, Russian, and Serbian, which have different italic forms for lowercase б, г, д, п, and/or т: Bulgarian: [ бѓлт ] [ бѓлт ] Mora да ям стъкло и не ме боли. Russian: [ бѓлт ] [ бѓлт ] Я могу есть стекло, это мне не вредит. Serbian: [ бѓлт ] [ бѓлт ] Mory јести стакло а да ми не шоди. Credits, Tools, and Commentary Credits: The "I can eat glass" phrase and the initial collection of translations: Ethan Mollick. Transcription / conversion to UTF-8: Frank da Cruz. Albanian: Sindi Keesan. Afrikaans: Johan Fourie, Kevin Poalses. Anglo Saxon: Frank da Cruz. Arabic: Najib Tounsi. Armenian: Vaçe Kundakçi. Belarusian: Alexey Chernyak. Bengali: Somnath Purkayastha, Deepayan Sarkar. Bislama: Dan McGarry. Braille: Frank da Cruz. Bulgarian: Sindi Keesan, Guentcho Skordev, Vladimir Marinov. Burmese: "cetanapa". Cabo Verde Creole: Cláudio Alexandre Duarte. Catalán: Jordi Bancells. Chinese: Jack Soo, Wong Pui Lam. Chinook Jargon: David Robertson. Cornish: Chris Stephens. Croatian: Marjan Baće. Czech: Stanislav Pecha, Radovan Garabík. Dutch: Peter Gotink. Pim Blokland, Rob Daniel, Rob de Wit. Erzian: Jack Rueter. Esperanto: Franko Luin, Radovan Garabík. Estonian: Meelis Roos. Faroese: Jón Gaasedal. Farsi/Persian: Payam Elahi. Finnish: Sampsa Toivanen. French: Luc Carissimo, Anne Colin du Terrail, Sean M. Burke. Galician: Laura Probaos. Georgian: Giorgi Lebanidze. German: Christoph Páper, Otto Stoltz, Karl Pentzlin, David Krings, Frank da Cruz. Gothic: Aurélien Coudurier. Greek: Ariel Glenn, Constantine Stathopoulos, Siva Nataraja, Christos Georgiou. Hebrew: Jonathan Rosenne, Tal Barnea. Hausa: Malami Buba, Tom Gewecke. Hawaiian: na Hau'oli Motta, Anela de Rego, Kaliko Trapp. Hindi: Shirish Kalele, Nitin Dahra. Hungarian: András Rácz, Mark Holzhammer. Icelandic: Andrés Magnússon, Sveinn Baldursson. International Phonetic Alphabet (IPA): Siva

Nataraja / Vincent Ramos. Irish: Michael Everson, Marion Gunn, James Kass, Curtis Clark. Italian: Thomas De Bellis. Jamaican: Stephen J. Cherin. Japanese: Makoto Takahashi, Yurio Miyazawa. Karelian: Aleksandr Semakov. Khmer: Tola Sann. Kirchröadsj: Roger Stoffers. Kreyòl: Sean M. Burke. Korean: Jungshik Shin. Langenfelder Platt: David Krings. Lao: Tola Sann. Lëtzebuergesch: Stefaan Eeckels. Lingala: Denis Moyogo Jacquerye (Nkóta ya Kóngó míbalé). (Nkóta ya Kóngó míbalé) Lithuanian: Gediminas Grigas. Lojban: Edward Cherlin. Lusatian: Ronald Schaffhirt. Macedonian: Sindi Keesan. Malay: Zarina Mustapha. Maltese: Kenneth Joseph Vella. Manx: Éanna Ó Brádaigh. Marathi: Shirish Kalele. Marquesan: Kaliko Trapp. Middle English: Frank da Cruz. Milanese: Marco Cimarosti. Mongolian: Tom Gewecke. Napoletano: Diego Quintano. Navajo: Tom Gewecke. Nórdicg: Yýlyan Rott. Norwegian: Herman Ranes. Odenwälderisch: Alexander Heß. Old Irish: Michael Everson. Old Norse: Andrés Magnússon. Papiamentu: Bianca and Denise Zanardi. Pashto: N.R. Liwal. Pfälzisch: Dr. Johannes Sander. Picard: Philippe Mennecier. Polish: Juliusz Chroboczek, Paweł Przeradowski. Portuguese: "Cláudio" Alexandre Duarte, Bianca and Denise Zanardi, Pedro Palhoto Matos, Wagner Amaral. Québécois: Laurent Detillieux. Roman: Pierpaolo Bernardi. Romanian: Juliusz Chroboczek, Ionel Mugurel. Romansch: Alexandre Suter. Ruhrdeutsch: "Timwi". Russian: Alexey Chernyak, Serge Nesterovitch. Sami: Anne Colin du Terrail, Luc Carissimo. Sanskrit: Siva Nataraja / Vincent Ramos. Sächsisch: André Müller. Schwäbisch: Otto Stolz. Scots: Jonathan Riddell. Serbian: Sindi Keesan, Ranko Narancic, Boris Daljevic, Szilvia Csorba, O. Dag. Slovak: G. Adam Stanislav, Radovan Garabík. Slovenian: Albert Kolar. Spanish: Aleida Muñoz, Laura Probaos. Swahili: Ronald Schaffhirt. Swedish: Christian Rose, Bengt Larsson. Taiwanese: Henry H. Tan-Tenn. Tagalog: Jim Soliven. Tamil: Vasee Vaseeharan. Tibetan: D. Germano, Tom Gewecke. Thai: Alan Wood's wife. Turkish: Väge Kundakçı, Tom Gewecke, Merlign Olnon. Ukrainian: Michael Zajac. Urdu: Mustafa Ali. Vietnamese: Dixon Au, [James] Đỗ Bá Phuorraine 杜伯福. Walloon: Pablo Saratzaga. Welsh: Geiriadur Prifysgol Cymru (Andrew). Yiddish: Mark David. Zeneise: Angelo Pavese. Tools Used to Create This Web Page: The UTF8-aware Kermit 95 terminal emulator on Windows, to a Unix host with the EMACS text editor. Kermit 95 displays UTF-8 and also allows keyboard entry of arbitrary Unicode BMP characters as 4 hex digits, as shown HERE. Hex codes for Unicode values can be found in The Unicode Standard (recommended) and the online code charts. When submissions arrive by email encoded in some other character set (Latin-1, Latin-2, KOI, various PC code pages, JEUC, etc), I use the TRANSLATE command of C-Kermit on the Unix host (where I read my mail) to convert the character set to UTF-8 (I could also use Kermit 95 for this; it has the same TRANSLATE command). That's it -- no "Web authoring" tools, no locales, no "smart" anything. It's just plain text, nothing more. By the way, there's nothing special about EMACS -- any text editor will do, providing it allows entry of arbitrary 8-bit bytes as text, including the 0x80-0x9F "C1" range. EMACS 21.1 actually supports UTF-8; earlier versions don't know about it and display the octal codes: either way is OK for this purpose. Commentary: Date: Wed, 27 Feb 2002 13:21:59 +0100 From: "Bruno DEDOMINICIS" <b.dedominicis@cite-sciences.fr> Subject: Je peux manger du verre, cela ne me fait pas mal. I just found out your website and it makes me feel like proposing an interpretation of the choice of this peculiar phrase. Glass is transparent and can hurt as everyone knows. The relation between people and civilisations is sometimes effusional and more often rude. The concept of breaking frontiers through globalization, in a way, is also an attempt to deny any difference. Isn't "transparency" the flag of modernity? Nothing should be hidden any more, authority is obsolete, and the new powers are supposed to reign through loving and smiling and no more through coercion... Eating glass without pain sounds like a very nice metaphor of this attempt. That is, frontiers should become glass transparent first, and be denied by incorporating them. On the reverse, it shows that through globalization, frontiers undergo a process of displacement, that is, when they are not any more speakable, they become repressed from the speech and are therefore incorporated and might become painful symptoms, as for example what happens when one tries to eat glass. The frontiers that used to separate bodies one from another tend to divide bodies from within and make them suffer.... The chosen phrase then appears as a denial of the symptom that might result from the destitution of traditional frontiers. Best, Bruno De Dominicis, Paris, France Other Unicode pages onsite: Peace in All Languages Frank's Compulsive Guide to Postal Addresses (especially the Index) Representing Middle English on the Web with UTF-8 The Kermit Bibliography (in UTF-8) Interchange of Non-English Computer Text (UTF-8 math and box-drawing) Unicode Table (in UTF-8) Unicode samplers and resources offsite: Unicode Code Converter Unicode Code Conversion (converts among different Unicode encoding forms and notations). Michael Everson's Bibliography of Typography and Scripts Does your browser support Unicode English? (James Kass) I don't know, I only work here Anyone can be provincial! Transcriptions of "Unicode" Example Unicode Usage for Business Applications UTF-8 and Unicode FAQ for Unix/Linux Unicode fonts: Code 2000 (James Kass) Unicode Fonts for Windows Computers (Alan Wood) Unicode Fonts and Tools for X11 (Markus Kuhn) Everson Mono (Michael Everson) Agfa Monotype [ Kermit 95 ] [ K95 Screen Shots ] [ C-Kermit ] [ Kermit Home ] [ Display Problems? ] [ The Unicode Consortium ] UTF-8 Sampler / The Kermit Project / Columbia University / kermit@columbia.edu