



Network Working Group  
Internet-Draft  
Expires: October 11, 2002

F. Yergeau  
Alis Technologies  
April 12, 2002

UTF-8, a transformation format of ISO 10646  
draft-yergeau-rfc2279bis-00

## Status of this Memo

This document is an Internet-Draft and is in full conformance with all provisions of [Section 10 of RFC2026](#).

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

The list of current Internet-Drafts can be accessed at <http://www.ietf.org/ietf/lid-abstracts.txt>.

The list of Internet-Draft Shadow Directories can be accessed at <http://www.ietf.org/shadow.html>.

This Internet-Draft will expire on October 11, 2002.

## Copyright Notice

Copyright (C) The Internet Society (2002). All Rights Reserved.

## Abstract

<1>

ISO/IEC 10646-1 defines a multi-octet character set called the Universal Character Set (UCS) which encompasses most of the world's writing systems. Multi-octet characters, however, are not compatible with many current applications and protocols, and this has led to the development of UTF-8, the object of this memo. UTF-8 has the characteristic of preserving the full US-ASCII range, providing compatibility with file systems, parsers and other software that rely on US-ASCII values but are transparent to other values. This memo updates and replaces [RFC 2279](#).

<2>

Discussion of this draft should take place on the ietf-

Yergeau

Expires October 11, 2002

[Page 1]

Internet-Draft

UTF-8

April 2002

charsets@iana.org mailing list.

## Table of Contents

<a href="#">1.</a>	Introduction . . . . .	<a href="#">3</a>
<a href="#">2.</a>	Notational conventions . . . . .	<a href="#">5</a>
<a href="#">3.</a>	UTF-8 definition . . . . .	<a href="#">6</a>
<a href="#">4.</a>	Versions of the standards . . . . .	<a href="#">8</a>
<a href="#">5.</a>	Byte order mark (BOM) . . . . .	<a href="#">9</a>
<a href="#">6.</a>	Examples . . . . .	<a href="#">10</a>
<a href="#">7.</a>	MIME registration . . . . .	<a href="#">11</a>
<a href="#">8.</a>	Security Considerations . . . . .	<a href="#">12</a>
	Bibliography . . . . .	<a href="#">13</a>
	Author's Address . . . . .	<a href="#">13</a>
<a href="#">A.</a>	Acknowledgements . . . . .	<a href="#">14</a>
<a href="#">B.</a>	Changes from <a href="#">RFC 2279</a> . . . . .	<a href="#">15</a>
	Full Copyright Statement . . . . .	<a href="#">16</a>

Internet-Draft

UTF-8

April 2002

## 1. Introduction

&lt;3&gt;

ISO/IEC 10646 [[ISO.10646-1](#)] defines a multi-octet character set called the Universal Character Set (UCS), which encompasses most of the world's writing systems. The same set of characters is defined by the Unicode standard [[UNICODE](#)], which further defines additional character properties and other application details of great interest to implementors. Up to the present time, changes in Unicode and amendments and additions to ISO/IEC 10646 have tracked each other, so that the character repertoires and code point assignments have remained in sync. The relevant standardization committees have committed to maintain this very useful synchronism.

&lt;4&gt;

ISO/IEC 10646 and Unicode define several encoding forms of their common repertoire: UTF-8, UCS-2, UTF-16, UCS-4 and UTF-32. In an encoding form, each character is represented as one or more encoding units. All standard UCS encoding forms except UTF-8 have an encoding unit larger than one octet, making them hard to use in many current applications and protocols that assume 8 or even 7 bit characters.

&lt;5&gt;

UTF-8, the object of this memo, has a one-octet encoding unit. It uses all bits of an octet, but has the quality of preserving the full US-ASCII [[US-ASCII](#)] range: US-ASCII characters are encoded in one octet having the normal US-ASCII value, and any octet with such a value can only stand for an US-ASCII character, and nothing else.

&lt;6&gt;

UTF-8 encodes UCS characters as a varying number of octets, where the number of octets, and the value of each, depend on the integer value assigned to the character in ISO/IEC 10646 (the character number, a.k.a. code point or Unicode scalar value). This encoding form has the following characteristics (all values are in hexadecimal):

&lt;7&gt;

- o Character numbers from U+0000 to U+007F (US-ASCII repertoire)

correspond to octets 00 to 7F (7 bit US-ASCII values). A direct consequence is that a plain ASCII string is also a valid UTF-8 string.

<8>

- o US-ASCII octet values do not appear otherwise in a UTF-8 encoded character stream. This provides compatibility with file systems or other software (e.g. the printf() function in C libraries) that parse based on US-ASCII values but are transparent to other values.

<9>

- o Round-trip conversion is easy between UTF-8 and other encoding forms.

<10>

- o The first octet of a multi-octet sequence indicates the number of octets in the sequence.

Yergeau

Expires October 11, 2002

[Page 3]

---

Internet-Draft

UTF-8

April 2002

<11>

- o The octet values FE and FF never appear.

<12>

- o Character boundaries are easily found from anywhere in an octet stream.

<13>

- o The lexicographic sorting order of strings is preserved. Of course this is of limited interest since a sort order based on character numbers is not culturally valid.

<14>

- o The Boyer-Moore fast search algorithm can be used with UTF-8 data.

<15>

- o UTF-8 strings can be fairly reliably recognized as such by a simple algorithm, i.e. the probability that a string of characters in any other encoding appears as valid UTF-8 is low, diminishing with increasing string length.

<16>

UTF-8 was originally a project of the X/Open Joint Internationalization Group XOIIG with the objective to specify a File System Safe UCS Transformation Format [[FSS UTF](#)] that is compatible with UNIX systems, supporting multilingual text in a single encoding. The original authors were Gary Miller, Greger Leijonhufvud and John Entenmann. Later, Ken Thompson and Rob Pike did significant work for the formal UTF-8.

## [2](#). Notational conventions

<17>

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in [[RFC2119](#)].

<18>

UCS characters are designated by the U+HHHH notation, where HHHH is a string of from 4 to 6 hexadecimal digits representing the character number in ISO/IEC 10646.

### [3.](#) UTF-8 definition

<19>

UTF-8 is defined by Annex D of ISO/IEC 10646-1 [[ISO.10646-1](#)].

Descriptions and formulae can also be found in the Unicode Standard [[UNICODE](#)] and in [[FSS UTF](#)].

<20>

In UTF-8, characters are encoded using sequences of 1 to 6 octets. If the repertoire is restricted to the range U+0000 to U+10FFFF (the Unicode repertoire), then only sequences of one to four octets will occur. The only octet of a "sequence" of one has the higher-order bit set to 0, the remaining 7 bits being used to encode the character number. In a sequence of n octets, n>1, the initial octet has the n

higher-order bits set to 1, followed by a bit set to 0. The remaining bit(s) of that octet contain bits from the number of the character to be encoded. The following octet(s) all have the higher-order bit set to 1 and the following bit set to 0, leaving 6 bits in each to contain bits from the character to be encoded.

<21>

The table below summarizes the format of these different octet types. The letter x indicates bits available for encoding bits of the character number.

Char. number range (hexadecimal)	UTF-8 octet sequence (binary)
0000 0000-0000 007F	0xxxxxxx
0000 0080-0000 07FF	110xxxxx 10xxxxxx
0000 0800-0000 FFFF	1110xxxx 10xxxxxx 10xxxxxx
0001 0000-001F FFFF	11110xxx 10xxxxxx 10xxxxxx 10xxxxxx
0020 0000-03FF FFFF	111110xx 10xxxxxx 10xxxxxx 10xxxxxx 10xxxxxx
0400 0000-7FFF FFFF	1111110x 10xxxxxx ... <a href="#">10xxxxxx</a>

<22>

Encoding a character to UTF-8 proceeds as follows:

<23>

1. Determine the number of octets required from the character number and the first column of the table above. It is important to note that the rows of the table are mutually exclusive, i.e. there is only one valid way to encode a given character.

<24>

2. Prepare the high-order bits of the octets as per the second column of the table.

<25>

3. Fill in the bits marked x from the bits of the character number, expressed in binary. Start from the lower-order bits of the character number and put them first in the last octet of the sequence, then the next to last, etc. until all x bits are filled in.

<26>

The definition of UTF-8 prohibits encoding character numbers between U+D800 and U+DFFF, which are reserved for use with the UTF-16 encoding form (as surrogate pairs) and do not directly represent characters. When encoding in UTF-8 from UTF-16 data, it is necessary

to first decode the UTF-16 data to obtain character numbers, which are then encoded in UTF-8 as described above.

<27>

Decoding a UTF-8 character proceeds as follows:

<28>

1. Initialize a binary number with all bits set to 0. Up to 31 bits may be needed (up to 21 if the repertoire is known to be restricted to the Unicode repertoire).

<29>

2. Determine which bits encode the character number from the number of octets in the sequence and the second column of the table above (the bits marked x).

<30>

3. Distribute the bits from the sequence to the binary number, first the lower-order bits from the last octet of the sequence and proceeding to the left until no x bits are left. The binary number is now equal to the character number.

<31>

Implementations of the decoding algorithm above MUST protect against decoding invalid sequences. For instance, a naive implementation may decode the overlong UTF-8 sequence C0 80 into the character U+0000, or the surrogate pair ED A1 8C ED BE B4 into U+233B4. Decoding invalid sequences may have security consequences or cause other problems. See Security Considerations ([Section 8](#)) below.



#### [4. Versions of the standards](#)

<32>

ISO/IEC 10646 is updated from time to time by publication of amendments and additional parts; similarly, different versions of the Unicode standard are published over time. Each new version obsoletes and replaces the previous one, but implementations, and more significantly data, are not updated instantly.

<33>

In general, the changes amount to adding new characters, which does not pose particular problems with old data. Amendment 5 to ISO/IEC 10646, however, has moved and expanded the Korean Hangul block, thereby making any previous data containing Hangul characters invalid under the new version. Unicode 2.0 has the same difference from Unicode 1.1. The official justification for allowing such an incompatible change was that no implementations and no data containing Hangul existed, a statement that is likely to be true but remains unprovable. The incident has been dubbed the "Korean mess", and the relevant committees have pledged to never, ever again make such an incompatible change (see Unicode Consortium Policies [[1](#)]).

<34>

New versions, and in particular any incompatible changes, have consequences regarding MIME character encoding labels, to be discussed in [section 5](#).

Internet-Draft

UTF-8

April 2002

## 5. Byte order mark (BOM)

&lt;35&gt;

The Unicode Standard and ISO 10646 define the character "ZERO WIDTH NO-BREAK SPACE" (U+FEFF), which is also known informally as "BYTE ORDER MARK" (abbreviated "BOM"). The latter name hints at a second possible usage of the character, in addition to its normal use as a genuine "ZERO WIDTH NO-BREAK SPACE" within text. This usage, suggested by Unicode [section 2.7](#) and ISO/IEC 10646 Annex H (informative), is to prepend a U+FEFF character to a stream of Unicode characters as a "signature"; a receiver of such a serialized stream may then use the initial character both as a hint that the stream consists of Unicode characters, as a way to recognize which UCS encoding is involved and, with encodings having a multi-octet encoding unit, as a way to recognize the serialization order of the octets. UTF-8 having a single-octet encoding unit, this last function is useless and the BOM will always appear as the octet sequence EF BB BF.

&lt;36&gt;

It is important to understand that the character U+FEFF appearing at any position other than the beginning of a stream MUST be interpreted with the semantics for the zero-width non-breaking space, and MUST NOT be interpreted as a byte-order mark. The contrapositive of that statement is not always true: the character U+FEFF in the first position of a stream MAY be interpreted as a zero-width non-breaking space, and is not always a byte-order mark. For example, if a process splits a UCS string into many parts, a part might begin with U+FEFF because there was a zero-width non-breaking space at the beginning of that substring.

&lt;37&gt;

The Unicode standard further suggests that an initial U+FEFF character may be stripped before processing the text, the rationale being that such a character in initial position may be an artifact of the encoding (an encoding signature), not a genuine intended "ZERO WIDTH NO-BREAK SPACE". Note that such stripping might affect an external process at a different layer (such as a digital signature or a count of the characters) that is relying on the presence of all characters in the stream.

&lt;38&gt;

In particular, in UTF-8 plain text it is likely, but not certain, that an initial octet sequence of EF BB BF is a signature. When concatenating two strings, it is important to strip out those signatures, because otherwise the resulting string may contain an

unintended "ZERO WIDTH NO-BREAK SPACE" at the connection point.

## 6. Examples

<39>

The character sequence "A<NOT IDENTICAL TO><ALPHA>." (U+0041, U+2262, U+0391, U+002E) is encoded in UTF-8 as follows:

```
--+-----+-----+--  
41 E2 89 A2 CE 91 2E  
--+-----+-----+--
```

<40>

The character sequence representing the Hangul characters for the Korean word "hangugo" (U+D55C, U+AD6D, U+C5B4) is encoded in UTF-8 as follows:

```
-----+-----+-----  
ED 95 9C EA B5 AD EC 96 B4  
-----+-----+-----
```

<41>

The character sequence representing the Han characters for the Japanese word "nihongo" (U+65E5, U+672C, U+8A9E) is encoded in UTF-8 as follows:

```
-----+-----+-----  
E6 97 A5 E6 9C AC E8 AA 9E  
-----+-----+-----
```

<42>

The character U+233B4 (a Chinese character meaning 'stump of tree'), prepended with a UTF-8 BOM, is encoded in UTF-8 as follows:

```
-----+-----+-----  
EF BB BF F0 A3 8E B4  
-----+-----+-----
```

Internet-Draft

UTF-8

April 2002

## [7.](#) MIME registration

&lt;43&gt;

This memo is meant to serve as the basis for registration of a MIME character set parameter (charset) [[RFC2978](#)]. The proposed charset parameter value is "UTF-8". This string labels media types containing text consisting of characters from the repertoire of ISO/IEC 10646 including all amendments at least up to amendment 5 (Korean block), encoded to a sequence of octets using the encoding scheme outlined above. UTF-8 is suitable for use in MIME content types under the "text" top-level type.

&lt;44&gt;

It is noteworthy that the label "UTF-8" does not contain a version identification, referring generically to ISO/IEC 10646. This is intentional, the rationale being as follows:

&lt;45&gt;

A MIME charset label is designed to give just the information needed to interpret a sequence of bytes received on the wire into a sequence of characters, nothing more (see [[RFC2045](#)], [section 2.2](#)). As long as a character set standard does not change incompatibly, version numbers serve no purpose, because one gains nothing by learning from the tag that newly assigned characters may be received that one doesn't know about. The tag itself doesn't teach anything about the new characters, which are going to be received anyway.

&lt;46&gt;

Hence, as long as the standards evolve compatibly, the apparent advantage of having labels that identify the versions is only that, apparent. But there is a disadvantage to such version-dependent labels: when an older application receives data accompanied by a

newer, unknown label, it may fail to recognize the label and be completely unable to deal with the data, whereas a generic, known label would have triggered mostly correct processing of the data, which may well not contain any new characters.

<47>

Now the "Korean mess" (ISO/IEC 10646 amendment 5) is an incompatible change, in principle contradicting the appropriateness of a version independent MIME charset label as described above. But the compatibility problem can only appear with data containing Korean Hangul characters encoded according to Unicode 1.1 (or equivalently ISO/IEC 10646 before amendment 5), and there is arguably no such data to worry about, this being the very reason the incompatible change was deemed acceptable.

<48>

In practice, then, a version-independent label is warranted, provided the label is understood to refer to all versions after Amendment 5, and provided no incompatible change actually occurs. Should incompatible changes occur in a later version of ISO/IEC 10646, the MIME charset label defined here will stay aligned with the previous version until and unless the IETF specifically decides otherwise.

## [8. Security Considerations](#)

<49>

Implementors of UTF-8 need to consider the security aspects of how they handle illegal UTF-8 sequences. It is conceivable that in some circumstances an attacker would be able to exploit an incautious UTF-8 parser by sending it an octet sequence that is not permitted by the UTF-8 syntax.

<50>

A particularly subtle form of this attack could be carried out against a parser which performs security-critical validity checks against the UTF-8 encoded form of its input, but interprets certain illegal octet sequences as characters. For example, a parser might prohibit the NUL character when encoded as the single-octet sequence 00, but allow the illegal two-octet sequence C0 80 and interpret it as a NUL character. Another example might be a parser which prohibits the octet sequence 2F 2E 2E 2F ("/./"), yet permits the illegal octet sequence 2F C0 AE 2E 2F. This last exploit has actually been used in a widespread virus attacking Web servers in 2001; the security threat is thus very real.

## Bibliography

- [FSS\_UTF] X/Open Company Ltd., "X/Open CAE Specification C501 -- File System Safe UCS Transformation Format (FSS\_UTF)", ISBN 1-85912-082-2, April 1995.
- [ISO.10646-1] International Organization for Standardization, "Information Technology - Universal Multiple-octet coded Character Set (UCS) - Part 1: Architecture and Basic Multilingual Plane", ISO Standard 10646-1, 2000.
- [RFC2045] Freed, N. and N. Borenstein, "Multipurpose Internet Mail Extensions (MIME) Part One: Format of Internet Message Bodies", [RFC 2045](#), November 1996.

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.
- [RFC2978] Freed, N. and J. Postel, "IANA Charset Registration Procedures", [BCP 19](#), [RFC 2978](#), October 2000.
- [UNICODE] The Unicode Consortium, "The Unicode Standard -- Version 3.0", ISBN 0-201-61633-5, 2000, <[http://www.unicode.org/unicode/standard/versions/enumeratedversions.html#Unicode\\_3\\_0\\_0](http://www.unicode.org/unicode/standard/versions/enumeratedversions.html#Unicode_3_0_0)>.
- [US-ASCII] American National Standards Institute, "Coded Character Set - 7-bit American Standard Code for Information Interchange", ANSI X3.4, 1986.
- [1] <<http://www.unicode.org/unicode/standard/policies.html>>

#### Author's Address

Fran ois Yergeau  
Alis Technologies  
100, boul. Alexis-Nihon, bureau 600  
Montr al, QC H4M 2P2  
Canada

Phone: +1 514 747 2547  
Fax: +1 514 747 2561  
EMail: fyergeau@alis.com

#### [Appendix A](#). Acknowledgements

<59>

The following have participated in the drafting and discussion of this memo: James E. Agenbroad, Andries Brouwer, Martin J. D erst, Ned Freed, David Goldsmith, Edwin F. Hart, Kent Karlsson, Markus Kuhn, Michael Kung, Alain LaBont , John Gardiner Myers, Murray Sargent, Keld Simonsen and Arnold Winkler.





## [Appendix B](#). Changes from [RFC 2279](#)

<60>

<61>

- o Significantly shortened Introduction. No more mention of UTF-1 or UTF-7, of Transformation Formats.

<62>

- o Straightened out terminology. UTF-8 now described in terms of an encoding form of the character number. UCS-2 and UCS-4 almost disappeared.

<63>

- o Note warning against decoding of invalid sequences turned into a normative MUST NOT.

<64>

- o New section about the BOM, mostly extracted and slightly adapted from [RFC 2781](#).

<65>

- o Updated a couple of references (10646-1:2000, Unicode 3, [RFC 2978](#)).

<66>

- o Added TOC.

<67>

- o Removed suggested UNICODE-1-1-UTF-8 MIME charset registration.

<68>

- o New "Notational conventions" section about [RFC 2119](#) and U+HHHH notation.

<69>

- o Pointer to Unicode Consortium Policies added in "Versions of the standards" section.

<70>

- o Added a fourth example with a non-BMP character and a BOM.

---

Internet-Draft

UTF-8

April 2002

### Full Copyright Statement

Copyright (C) The Internet Society (2002). All Rights Reserved.

This document and translations of it may be copied and furnished to others, and derivative works that comment on or otherwise explain it or assist in its implementation may be prepared, copied, published and distributed, in whole or in part, without restriction of any kind, provided that the above copyright notice and this paragraph are included on all such copies and derivative works. However, this document itself may not be modified in any way, such as by removing the copyright notice or references to the Internet Society or other Internet organizations, except as needed for the purpose of developing Internet standards in which case the procedures for copyrights defined in the Internet Standards process must be followed, or as required to translate it into languages other than English.

The limited permissions granted above are perpetual and will not be revoked by the Internet Society or its successors or assigns.

This document and the information contained herein is provided on an "AS IS" basis and THE INTERNET SOCIETY AND THE INTERNET ENGINEERING TASK FORCE DISCLAIMS ALL WARRANTIES, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE INFORMATION HEREIN WILL NOT INFRINGE ANY RIGHTS OR ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

### Acknowledgement

Funding for the RFC Editor function is currently provided by the Internet Society.

Yergeau

Expires October 11, 2002

[Page 16]