Network Working Group Internet-Draft Obsoletes: <u>5987</u> (if approved) Intended status: Standards Track Expires: June 10, 2012

Indicating Character Encoding and Language for HTTP Header Field Parameters draft-reschke-rfc5987bis-03

Abstract

By default, message header field parameters in Hypertext Transfer Protocol (HTTP) messages cannot carry characters outside the ISO-8859-1 character set. <u>RFC 2231</u> defines an encoding mechanism for use in Multipurpose Internet Mail Extensions (MIME) headers. This document specifies an encoding suitable for use in HTTP header fields that is compatible with a profile of the encoding defined in <u>RFC</u> <u>2231</u>.

Editorial Note (To be removed by RFC Editor before publication)

Distribution of this document is unlimited. Although this is not a work item of the HTTPbis Working Group, comments should be sent to the Hypertext Transfer Protocol (HTTP) mailing list at ietf-http-wg@w3.org [1], which may be joined by sending a message with subject "subscribe" to ietf-http-wg-request@w3.org [2].

Discussions of the HTTPbis Working Group are archived at http://lists.w3.org/Archives/Public/ietf-http-wg/.

XML versions, latest edits and the issues list for this document are available from <<u>http://greenbytes.de/tech/webdav/#draft-reschke-rfc5987bis</u>>. A collection of test cases is available at <<u>http://greenbytes.de/tech/tc2231/</u>>.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of <u>BCP 78</u> and <u>BCP 79</u>.

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <u>http://datatracker.ietf.org/drafts/current/</u>. 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."

This Internet-Draft will expire on June 10, 2012.

Copyright Notice

Copyright (c) 2011 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to <u>BCP 78</u> and the IETF Trust's Legal Provisions Relating to IETF Documents (<u>http://trustee.ietf.org/license-info</u>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Expires June 10, 2012 [Page 2]

Internet-Draft

Table of Contents

| $\underline{1}$. Introduction | . 4 | | | | | | |
|--|-------------|--|--|--|--|--|--|
| <u>2</u> . Notational Conventions | . <u>4</u> | | | | | | |
| <u>3</u> . Comparison to <u>RFC 2231</u> and Definition of the Encoding | . <u>4</u> | | | | | | |
| <u>3.1</u> . Parameter Continuations | . <u>5</u> | | | | | | |
| 3.2. Parameter Value Character Encoding and Language | | | | | | | |
| Information | . <u>5</u> | | | | | | |
| <u>3.2.1</u> . Definition | . <u>5</u> | | | | | | |
| <u>3.2.2</u> . Examples | | | | | | | |
| <u>3.3</u> . Language Specification in Encoded Words <u>8</u> | | | | | | | |
| $\underline{4}$. Guidelines for Usage in HTTP Header Field Definitions <u>8</u> | | | | | | | |
| 4.1. When to Use the Extension | | | | | | | |
| <u>4.2</u> . Error Handling | . <u>9</u> | | | | | | |
| 5. Security Considerations | . <u>10</u> | | | | | | |
| <u>6</u> . Acknowledgements | . <u>10</u> | | | | | | |
| <u>7</u> . References | . <u>10</u> | | | | | | |
| <u>7.1</u> . Normative References | . <u>10</u> | | | | | | |
| <u>7.2</u> . Informative References | . <u>11</u> | | | | | | |
| Appendix A. Changes from <u>RFC 5987</u> | . <u>12</u> | | | | | | |
| Appendix B. Implementation Report | . <u>12</u> | | | | | | |
| Appendix C. Change Log (to be removed by RFC Editor before | | | | | | | |
| publication) | . <u>13</u> | | | | | | |
| <u>C.1</u> . Since <u>RFC5987</u> | . <u>13</u> | | | | | | |
| <u>C.2</u> . Since <u>draft-reschke-rfc5987bis-00</u> | . <u>13</u> | | | | | | |
| <u>C.3</u> . Since <u>draft-reschke-rfc5987bis-01</u> | . <u>13</u> | | | | | | |
| <u>C.4</u> . Since <u>draft-reschke-rfc5987bis-02</u> | . <u>13</u> | | | | | | |
| Appendix D. Resolved issues (to be removed by RFC Editor | | | | | | | |
| before publication) | . <u>13</u> | | | | | | |
| <u>D.1</u> . terminology | . <u>13</u> | | | | | | |
| Appendix E. Open issues (to be removed by RFC Editor prior to | | | | | | | |
| publication) | . <u>13</u> | | | | | | |
| <u>E.1</u> . edit | . <u>14</u> | | | | | | |
| <u>E.2</u> . parmsyntax | . <u>14</u> | | | | | | |
| <u>E.3</u> . valuesyntax | | | | | | | |
| | • <u>14</u> | | | | | | |

1. Introduction

By default, message header field parameters in HTTP ([RFC2616]) messages cannot carry characters outside the ISO-8859-1 coded character set ([ISO-8859-1]). RFC 2231 ([RFC2231]) defines an encoding mechanism for use in MIME headers. This document specifies an encoding suitable for use in HTTP header fields that is compatible with a profile of the encoding defined in <u>RFC 2231</u>.

This document obsoletes [RFC5987] and moves it to "historic" status; the changes are summarized in Appendix A.

Note: in the remainder of this document, <u>RFC 2231</u> is only referenced for the purpose of explaining the choice of features that were adopted; they are therefore purely informative.

Note: this encoding does not apply to message payloads transmitted over HTTP, such as when using the media type "multipart/form-data" ([<u>RFC2388</u>]).

2. Notational Conventions

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].

This specification uses the ABNF (Augmented Backus-Naur Form) notation defined in [RFC5234]. The following core rules are included by reference, as defined in [RFC5234], Appendix B.1: ALPHA (letters), DIGIT (decimal 0-9), HEXDIG (hexadecimal 0-9/A-F/a-f), and LWSP (linear whitespace).

This specification uses terminology defined in [<u>RFC6365</u>], namely: "character encoding scheme" (below abbreviated to "character encoding"), "charset" and "coded character set".

Note that this differs from RFC 2231, which uses the term "character set" for "character encoding scheme".

3. Comparison to RFC 2231 and Definition of the Encoding

RFC 2231 defines several extensions to MIME. The sections below discuss if and how they apply to HTTP header fields.

In short:

o Parameter Continuations aren't needed (Section 3.1),

[Page 4]

- o Character Encoding and Language Information are useful, therefore a simple subset is specified (Section 3.2), and
- o Language Specifications in Encoded Words aren't needed (<u>Section 3.3</u>).

3.1. Parameter Continuations

Section 3 of [RFC2231] defines a mechanism that deals with the length limitations that apply to MIME headers. These limitations do not apply to HTTP ([RFC2616], Section 19.4.7).

Thus, parameter continuations are not part of the encoding defined by this specification.

3.2. Parameter Value Character Encoding and Language Information

Section 4 of [RFC2231] specifies how to embed language information into parameter values, and also how to encode non-ASCII characters, dealing with restrictions both in MIME and HTTP header field parameters.

However, <u>RFC 2231</u> does not specify a mandatory-to-implement character encoding, making it hard for senders to decide which encoding to use. Thus, recipients implementing this specification MUST support the "UTF-8" character encoding [RFC3629].

Furthermore, RFC 2231 allows the character encoding information to be left out. The encoding defined by this specification does not allow that.

3.2.1. Definition

The syntax for parameters is defined in <u>Section 3.6 of [RFC2616]</u> (with RFC 2616 implied LWS translated to RFC 5234 LWSP):

| parameter | = | attribute LWSP "=" LWSP value |
|------------------------|--------|---|
| attribute value | = | token token / quoted-string |
| quoted-string token | = = | <quoted-string, 2.2="" [rfc2616],="" defined="" in="" section=""> <token, 2.2="" [rfc2616],="" defined="" in="" section=""></token,></quoted-string,> |

In order to include character encoding and language information, this specification modifies the RFC 2616 grammar to be:

```
Charset/Language Encoding in HTTP December 2011
Internet-Draft
    parameter = reg-parameter / ext-parameter
    reg-parameter = parmname LWSP "=" LWSP value
    ext-parameter = parmname "*" LWSP "=" LWSP ext-value
    parmname
                 = 1*attr-char
    ext-value = charset "'" [ language ] "'" value-chars
                  ; like RFC 2231's <extended-initial-value>
                  ; (see [RFC2231], Section 7)
                = "UTF-8" / mime-charset
    charset
    mime-charset = 1*mime-charsetc
    mime-charsetc = ALPHA / DIGIT
                  / "!" / "#" / "$" / "%" / "&"
                  / "+" / "-" / "^" / "_" / "`"
                  / "{" / "}" / "~"
                  ; as <mime-charset> in Section 2.3 of [RFC2978]
                  ; except that the single guote is not included
                  ; SHOULD be registered in the IANA charset registry
    language
                  = <Language-Tag, defined in [RFC5646], Section 2.1>
    value-chars
                  = *( pct-encoded / attr-char )
    pct-encoded
                  = "%" HEXDIG HEXDIG
                  ; see [RFC3986], Section 2.1
    attr-char
                  = ALPHA / DIGIT
                  / "!" / "#" / "$" / "&" / "+" / "-" / "."
                  / "^" / " " / "`" / "|" / "~"
                  ; token except ( "*" / "'" / "%" )
```

Thus, a parameter is either a regular parameter (reg-parameter), as previously defined in <u>Section 3.6 of [RFC2616]</u>, or an extended parameter (ext-parameter).

Extended parameters are those where the left-hand side of the assignment ends with an asterisk character.

The value part of an extended parameter (ext-value) is a token that consists of three parts: the REQUIRED character encoding name (charset), the OPTIONAL language information (language), and a character sequence representing the actual value (value-chars), separated by single quote characters. Note that both character encoding names and language tags are restricted to the US-ASCII coded

character set, and are matched case-insensitively (see <u>[RFC2978]</u>, <u>Section 2.3</u> and <u>[RFC5646]</u>, <u>Section 2.1.1</u>).

Inside the value part, characters not contained in attr-char are encoded into an octet sequence using the specified character encoding. That octet sequence is then percent-encoded as specified in Section 2.1 of [RFC3986].

Producers MUST use the "UTF-8" ([<u>RFC3629</u>]) character encoding. Extension character encodings (mime-charset) are reserved for future use.

Note: recipients should be prepared to handle encoding errors, such as malformed or incomplete percent escape sequences, or nondecodable octet sequences, in a robust manner. This specification does not mandate any specific behavior, for instance, the following strategies are all acceptable:

- * ignoring the parameter,
- * stripping a non-decodable octet sequence,
- * substituting a non-decodable octet sequence by a replacement character, such as the Unicode character U+FFFD (Replacement Character).

Note: the <u>RFC 2616</u> token production (<u>[RFC2616]</u>, <u>Section 2.2</u>) differs from the production used in <u>RFC 2231</u> (imported from <u>Section 5.1 of [RFC2045]</u>) in that curly braces ("{" and "}") are excluded. Thus, these two characters are excluded from the attrchar production as well.

Note: the <mime-charset> ABNF defined here differs from the one in <u>Section 2.3 of [RFC2978]</u> in that it does not allow the single quote character (see also RFC Errata ID 1912 [<u>Err1912</u>]). In practice, no character encoding names using that character have been registered at the time of this writing.

Note: [<u>RFC5987</u>] did require support for ISO-8859-1, too; for compatibility with legacy code, recipients are encouraged to support this encoding as well.

[Page 7]

3.2.2. Examples

Non-extended notation, using "token":

foo: bar; title=Economy

Non-extended notation, using "quoted-string":

foo: bar; title="US-\$ rates"

Extended notation, using the Unicode character U+00A3 (POUND SIGN):

foo: bar; title*=utf-8'en'%C2%A3%20rates

Note: the Unicode pound sign character U+00A3 was encoded into the octet sequence C2 A3 using the UTF-8 character encoding, then percent-encoded. Also, note that the space character was encoded as %20, as it is not contained in attr-char.

Extended notation, using the Unicode characters U+00A3 (POUND SIGN) and U+20AC (EURO SIGN):

foo: bar; title*=UTF-8''%c2%a3%20and%20%e2%82%ac%20rates

Note: the Unicode pound sign character U+00A3 was encoded into the octet sequence C2 A3 using the UTF-8 character encoding, then percent-encoded. Likewise, the Unicode euro sign character U+20AC was encoded into the octet sequence E2 82 AC, then percent-encoded. Also note that HEXDIG allows both lowercase and uppercase characters, so recipients must understand both, and that the language information is optional, while the character encoding is not.

3.3. Language Specification in Encoded Words

<u>Section 5 of [RFC2231]</u> extends the encoding defined in [<u>RFC2047</u>] to also support language specification in encoded words. Although the HTTP/1.1 specification does refer to <u>RFC 2047</u> ([<u>RFC2616</u>], <u>Section</u> <u>2.2</u>), it's not clear to which header field exactly it applies, and whether it is implemented in practice (see <<u>http://tools.ietf.org/wg/httpbis/trac/ticket/111</u>> for details).

Thus, this specification does not include this feature.

<u>4</u>. Guidelines for Usage in HTTP Header Field Definitions

Specifications of HTTP header fields that use the extensions defined in <u>Section 3.2</u> ought to clearly state that. A simple way to achieve this is to normatively reference this specification, and to include

the ext-value production into the ABNF for that header field.

For instance:

Note: The Parameter Value Continuation feature defined in <u>Section</u> <u>3 of [RFC2231]</u> makes it impossible to have multiple instances of extended parameters with identical parmname components, as the processing of continuations would become ambiguous. Thus, specifications using this extension are advised to disallow this case for compatibility with <u>RFC 2231</u>.

4.1. When to Use the Extension

<u>Section 4.2 of [RFC2277]</u> requires that protocol elements containing human-readable text are able to carry language information. Thus, the ext-value production ought to be always used when the parameter value is of textual nature and its language is known.

Furthermore, the extension ought to also be used whenever the parameter value needs to carry characters not present in the US-ASCII ([USASCII]) coded character set (note that it would be unacceptable to define a new parameter that would be restricted to a subset of the Unicode character set).

4.2. Error Handling

Header field specifications need to define whether multiple instances of parameters with identical parmname components are allowed, and how they should be processed. This specification suggests that a parameter using the extended syntax takes precedence. This would allow producers to use both formats without breaking recipients that do not understand the extended syntax yet.

Example:

In this case, the sender provides an ASCII version of the title for legacy recipients, but also includes an internationalized version for recipients understanding this specification -- the latter obviously ought to prefer the new syntax over the old one.

Note: at the time of this writing, many implementations failed to ignore the form they do not understand, or prioritize the ASCII form although the extended syntax was present.

<u>5</u>. Security Considerations

The format described in this document makes it possible to transport non-ASCII characters, and thus enables character "spoofing" scenarios, in which a displayed value appears to be something other than it is.

Furthermore, there are known attack scenarios relating to decoding UTF-8.

See <u>Section 10 of [RFC3629]</u> for more information on both topics.

In addition, the extension specified in this document makes it possible to transport multiple language variants for a single parameter, and such use might allow spoofing attacks, where different language versions of the same parameter are not equivalent. Whether this attack is useful as an attack depends on the parameter specified.

6. Acknowledgements

Thanks to Martin Duerst and Frank Ellermann for help figuring out ABNF details, to Graham Klyne and Alexey Melnikov for general review, to Chris Newman for pointing out an <u>RFC 2231</u> incompatibility, and to Benjamin Carlyle, Roar Lauritzsen, and Eric Lawrence for implementer's feedback.

7. References

7.1. Normative References

| [RFC2119] | Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", <u>BCP 14</u> , <u>RFC 2119</u> , March 1997. |
|-----------|---|
| [RFC2616] | Fielding, R., Gettys, J., Mogul, J., Frystyk, H., Masinter, L., Leach, P., and T. Berners-Lee, "Hypertext Transfer Protocol HTTP/1.1", <u>RFC 2616</u> , June 1999. |
| [RFC2978] | Freed, N. and J. Postel, "IANA Charset Registration Procedures", <u>BCP 19</u> , <u>RFC 2978</u> , October 2000. |

[RFC3629] Yergeau, F., "UTF-8, a transformation format of ISO 10646", STD 63, <u>RFC 3629</u>, November 2003.

Internet-Draft Charset/Language Encoding in HTTP December 2011

- [RFC3986] Berners-Lee, T., Fielding, R., and L. Masinter, "Uniform Resource Identifier (URI): Generic Syntax", STD 66, <u>RFC 3986</u>, January 2005.
- [RFC5234] Crocker, D., Ed. and P. Overell, "Augmented BNF for Syntax Specifications: ABNF", STD 68, <u>RFC 5234</u>, January 2008.
- [RFC5646] Phillips, A., Ed. and M. Davis, Ed., "Tags for Identifying Languages", <u>BCP 47</u>, <u>RFC 5646</u>, September 2009.
- [USASCII] American National Standards Institute, "Coded Character Set -- 7-bit American Standard Code for Information Interchange", ANSI X3.4, 1986.

<u>7.2</u>. Informative References

- [Err1912] RFC Errata, "Errata ID 1912, <u>RFC 2978</u>", <<u>http://www.rfc-editor.org</u>>.
- [ISO-8859-1] International Organization for Standardization, "Information technology -- 8-bit single-byte coded graphic character sets -- Part 1: Latin alphabet No. 1", ISO/IEC 8859-1:1998, 1998.
- [RFC2045] Freed, N. and N. Borenstein, "Multipurpose Internet Mail Extensions (MIME) Part One: Format of Internet Message Bodies", <u>RFC 2045</u>, November 1996.
- [RFC2047] Moore, K., "MIME (Multipurpose Internet Mail Extensions) Part Three: Message Header Extensions for Non-ASCII Text", <u>RFC 2047</u>, November 1996.
- [RFC2231] Freed, N. and K. Moore, "MIME Parameter Value and Encoded Word Extensions: Character Sets, Languages, and Continuations", <u>RFC 2231</u>, November 1997.
- [RFC2277] Alvestrand, H., "IETF Policy on Character Sets and Languages", <u>BCP 18</u>, <u>RFC 2277</u>, January 1998.
- [RFC2388] Masinter, L., "Returning Values from Forms: multipart/ form-data", <u>RFC 2388</u>, August 1998.
- [RFC5987] Reschke, J., "Character Set and Language Encoding for Hypertext Transfer Protocol (HTTP) Header Field Parameters", <u>RFC 5987</u>, August 2010.

| Internet-Draft | Charset/Language Encoding in HTTP December 2011 |
|----------------|--|
| [RFC5988] | Nottingham, M., "Web Linking", <u>RFC 5988</u> , October 2010. |
| [RFC6266] | Reschke, J., "Use of the Content-Disposition Header Field in the Hypertext Transfer Protocol (HTTP)", <u>RFC 6266</u> , June 2011. |
| [RFC6365] | Hoffman, P. and J. Klensin, "Terminology Used in Internationalization in the IETF", <u>BCP 166</u> , <u>RFC 6365</u> , September 2011. |

URIS

- [1] <mailto:ietf-http-wg@w3.org>
- [2] <mailto:ietf-http-wg-request@w3.org?subject=subscribe>

Appendix A. Changes from <u>RFC 5987</u>

This section summarizes the changes compared to [RFC5987]:

- o The document title was changed to "Indicating Character Encoding and Language for HTTP Header Field Parameters".
- o The requirement to support the "ISO-8859-1" encoding was removed.

Appendix B. Implementation Report

The encoding defined in this document currently is used for two different HTTP header fields:

- o "Content-Disposition", defined in [RFC6266], and
- o "Link", defined in [<u>RFC5988</u>].

As the encoding is a profile/clarification of the one defined in [<u>RFC2231</u>] in 1997, many user agents already supported it for use in "Content-Disposition" when [<u>RFC5987</u>] got published.

Since the publication of [<u>RFC5987</u>], two more popular desktop user agents have added support for this encoding; see <<u>http://purl.org/</u> <u>NET/http/content-disposition-tests#encoding-2231-char</u>> for details. At this time, only one major desktop user agent (Safari) does not support it.

Note that the implementation in Internet Explorer 9 does not support the ISO-8859-1 character encoding; this document revision acknowledges that UTF-8 is sufficient for expressing all code points, and removes the requirement to support ISO-8859-1.

The "Link" header field, on the other hand, was only recently specified in [RFC5988]. At the time of this writing, no User Agent supported the "title*" parameter, using the encoding defined by this document, but implementation for Firefox was already in progress (see <<u>https://bugzilla.mozilla.org/show_bug.cgi?id=663057</u>>).

Appendix C. Change Log (to be removed by RFC Editor before publication)

<u>C.1</u>. Since <u>RFC5987</u>

Only editorial changes for the purpose of starting the revision process (obs5987).

C.2. Since draft-reschke-rfc5987bis-00

Resolved issues "iso-8859-1" and "title" (title simplified). Added and resolved issue "historic5987".

C.3. Since draft-reschke-rfc5987bis-01

Added issues "httpbis", "parmsyntax", "terminology" and "valuesyntax". Closed issue "impls".

<u>C.4</u>. Since <u>draft-reschke-rfc5987bis-02</u>

Resolved issue "terminology".

<u>Appendix D</u>. Resolved issues (to be removed by RFC Editor before publication)

Issues that were either rejected or resolved in this version of this document.

D.1. terminology

Type: edit

julian.reschke@greenbytes.de (2011-09-17): Try to be consistent with the terminology defined in <u>RFC 6365</u>.

Resolution (2011-12-08): Done (but abbreviating "character encoding scheme" to "character encoding").

<u>Appendix E</u>. Open issues (to be removed by RFC Editor prior to publication)

Expires June 10, 2012 [Page 13]

<u>E.1</u>. edit

Type: edit

julian.reschke@greenbytes.de (2011-04-15): Umbrella issue for editorial fixes/enhancements.

E.2. parmsyntax

Type: edit

<<u>http://lists.w3.org/Archives/Public/ietf-http-wg/20110ctDec/</u> 0159.html>

James.H.Manger@team.telstra.com (2011-11-02): Noted by James Manger: "Presumably <u>RFC5987</u> (or its predecessors) decided it was highly unlikely that any parameter names in use ended in "*" (though they are valid) so it could redefine the syntax of values for such names." - add a note that the notation indeed overloads parameter name syntax and clarify the use.

E.3. valuesyntax

Type: edit

<<u>http://lists.w3.org/Archives/Public/ietf-http-wg/20110ctDec/</u> 0159.html>

James.H.Manger@team.telstra.com (2011-11-02): Noted by James Manger: "Curiously, <u>RFC5987</u> disobeys the proposed recommendations for new parameters. It allows foo*=UTF-8''coll%C3%A8gues but not foo*="UTF-8''coll%C3%A8gues" That might be ok with a parser that understands token, quoted-string, and <u>RFC5987</u>, but presumably it will cause problems when <u>RFC5987</u> processing is done after a "standard httpbis parser" handles the token | quoted-string step. " - add a note clarifying that this is indeed a shortcoming of the format, and what it means for implementations.

E.4. httpbis

Type: edit

julian.reschke@greenbytes.de (2011-09-17): The document refers normatively to <u>RFC 2616</u>. Should it continue to do so, or should we wait for HTTPbis? This may affect edge case in the ABNF, such as the definition of linear white space or the characters allowed in "token".

Author's Address

Julian F. Reschke greenbytes GmbH Hafenweg 16 Muenster, NW 48155 Germany

EMail: julian.reschke@greenbytes.de URI: <u>http://greenbytes.de/tech/webdav/</u>