

Network Working Group
Internet Draft

Ned Freed, Innosoft
Jon Postel, ISI
<[draft-freed-charset-reg-03.txt](#)>

IANA Charset
Registration Procedures

September 1997

Status of this Memo

This document is an Internet-Draft. 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. Internet-Drafts may be updated, replaced, or obsoleted by other documents at any time. It is not appropriate to use Internet-Drafts as reference material or to cite them other than as a "working draft" or "work in progress".

To learn the current status of any Internet-Draft, please check the `1id-abstracts.txt` listing contained in the Internet-Drafts Shadow Directories on `ds.internic.net` (US East Coast), `nic.nordu.net` (Europe), `ftp.isi.edu` (US West Coast), or `munari.oz.au` (Pacific Rim).

1. Abstract

MIME [RFC-2045, [RFC-2046](#), [RFC-2047](#)] and various other modern Internet protocols are capable of using many different charsets. This in turn means that the ability to label different charsets is essential. This registration procedure exists solely to associate a specific name or names with a given charset and to give an indication of whether or not a given charset can be used in MIME text objects. In particular, the general applicability and appropriateness of a given registered charset is a protocol issue, not a registration issue, and is not dealt with by this registration procedure.

2. Definitions and Notation

The following sections define various terms used in this document.

2.1. Requirements Notation

This document occasionally uses terms that appear in capital letters. When the terms "MUST", "SHOULD", "MUST NOT", "SHOULD NOT", and "MAY" appear capitalized, they are being used to indicate particular requirements of this specification. A discussion of the meanings of these terms appears in [RFC-2119].

2.2. Character

A member of a set of elements used for the organisation, control, or representation of data.

2.3. Charset

The term "charset" (referred to as a "character set" in previous versions of this document) is used here to refer to a method of converting a sequence of octets into a sequence of characters. This conversion may also optionally produce additional control information such as directionality indicators.

Note that unconditional and unambiguous conversion in the other direction is not required, in that not all characters may be representable by a given charset and a charset may provide more than one sequence of octets to represent a particular sequence of characters.

This definition is intended to allow charsets to be defined in a variety of different ways, from simple single-table mappings such as US-ASCII to complex table switching methods such as those that use ISO 2022's techniques, to be used as charsets. However, the definition associated with a charset name must fully specify the mapping to be performed. In particular, use of external profiling information to determine the exact mapping is not permitted.

Expires March 1998

[Page 2]

HISTORICAL NOTE: The term "character set" was originally used in MIME to describe such straightforward schemes as US-ASCII and ISO-8859-1 which consist of a small set of characters and a simple one-to-one mapping from single octets to single characters. Multi-octet character encoding schemes and switching techniques make the situation much more complex. As such, the definition of this term was revised to emphasize both the conversion aspect of the process, and the term itself has been changed to "charset" to emphasize that it is not, after all, just a set of characters. A discussion of these issues as well as specification of standard terminology for use in the IETF appears in [RFC 2130](#).

[2.4.](#) Coded Character Set

A Coded Character Set (CCS) is a mapping from a set of abstract characters to a set of integers. Examples of coded character sets are ISO 10646 [[ISO-10646](#)], US-ASCII [[US-ASCII](#)], and the ISO-8859 series [[ISO-8859](#)].

[2.5.](#) Character Encoding Scheme

A Character Encoding Scheme (CES) is a mapping from a Coded Character Set or several coded character sets to a set of octets. A given CES is typically associated with a single CCS; for example, UTF-8 applies only to ISO 10646.

[3.](#) Registration Requirements

Registered charsets are expected to conform to a number of requirements as described below.

[3.1.](#) Required Characteristics

Registered charsets **MUST** conform to the definition of a "charset" given above. In addition, charsets intended for use in MIME content types under the "text" top-level type must conform to the restrictions on that type described in RFC [2045](#). **All registered charsets MUST note whether or not they are suitable for use in MIME.**

Expires March 1998

[Page 3]

All charsets which are constructed as a composition of a CCS and a CES MUST either include the CCS and CES they are based on in their registration or else cite a definition of their CCS and CES that appears elsewhere.

All registered charsets MUST be specified in a stable, openly available specification. Registration of charsets whose specifications aren't stable and openly available is forbidden.

3.2. New Charsets

This registration mechanism is not intended to be a vehicle for the definition of entirely new charsets. This is due to the fact that the registration process does NOT contain adequate review mechanisms for such undertakings.

As such, only charsets defined by other processes and standards bodies, or specific profiles of such charsets, are eligible for registration.

3.3. Naming Requirements

One or more names MUST be assigned to all registered charsets. Multiple names for the same charset are permitted, but if multiple names are assigned a single primary name for the charset MUST be identified. All other names are considered to be aliases for the primary name and use of the primary name is preferred over use of any of the aliases.

Each assigned name MUST uniquely identify a single charset. All charset names MUST be suitable for use as the value of a MIME content type charset parameter and hence MUST conform to MIME parameter value syntax. This applies even if the specific charset being registered is not suitable for use with the "text" media type.

Finally, charsets being registered for use with the "text" media type MUST have a primary name that conforms to the more restrictive syntax of the charset field in a MIME encoded-word [[RFC-2047](#)].

Expires March 1998

[Page 4]

3.4. Functionality Requirement

Charsets must function as actual charsets: Registration of things that are better thought of as a transfer encoding, as a media type, or as a collection of separate entities of another type, is not allowed. For example, although HTML could theoretically be thought of as a charset, it is really better thought of as a media type and as such it cannot be registered as a charset.

3.5. Usage and Implementation Requirements

Use of a large number of charsets in a given protocol may hamper interoperability. However, the use of a large number of undocumented and/or unlabelled charsets hampers interoperability even more.

A charset should therefore be registered ONLY if it adds significant functionality that is valuable to a large community, OR if it documents existing practice in a large community. Note that charsets registered for the second reason should be explicitly marked as being of limited or specialized use and should only be used in Internet messages with prior bilateral agreement.

3.6. Publication Requirements

Charset registrations can be published in RFCs, however, RFC publication is not required to register a new charset.

The registration of a charset does not imply endorsement, approval, or recommendation by the IANA, IESG, or IETF, or even certification that the specification is adequate. It is expected that applicability statements for particular applications will be published from time to time that recommend implementation of, and support for, charsets that have proven particularly useful in those contexts.

3.7. MIBenum Requirements

Each registered charset MUST also be assigned a unique enumerated integer value. These "MIBenum" values are defined

Expires March 1998

[Page 5]

by and used in the Printer MIB [[RFC-1759](#)].

A MIBenum value for each charset will be assigned by IANA at the time of registration.

[4.](#) Registration Procedure

The following procedure has been implemented by the IANA for review and approval of new charsets. This is not a formal standards process, but rather an administrative procedure intended to allow community comment and sanity checking without excessive time delay.

[4.1.](#) Present the Charset to the Community

Send the proposed charset registration to the "ietf-charsets@iana.org" mailing list. This mailing list has been established for the sole purpose of reviewing proposed charset registrations. Proposed charsets are not formally registered and must not be used; the "x-" prefix specified in [RFC 2045](#) can be used until registration is complete.

The intent of the public posting is to solicit comments and feedback on the definition of the charset and the name chosen for it over a two week period.

[4.2.](#) Charset Reviewer

When the two week period has passed and the registration proposer is convinced that consensus has been achieved, the registration application should be submitted to IANA and the charset reviewer. The charset reviewer, who is appointed by the IETF Applications Area Director(s), either approves the request for registration or rejects it. Rejection may occur because of significant objections raised on the list or objections raised externally. If the charset reviewer considers the registration sufficiently important and controversial, a last call for comments may be issued to the full IETF. The charset reviewer may also recommend standards track processing (before or after registration) when that appears appropriate and the level of specification of the charset is adequate.

Expires March 1998

[Page 6]

Decisions made by the reviewer must be posted to the ietf-charsets mailing list within 14 days. Decisions made by the reviewer may be appealed to the IESG.

4.3. IANA Registration

Provided that the charset registration has either passed review or has been successfully appealed to the IESG, the IANA will register the charset, assign a MIBenum value, and make its registration available to the community.

5. Location of Registered Charset List

Charset registrations will be posted in the anonymous FTP file "ftp://ftp.isi.edu/in-notes/iana/assignments/character-sets" and all registered charsets will be listed in the periodically issued "Assigned Numbers" RFC [currently [RFC-1700](#)]. The description of the charset may also be published as an Informational RFC by sending it to "rfc-editor@isi.edu" (please follow the instructions to RFC authors [[RFC-1543](#)]).

6. Registration Template

To: ietf-charsets@iana.org
Subject: Registration of new charset

Charset name(s):

(All names must be suitable for use as the value of a MIME content-type parameter.)

Published specification(s):

(A specification for the charset must be openly available that accurately describes what is being registered. If a charset is defined as a composition of a CCS and a CES then these definitions must either be included or referenced.)

Person & email address to contact for further information:

Expires March 1998

[Page 7]

7. Security Considerations

This registration procedure is not known to raise any sort of security considerations that are appreciably different from those already existing in the protocols that employ registered charsets.

8. References

[ISO-2022]

International Standard -- Information Processing --
Character Code Structure and Extension Techniques,
ISO/IEC 2022:1994, 4th ed.

[ISO-8859]

International Standard -- Information Processing -- 8-bit
Single-Byte Coded Graphic Character Sets

- Part 1: Latin Alphabet No. 1, ISO 8859-1:1987, 1st ed.
- Part 2: Latin Alphabet No. 2, ISO 8859-2:1987, 1st ed.
- Part 3: Latin Alphabet No. 3, ISO 8859-3:1988, 1st ed.
- Part 4: Latin Alphabet No. 4, ISO 8859-4:1988, 1st ed.
- Part 5: Latin/Cyrillic Alphabet, ISO 8859-5:1988, 1st ed.
- Part 6: Latin/Arabic Alphabet, ISO 8859-6:1987, 1st ed.
- Part 7: Latin/Greek Alphabet, ISO 8859-7:1987, 1st ed.
- Part 8: Latin/Hebrew Alphabet, ISO 8859-8:1988, 1st ed.
- Part 9: Latin Alphabet No. 5, ISO/IEC 8859-9:1989, 1st ed.

International Standard -- Information Technology -- 8-bit
Single-Byte Coded Graphic Character Sets

- Part 10: Latin Alphabet No. 6, ISO/IEC 8859-10:1992, 1st ed.

[ISO-10646]

ISO/IEC 10646-1:1993(E), "Information technology --
Universal Multiple-Octet Coded Character Set (UCS) --
Part 1: Architecture and Basic Multilingual Plane",
JTC1/SC2, 1993.

Expires March 1998

[Page 8]

[RFC-1590]

Postel, J., "Media Type Registration Procedure", [RFC 1590](#), USC/Information Sciences Institute, March 1994.

[RFC-1700]

Reynolds, J. and Postel, J., "Assigned Numbers", STD 2, [RFC 1700](#), USC/Information Sciences Institute, October 1994.

[RFC-1759]

Smith, R., Wright, F., Hastings, T., Zilles, S., Gyllenskog, J., "Printer MIB", [RFC 1759](#), March 1995.

[RFC-2045]

Freed, N. and Borenstein, N., "Multipurpose Internet Mail Extensions (MIME) Part One: Format of Internet Message Bodies", [RFC 2045](#), Bellcore, Innosoft, November 1996.

[RFC-2046]

Freed, N. and Borenstein, N., "Multipurpose Internet Mail Extensions (MIME) Part Two: Media Types", [RFC 2046](#), Bellcore, Innosoft, November 1996.

[RFC-2047]

Moore, K., "Multipurpose Internet Mail Extensions (MIME) Part Three: Representation of Non-Ascii Text in Internet Message Headers", [RFC 2047](#), University of Tennessee, November 1996.

[RFC-2119]

Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [RFC 2119](#), March 1997.

[RFC-2130]

Weider, C., Preston, C., Simonsen, K., Alvestrand, H., Atkinson, R., Crispin, M., Svanberg, P., "Report from the IAB Character Set Workshop", [RFC 2130](#), April 1997.

Expires March 1998

[Page 9]

[US-ASCII]

Coded Character Set -- 7-Bit American Standard Code for
Information Interchange, ANSI X3.4-1986.

9. Authors' Addresses

Ned Freed

Innosoft International, Inc.

1050 Lakes Drive

West Covina, CA 91790

USA

tel: +1 626 919 3600

fax: +1 626 919 3614

email: ned.freed@innosoft.com

Jon Postel

USC/Information Sciences Institute

4676 Admiralty Way

Marina del Rey, CA 90292

USA

tel: +1 310 822 1511

fax: +1 310 823 6714

email: Postel@ISI.EDU

[Appendix A](#) -- IANA and RFC Editor To-Do List

VERY IMPORTANT NOTE: This appendix is intended to communicate various editorial and procedural tasks the IANA and the RFC Editor should undertake prior to publication of this document as an RFC. This appendix should NOT appear in the actual RFC version of this document!

This document refers to the character set mailing list ietf-charsets@iana.org. This alias needs to be established and should initially point to ietf-charsets@innosoft.com.

Expires March 1998

[Page 10]