

TELNET CHARSET Option

Status of this Memo

This memo defines an Experimental Protocol for the Internet community. This memo does not specify an Internet standard of any kind. Discussion and suggestions for improvement are requested. Distribution of this memo is unlimited.

Abstract

This document specifies a mechanism for passing character set and translation information between a TELNET client and server. Use of this mechanism enables an application used by a TELNET user to send and receive data in the correct character set.

Either side can (subject to option negotiation) at any time request that a (new) character set be used.

1. Command Names and Codes

CHARSET.....42

REQUEST	01
ACCEPTED	02
REJECTED	03
TTABLE-IS	04
TTABLE-REJECTED	05
TTABLE-ACK	06
TTABLE-NAK	07

As a convenience, standard TELNET text and codes for commands used in this document are reproduced here (excerpted from [1]):

All TELNET commands consist of at least a two byte sequence: the "Interpret as Command" (IAC) escape character followed by the code for the command. The commands dealing with option negotiation are three byte sequences, the third byte being the code for the option referenced. ... [O]nly the IAC need be doubled to be sent as data, and the other 255 codes may be passed transparently. The following are [some of] the defined TELNET commands. Note that these codes and code sequences have the indicated meaning only when immediately preceded by an IAC.

NAME	CODE	MEANING
SE	240	End of subnegotiation parameters.
SB	250	Indicates that what follows is subnegotiation of the indicated option.
WILL	251	Indicates the desire to begin performing, or confirmation that you are now performing, the indicated option.
WON'T	252	Indicates the refusal to perform, or continue performing, the indicated option.
DO	253	Indicates the request that the other party perform, or confirmation that you are expecting the other party to perform, the indicated option.
DON'T	254	Indicates the demand that the other party stop performing, or confirmation that you are no longer expecting the other party to perform, the indicated option.
IAC	255	Data Byte 255.

2. Command Meanings

A very simple meta-syntax is used, where most tokens represent previously defined items (such as IAC); angle-brackets ("`<>`") are used for items to be further defined; curly-braces ("`{}`") are used around optional items; ellipses represent repeated sequences of items; and quotes are used for literal strings.

IAC WILL CHARSET

The sender REQUESTS permission to, or AGREES to, use CHARSET option subnegotiation to choose a character set.

IAC WON'T CHARSET

The sender REFUSES to use CHARSET option subnegotiation to choose a character set.

IAC DO CHARSET

The sender REQUESTS that, or AGREES to have, the other side use CHARSET option subnegotiation to choose a character set.

IAC DON'T CHARSET

The sender DEMANDS that the other side not use the CHARSET option subnegotiation.

IAC SB CHARSET REQUEST { "[TTABLE]" <Version> } <char set list> IAC SE

Char set list:

<sep> <character set> { ... <sep> <character set> }

This message initiates a new CHARSET subnegotiation. It can only be sent by a side that has received a DO CHARSET message and sent a WILL CHARSET message (in either order).

The sender requests that all text sent to and by it be encoded in one of the specified character sets.

If the string [TTABLE] appears, the sender is willing to accept a mapping (translation table) between any character set listed in <char set list> and any character set desired by the receiver.

<Version> is an octet whose binary value is the highest version level of the TTABLE-IS message which can be sent in response. This field must not be zero. See the TTABLE-IS message for the permitted version values.

<Char set list> is a sequence of 7-BIT ASCII printable characters. The first octet defines the separator character (which must not appear within any character set). It is terminated by the IAC SE sequence. Case is not significant. It consists of one or more character sets. The character sets should appear in order of preference (most preferred first).

<Sep> is a separator octet, the value of which is chosen by the sender. Examples include a space or a semicolon. Any value other than IAC is allowed. The obvious choice is a space or any other punctuation symbol which does not appear in any of the character set names.

<Character set> is a sequence of 7-BIT ASCII printable characters. Case is not significant.

If a requested character set name does not start with "X-" or "x-", it MUST be registered with the Internet Assigned Number Authority (IANA) [2].

The receiver responds in one of four ways:

If the receiver is already sending text to and expecting text from the sender to be encoded in one of the specified character sets, it sends a positive acknowledgment (CHARSET ACCEPTED); it MUST NOT ignore the message. (Although ignoring the message is perhaps suggested by some interpretations of the relevant RFCs ([1], [3]), in the interests of determinacy it is not permitted. This ensures that the issuer does not need to time out and infer a response, while avoiding (because there is no response to a positive acknowledgment) the non-terminating subnegotiation which is the rationale in the RFCs for the non-response behavior.)

If the receiver is capable of handling at least one of the specified character sets, it can respond with a positive acknowledgment for one of the requested character sets. Normally, it should pick the first set it is capable of handling but may choose one based on its own preferences. After doing so, each side MUST encode subsequent text in the specified character set.

If the string [TTABLE] is present, and the receiver prefers to use a character set not included in <char set list>, and is capable of doing so, it can send a translate table (TTABLE-IS) response.

If the receiver is not capable of handling any of the specified character sets, it sends a negative acknowledgment (CHARSET REJECTED).

Because it is not valid to reply to a CHARSET REQUEST message with another CHARSET REQUEST message, if a CHARSET REQUEST message is received after sending one, it means that both sides have sent them simultaneously. In this case, the server side MUST issue a negative acknowledgment. The client side MUST respond to the one from the server.

IAC SB CHARSET ACCEPTED <Charset> IAC SE

This is a positive acknowledgment response to a CHARSET REQUEST message; the receiver of the CHARSET REQUEST message acknowledges its receipt and accepts the indicated character set.

<Charset> is a character sequence identical to one of the character sets in the CHARSET REQUEST message. It is terminated by the IAC SE sequence.

Text messages which follow this response must now be coded in the indicated character set. This message terminates the current CHARSET subnegotiation.

IAC SB CHARSET REJECTED IAC SE

This is a negative acknowledgment response to a CHARSET REQUEST message; the receiver of the CHARSET REQUEST message acknowledges its receipt but refuses to use any of the requested character sets. Messages can not be sent in any of the indicated character sets. This message can also be sent by the sender of a TTABLE-IS message, if multiple TTABLE-NAK messages were sent in response. This message terminates the current CHARSET subnegotiation.

IAC SB CHARSET TTABLE-IS <version> <syntax for version> IAC SE

In response to a CHARSET REQUEST message in which [TTABLE] was specified, the receiver of the CHARSET REQUEST message acknowledges its receipt and is transmitting a pair of tables which define the mapping between specified character sets.

<Version> is an octet whose binary value is the version level of this TTABLE-IS message. Different versions have different syntax. The lowest version level is one (zero is not valid). The current highest version level is also one. This field is provided so that future versions of the TTABLE-SEND message can be specified, for example, to handle character sets for which there is no simple one-to-one character-for-character translation. This might include some forms of multi-octet character sets for which translation algorithms or subsets need to be sent.

Syntax for Version 1:

<sep> <char set name 1> <sep> < char size 1> < char count 1> <char set name 2> <sep> <char size 2> <char count 2> <map 1> <map 2>

<Sep> is a separator octet, the value of which is chosen by the sender. Examples include a space or a semicolon. Any value other than IAC is allowed. The obvious choice is a space or any other punctuation symbol which does not appear in either of the character set names.

<Char set name 1> and <Char set name 2> are sequences of 7-BIT ASCII printable characters which identify the two character sets for which a mapping is being specified. Each is terminated by <sep>. Case is not significant. If a character set name does not

start with "X-" or "x-", it MUST be registered with IANA. <Char set name 1> MUST be chosen from the <char set list> in the CHARSET REQUEST message. <Char set name 2> can be arbitrarily chosen. Text on the wire MUST be encoded using <char set name 2>.

<Char size 1> and <char size 2> are single octets each. The binary value of the octet is the number of bits nominally required for each character in the corresponding table. It SHOULD be a multiple of eight.

<Char count 1> and <char count 2> are each three-octet binary fields in Network Byte Order [6]. Each specifies how many characters (of the maximum $2^{**}<\text{char size}>$) are being transmitted in the corresponding map.

<Map1> and <Map 2> each consist of the corresponding <char count> number of characters. These characters form a mapping from all or part of the characters in one of the specified character sets to the correct characters in the other character set. If the indicated <char count> is less than $2^{**}<\text{char size}>$, the first <char count> characters are being mapped, and the remaining characters are assumed to not be changed (and thus map to themselves). That is, each map contains characters 0 through <char count> -1. <Map 1> maps from <char set name 1> to <char set name 2>. <Map 2> maps from <char set name 2> to <char set name 1>. Translation between the character sets is thus an obvious process of using the binary value of a character as an index into the appropriate map. The character at that index replaces the original character. If the index exceeds the <char count> for the map, no translation is performed for the character.

[Note to implementers: since TELNET works in octets, it is possible for octets of value 255 to appear "spontaneously" when using multi-octet or non-8-bit characters. All octets of value 255 (other than IAC) MUST be quoted to conform with TELNET requirements. This applies even to octets within a table, or text in a multi-octet character set.]

IAC SB CHARSET TTABLE-ACK IAC SE

The sender acknowledges the successful receipt of the translate table. Text messages which follow this response must now be coded in the character set specified as <char set name 2> of the TTABLE-IS message. This message terminates the current CHARSET subnegotiation.

IAC SB CHARSET TTABLE-NAK IAC SE

The sender reports the unsuccessful receipt of the translate table and requests that it be resent. If subsequent transmission attempts also fail, a TTABLE-REJECTED or CHARSET REJECTED message (depending on which side sends it) should be sent instead of additional futile TTABLE-IS and TTABLE-NAK messages.

IAC SB CHARSET TTABLE-REJECTED IAC SE

In response to a TTABLE-IS message, the receiver of the TTABLE-IS message acknowledges its receipt and indicates it is unable to handle it. This message terminates the current CHARSET subnegotiation.

Any system which supports the CHARSET option MUST fully support the CHARSET REQUEST, ACCEPTED, REJECTED, and TTABLE-REJECTED subnegotiation messages. It MAY optionally fully support the TTABLE-IS, TTABLE-ACK, and TTABLE-NAK messages. If it does fully support the TTABLE-IS message, it MUST also fully support the TTABLE-ACK and TTABLE-NAK messages.

3. Default

WON'T CHARSET

DON'T CHARSET

4. Motivation for the Option

Many TELNET sessions need to transmit data which is not in 7-bit ASCII. This is usually done by negotiating BINARY, and using local conventions (or terminal type kluges) to determine the character set of the data. However, such methods tend not to interoperate well, and have difficulties when multiple character sets need to be supported by different sessions.

Many computer systems now utilize a variety of character sets. Increasingly, a server computer needs to document character sets or translate transmissions and receptions using different pairs of character sets on a per-application or per-connection basis. This is becoming more common as client and server computers become more geographically disperse. (And as servers are consolidated into ever-larger hubs, serving ever-wider areas.) In order for files, databases, etc. to contain correct data, the server must determine the character set in which the user is sending, and the character set in which the application expects to receive.

In some cases, it is sufficient to determine the character set of the end user (because every application on the server expects to use the same character set, or because applications can handle the user's character set), but in other cases different server applications expect to use different character sets. In the former case, an initial CHARSET subnegotiation suffices. In the latter case, the server may need to initiate additional CHARSET subnegotiations as the user switches between applications.

At a minimum, the option described in this memo allows both sides to be clear as to which character set is being used. A minimal implementation would have the server send DO CHARSET, and the client send WILL CHARSET and CHARSET REQUEST. The server could then communicate the client's character set to applications using whatever means are appropriate. Such a server might refuse subsequent CHARSET REQUEST messages from the client (if it lacked the ability to communicate changed character set information to applications, for example). Another system might have a method whereby various applications could communicate to the TELNET server their character set needs and abilities, which the server would handle by initiating new CHARSET REQUEST negotiations as appropriate.

In some cases, servers may have a large set of clients which tend to connect often (such as daily) and over a long period of time (such as years). The server administrators may strongly prefer that the servers not do character set translation (to save CPU cycles when serving very large numbers of users). To avoid manually configuring each copy of the user TELNET software, the administrators might prefer that the software supports translate tables. (If the client software received a translate table from the server and stored it, the table would only need to be sent once.)

5. Description of the Option

When the client TELNET program is able to determine the user's character set it should offer to specify the character set by sending IAC WILL CHARSET.

If the server system is able to make use of this information, it replies with IAC DO CHARSET. The client TELNET is then free to request a character set in a subnegotiation at any time.

Likewise, when the server is able to determine the expected character set(s) of the user's application(s), it should send IAC DO CHARSET to request that the client system specify the character set it is using. Or the server could send IAC WILL CHARSET to offer to specify the character sets.

Once a character set has been determined, the server can either perform the translation between the user and application character sets itself, or request by additional CHARSET subnegotiations that the client system do so.

Once it has been established that both sides are capable of character set negotiation (that is, each side has received either a WILL CHARSET or a DO CHARSET message, and has also sent either a DO CHARSET or a WILL CHARSET message), subnegotiations can be requested at any time by whichever side has sent a WILL CHARSET message and also received a DO CHARSET message (this may be either or both sides). Once a CHARSET subnegotiation has started, it must be completed before additional CHARSET subnegotiations can be started (there must never be more than one CHARSET subnegotiation active at any given time). When a subnegotiation has completed, additional subnegotiations can be started at any time.

If either side violates this rule and attempts to start a CHARSET subnegotiation while one is already active, the other side MUST reject the new subnegotiation by sending a CHARSET REJECTED message.

Receipt of a CHARSET REJECTED or TTABLE-REJECTED message terminates the subnegotiation, leaving the character set unchanged. Receipt of a CHARSET ACCEPTED or TTABLE-ACK message terminates the subnegotiation, with the new character set in force.

In some cases, both the server and the client systems are able to perform translations and to send and receive in the character set(s) expected by the other side. In such cases, either side can request that the other use the character set it prefers. When both sides simultaneously make such a request (send CHARSET REQUEST messages), the server MUST reject the client's request by sending a CHARSET REJECTED message. The client system MUST respond to the server's request. (See the CHARSET REQUEST description, above.)

When the client system makes the request first, and the server is able to handle the requested character set(s), but prefers that the client system instead use the server's (user application) character set, it may reject the request, and issue a CHARSET REQUEST of its own. If the client system is unable to comply with the server's preference and issues a CHARSET REJECTED message, the server can issue a new CHARSET REQUEST message for one of the previous character sets (one of those which the client system originally requested). The client system would obviously accept this character set.

While a CHARSET subnegotiation is in progress, data SHOULD be queued. Once the CHARSET subnegotiation has terminated, the data can be sent (in the correct character set).

Note that regardless of CHARSET negotiation, translation only applies to text (not commands), and only occurs when in BINARY mode [4]. If not in BINARY mode, all data is assumed to be in NVT ASCII [1].

Also note that the CHARSET option should be used with the END OF RECORD option [5] for block-mode terminals in order to be clear on what character represents the end of each record.

As an example of character set negotiation, consider a user on a workstation using TELNET to communicate with a server. In this example, the workstation normally uses the Cyrillic (ASCII) character set [2] but is capable of using EBCDIC-Cyrillic [2], and the server normally uses EBCDIC-Cyrillic. The server could handle the (ASCII) Cyrillic character set, but prefers that instead the client system uses the EBCDIC-Cyrillic character set. (This and the following examples do not show the full syntax of the subnegotiation messages.)

CLIENT	SERVER
WILL CHARSET	WILL CHARSET
DO CHARSET	DO CHARSET
CHARSET REQUEST Cyrillic EBCDIC-Cyrillic	
	CHARSET ACCEPTED EBCDIC- Cyrillic

Now consider a case where the workstation can't handle EBCDIC-Cyrillic, but can accept a translate table:

CLIENT	SERVER
WILL CHARSET	WILL CHARSET
DO CHARSET	DO CHARSET
CHARSET REQUEST [TTABLE] 1 Cyrillic	
	CHARSET TTABLE-IS 1 Cyrillic EBCDIC-Cyrillic
CHARSET TTABLE-ACK	

For another example, consider a case similar to the previous case, but now the user switches server applications in the middle of the session (denoted by ellipses), and the new application requires a different character set:

CLIENT	SERVER
WILL CHARSET	WILL CHARSET
DO CHARSET	DO CHARSET
CHARSET REQUEST [TTABLE] 1 Cyrillic EBCDIC-INT	
	CHARSET TTABLE-IS 1 Cyrillic EBCDIC-Cyrillic
CHARSET TTABLE-ACK	
.
	CHARSET REQUEST EBCDIC-INT
CHARSET ACCEPTED EBCDIC-INT	

6. Security Considerations

Security issues are not discussed in this memo.

7. References

- [1] Postel, J. and J. Reynolds, "Telnet Protocol Specification", STD 8, RFC 854, ISI, May 1983.
- [2] Reynolds, J., and J. Postel, "Assigned Numbers", STD 2, RFC 1700, ISI, October 1994.
- [3] Postel, J. and J. Reynolds, "Telnet Option Specifications", STD 8, RFC 855, ISI, May 1983.
- [4] Postel, J. and J. Reynolds, "Telnet Binary Transmission", STD 27, RFC 856, ISI, May 1983.
- [5] Postel, J., "Telnet End-Of-Record Option", RFC 885, ISI, December 1983.
- [6] Postel, J., "Internet Official Protocol Standards", STD 1, RFC 1920, IAB, March 1996.

8. Author's Address

Randall Gellens
Unisys Corporation
25725 Jeronimo Road
Mail Stop 237
Mission Viejo, CA 92691
USA

Phone: +1.714.380.6350
Fax: +1.714.380.5912
EMail: Randy@MV.Unisys.Com

