

C0 and C1 stability for Unicode and 10646

Kent Karlsson
Stockholm
2022-02-11

Individual contribution

Introduction

Lately there have been a number of quite destructive, from a standard and standards use perspective, suggestions and proposals for the so-called C0 and C1 areas for “control codes”. Effectively, some propose to regard these areas as wholly private use, and one can assign whatever ill-designed set of control codes there; some do not even have line feed. This is very strange as this goes against the very grain of Unicode and ISO/IEC 10646. Unfortunately, neither Unicode nor 10646 in any way discourages this. These areas cannot be regarded as some kind of private-use areas. There are many important characters allocated in these areas, that are never expected to change, characters that are either near-universally supported, or has special properties in Unicode (for bidi and line-breaking at least). Just about every process that handles Unicode (or indeed many other encodings) rely on that several of the C0 (in particular) and C1 characters do not suddenly get some other semantics.

No-one is expecting that U+0009 would be anything else than CHARACTER TABULATION, no-one is expecting that U+000A be anything else than LINE FEED, etc., in a Unicode encoding, no matter how weird or old the “original encoding” (to have a mapping to Unicode, perhaps proposing addition of some so-called “printable” characters to the Unicode standard). But there are weird encodings. Teletext has no line feed, no carriage return, no character tabulation, ... Instead, the *entire* C0 (yes, C0, not C1) is repurposed to have graphic characters with control functionality (code page switching, colour change, ...).

True, most C0 and C1 (“control”) characters are today generally uninterpreted. And some of them are really obsolete, especially in a Unicode/10646 context. Like SHIFT IN and SHIFT OUT, which make sense only in a code page switching context and are explicitly forbidden to have any effect by ISO/IEC 10646, or the INFORMATION SEPARATORS, which were intended for a kind of data structure linearisation that is, AFAIK, not used anymore (but they do have non-default properties in the Unicode character database).

There are actually some “private-use” control codes in the C0/C1 area: DEVICE CONTROLS, though DC1 and DC3 are de-facto standardized for flow control in terminal (emulators), leaving DC2 and DC4, and then also PRIVATE USE 1 and 2.

But there are some character encodings that are “a bit crazy” when it comes to control codes. They override all of C0 (or C1) with something that cannot even be regarded as pure control characters. In particular Teletext (but there are apparently others, also apart from EBCDIC encodings). Teletext has in the C0 range “control codes” where most of them serve three purposes at once 1) a graphic

character (usually a SPACE, but under some conditions it could be a “mosaic” character), 2) code page switching often to the default “alphanumeric” code page (set elsewhere in the Teletext protocol), or the “mosaic” code page (but there are some that do other code page switching), 3) sets a colour styling for the text. The Teletext encodings have no CHARACTER TABULATION, HT is just not used. And in Teletext there is no LINE FEED (or any other NLF), the Teletext *protocol* instead has an explicit line number (within the Teletext page) for each line transmitted.

L2/21-235 ‘Proposal to add further characters from legacy computers and teletext’ suggests that:

“Control characters from microcomputer platforms and teletext were also determined to be out of scope for the UCS. These characters were located in what would today be considered the C0 control range (0x00–0x1F) or the C1 control range (0x7F–0x9F). Processes that need to interchange these codes should simply interchange the binary C0 or C1 value, extended to the UCS code space but without further mapping. Emulators should treat these control codes as appropriate for the targeted environment.”

This must be the *absolutely worst* suggestion in the history of Unicode by an unspeakably enormous margin. Literally! Firstly, there is no such thing as an (isolated) target environments in this solar system. In particular, Teletext is not an isolated target system at all; it is very widespread and has connections to web pages and smart phone apps. The suggestion is detrimental to interoperability, for all systems that use Unicode/10646 and to the Unicode/10646 standards themselves. Secondly, no-one, *absolutely no-one*, has “processes” that expects the C0/C1 areas to be private-use areas in any way whatsoever.

Teletext

The suggestion above is an especially bad idea for converting Teletext text to Unicode. Teletext has “control codes” in the C0 area that aren’t quite control codes, they are spacing *graphic* characters, with two other functionalities. They are *usually* rendered as a SPACE, but in some circumstances rendered as a “mosaic” character. In Teletext, all codepages (they are defined in the Teletext standard) are 7-bit using the 8th bit as parity). For instance, 0x0A, which is used for “End box” spacing “control” (usually rendered as a SPACE, and used for subtitling via Teletext, as well as news flash functionality), would be mapped to U+000A, LINE FEED, which absolutely no-one expects to be anything else than LINE FEED in any context. And 0x00 which is used for “Alpha black”, which switches to the primary G0 codepage, sets the foreground colour to black, and is (usually) displayed as a SPACE; 0x00 would be turned into U+0000, NULL, which absolutely no-one expects to be anything else than NULL, which in addition is often used as string terminator; which is inaccurate but extremely widespread. Indeed, the entire C0 is repurposed in Teletext (the one called ESC is not an escape character in the normal sense, it switches to an alternative G0/G2, so it is more like SI).

So, what about LINE FEED (or CARRIAGE RETURN)? Those C0 codes are mapped to something entirely different in Teletext. How represent multiple lines (in a page)? Is there some other code used for LINE FEED or similar? No. In addition to page numbers (and subpage numbers), the Teletext *protocol* has explicit line numbers (within the page). (The Teletext protocol is a quite complicated protocol; not a design I would recommend...)

So, is this a disaster for Teletext when (literally) outside of the box (i.e., TV set)? No, of course not. Nobody, absolutely nobody, follows the disastrous suggestion quoted above. Teletext is displayed via several handfuls of web pages (and smartphone apps). Generally, the Teletext pages are converted to HTML pages, with styled text for the text in the pages; though some use generated images, more and more it seems, but not a technical necessity (but images may currently provide better fidelity than styled HTML text).

This image is taken from one such web page:



Teletext colours

L2/21-235 'Proposal to add further characters from legacy computers and teletext' states:

"4. Teletext. Teletext was a service invented in the United Kingdom in the early 1970s for broadcasting pages of information, generally text and simple block graphics, to analog television receivers via the vertical blanking interval. Teletext found its greatest popularity in Europe, where it was commonplace until the adoption of digital television; almost all analog television sets sold in Europe since the early 1980s had built-in teletext decoders."

While most of that is true, the use of *past tense is not quite accurate*, Teletext is still in use. Nor is the use of "until". Teletext is supported in all manufactured TV sets and "TV boxes", also for digital television. In addition, Teletext pages are commonly converted to web pages (using HTML, usually this was with styled text, rarely with images, but the latter has become more common), and smart phone apps (again using HTML). Teletext pages can also be converted to ECMA-48 (fast-tracked ISO version is ISO/IEC 6429) using control sequences for the colour changes, both the ones inline and the ones out of line (in the Teletext protocol, for additional colours).








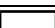
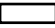









Below is an excerpt from [Teletext] regarding colours (*CLUT* is colour lookup table, 2-bit reference number, *Entry number* (in the CLUT) has a 3-bit reference number). (Leaving out the details of how these can be used in the Teletext protocol, as that is out of scope for this proposal, and quite intricate). Note that there is one added (not part of the quote) column, for colour swatches. Some of the colour entries can be modified (at conformity/presentation level 3.5), and thus corresponds to an ECMA-48 colour palette entry; a converter to HTML/CSS need to keep more directly track of the settings of the


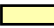












Teletext CLUTs, and that approach of course can be used for a converter to ECMA-48. *Some* of the colours can be invoked by using inline spacing “controls”, but only within the Teletext protocol. We will basically ignore the complications of “conformity levels” (“presentation levels”). It is obvious that the Teletext protocol is overcomplicated, and nothing anyone should try to mimic in any replacement, in particular not in Unicode.

The quote below is included here just to establish that Teletext (at conformity levels greater than 2) has more colours than the ones available via the spacing “control” codes (which aren’t really control codes but serves three functions in one). The additional colours are used via a kind of out of line formatting objects, that reference substrings of the page. At conformity level 3.5 there is a mechanism for redefining the colours in this colour palette. All text colour settings (background and foreground) are reset to default colours (black (or transparent) background, white foreground/text) at end of line. For details, see [Teletext].

12.4 Colour Map

Table 30: Colour Map

CLUT	Entry Number	Default Colour	Default Values			Swatches	Comments
0			R	G	B		
	0	Black	0	0	0		Fixed at Levels 1, 1.5 and 2.5 Re-definable using X/28/4 or M/29/4 at Level 3.5
	1	Red	15	0	0		
	2	Green	0	15	0		
	3	Yellow	15	15	0		
	4	Blue	0	0	15		
	5	Magenta	15	0	15		
	6	Cyan	0	15	15		
	7	White	15	15	15		
1	0	Transparent	-	-	-		Valid at Levels 2.5 and 3.5 (fixed)
	1	Half red	7	0	0		Valid at Levels 2.5 and 3.5 Fixed at Level 2.5 Re-definable using X/28/4 or M/29/4 at Level 3.5
	2	Half green	0	7	0		
	3	Half yellow	7	7	0		
	4	Half blue	0	0	7		
	5	Half magenta	7	0	7		
	6	Half cyan	0	7	7		
	7	Grey	7	7	7		
2	0		15	0	5		
	1		15	7	0		

	2		0	15	7		Valid at Levels 2.5 and 3.5 Re-definable using X/28/0 Format 1 or M/29/0
	3		15	15	11		
	4		0	12	10		
	5		5	0	0		
	6		6	5	2		
	7		12	7	7		
3	0		3	3	3		Valid at Levels 2.5 and 3.5 Re-definable using X/28/0 Format 1 or M/29/0
	1		15	7	7		
	2		7	15	7		
	3		15	15	7		
	4		7	7	15		
	5		15	7	15		
	6		7	15	15		
	7		13	13	13		
NOTE:	The individual R, G, B levels are variable in 16 equally spaced steps. A value of 0 represents zero intensity and a value of 15 (decimal) represents full intensity. The levels are not gamma corrected.						

Teletext styles (other than colour)

Teletext in origin has only fixed-width, upright, “normal” weight style. But at conformity level (“presentation level”) 3.5 there is also support for bold, italics and proportional spacing. There is also support for underlined text, but not quoting that section here.

This text styling is done via out of line formatting objects, that reference a starting point in the line.

01110	Font Style							
	This command specifies the appearance of the text with respect to italics, bold and proportional spacing. The address field defines the column at which the style(s) starts. The effect of this attribute persists to the end of a display row unless overridden by a further Font Style command.							
		D6	D5	D4	D3	D2	D1	D0
		R2	R1	R0	Reserved	Italics	Bold	Proportional Spacing
The functions controlled by bits D0, D1 and D2 are enabled when the bit is set to '1', and cancelled when it is set to '0'. Bits D6 - D4 allow the font style to be extended to the following 0 to 7 (maximum) rows, at the same column positions.								

The quote above is included here just to establish that Teletext (at conformity levels greater than 3) has more styles than just upright, fixed width, “normal” weight. This styling is not mediated by control characters, but out of line with the text, in formatting “objects”. For details see [Teletext].

Teletext protocol and code page switching

Teletext relies heavily on code page switching via the “control” characters (again, three functions in one code), and setting of code pages in the Teletext protocol. All the code pages are 7-bit, using the 8th bit for parity. The G2 page is a “mosaic graphics” page (in two variants, connected and separated “mosaic” pieces, there is a “control” code for switching between them). The “alphanumerics” page is set in the Teletext protocol to one of the “language variants” of ISO/IEC 646, or a Hebrew page, or an Arabic page (well, there are two, primary and secondary, with ESC/SI to switch between them).

Teletext “control codes” to ECMA-48 mapping

Teletext pages are commonly converted to HTML/CSS. This could be done (and was done) by using the styling mechanisms of CSS. But now they are often converted to images (embedded in HTML), presumably for greater fidelity with the “look” of Teletext pages on TV sets.

But an alternative is to convert the Teletext styling to ECMA-48 (ISO/IEC 6429) styling. Most of which do not really require any additions to ECMA-48, except for specifying the colours better. There is one functionality in particular that *does* require a substantial addition to ECMA-48 styling: for the *Start box* and *End box*. These are used for the subtitling (and news flash) functionality of Teletext, and has no close corresponding control sequences in ECMA-48. (An addition is also needed for more directly mapping Teletext’s strange way of setting background colour.) Here is a quote from my proposal to extend the ECMA-48 styling functionality, focussing on converting Teletext text to use ECMA-48 styling:

11. Converting Teletext styling to ECMA-48 styling

Teletext is still in common use around the world, especially for optional subtitling, though the use of Teletext for news pages has declined or even been abandoned. Teletext allows for certain style settings, mostly to do with colour. There are also control bits in the Teletext protocol for handling bold, italic and underline. The colour codes also select which set of characters to use (either alphanumeric or “mosaic” characters, the latter are Teletext specific symbols used to build up larger (and crude...) graphics). Note that the coding for “alphanumeric” characters have multiple variants for G0/G2, various Latin subsets, Greek, Cyrillic, Arabic and Hebrew. These are selected by control bits in the Teletext protocol, in addition **ESC** can be used to switch between a primary and a secondary G0+G2 set. Further, the “mosaic” characters (G1) have two variants, a “contiguous” form (default, and selected by **0x19**) and a “separated” form (selected by **0x1A**). From a modern encoding point of view the latter are separate characters (not yet in Unicode). *Teletext also allows for dynamically defined (“bitmapped”) fonts of unspecified charset.*

The table is a rough sketch only. Teletext overrides (italics, bold, underline, more colours, prop. font, G3 chars.) are not covered in the sketch mapping here, though the functionality is covered. Character sets are specified in control bits in the Teletext protocol, that is not covered here.

Teletext	ECMA-48 (as extended here)
0x00 Alpha black (and sets “alpha mode”, G0) (default)	SP? CSI 30m
0x01 Alpha red (and sets “alpha mode”, G0)	SP? CSI 91m
0x02 Alpha green (and sets “alpha mode”, G0)	SP? CSI 92m
0x03 Alpha yellow (and sets “alpha mode”, G0)	SP? CSI 93m

0x04 Alpha blue (and sets “alpha mode”, G0)	SP? CSI 94m
0x05 Alpha magenta (and sets “alpha mode”, G0)	SP? CSI 95m
0x06 Alpha cyan (and sets “alpha mode”, G0)	SP? CSI 96m
0x07 Alpha white (and sets “alpha mode”, G0)	SP? CSI 37m
0x08 Flash	SP? CSI 5m or SP? CSI 6m
0x09 Steady (default)	CSI 25m SP?
0x0A End box (default)	SP? CSI 112:0m (allow 100% fading)
0x0B Start box (news flashes; also used for subtitles)	SP? CSI 112:1m (block fading)
0x0C Normal size (default)	SP? CSI 73:1:1m
0x0D Double height (note: background & foreground overshadows text on the next line)	SP? CSI 73:2:1m (overshadowing is implementation defined)
0x0E Double width	SP? CSI 73:1:2m
0x0F Double size (note: extends down one “row” (overshadows text on the next line))	SP? CSI 73:2:2m (might not overshadow text on the next line)
0x10 Mosaics black (and sets “mosaics mode”, G1)	SP? CSI 30m
0x11 Mosaics red (and sets “mosaics mode”, G1)	SP? CSI 91m
0x12 Mosaics green (and sets “mosaics mode”, G1)	SP? CSI 92m
0x13 Mosaics yellow (and sets “mosaics mode”, G1)	SP? CSI 93m
0x14 Mosaics blue (and sets “mosaics mode”, G1)	SP? CSI 94m
0x15 Mosaics magenta (and sets “mosaics mode”, G1)	SP? CSI 95m
0x16 Mosaics cyan (and sets “mosaics mode”, G1)	SP? CSI 96m
0x17 Mosaics white (and sets “mosaics mode”, G1)	SP? CSI 37m
0x18 Conceal (ignored on the currently displayed page when user presses a ‘reveal’ button)	CSI 38:8m SP? (not CSI 8m)
0x19 (“mosaics mode”) Contiguous mosaic graphics (default); note: does <i>not</i> unset underlined in “alpha mode”	(converter change) SP?
0x1A (“mosaics mode”) Separated mosaic graphics; note: does <i>not</i> set underlined in “alpha mode”	(converter change) SP?
0x1B Escape (toggles between two G0 code pages)	SP? (converter change)
0x1C Black background (default)	CSI 40m SP?
0x1D New background (foreground to background)	CSI 48:6m SP?
0x1E Hold mosaics (details of hold/release mosaics is beyond the scope of this paper)	(SP? := <mos. char>; SP?
0x1F Release mosaics (details of hold/release mosaics is beyond the scope of this paper)	SP?; SP? := SP
<i>line break</i> (implicit; Teletext has numbered “rows”) (style settings are reset at the beginning of each row)	<i>line break</i> (CRLF/CR/LF), CSI 0m; SP? := SP

As you can see, the “triple functionality” control codes can be converted to using Unicode with ECMA-48 control sequences for the colour changes as well as “subtitle boxing” (with a proposed extension to ECMA-48). There is no need for new control characters or redefining any control characters. The Teletext “control” characters are so deeply rooted in the Teletext *protocol*, which uses code page switching extensively, that they do not even make sense outside of that protocol and must be converted anyway.

I do not know the storage format that broadcaster’s Teletext systems use. But it is surely not in the Teletext *protocol*, but something *converted* to that protocol *when the pages are sent*. The storage

formats may well use “escapes” (in the programming language sense; like \g for alpha green, say). There is no need for “Teletext control codes” in the storage formats. The entire idea of using Teletext control codes outside of the Teletext *protocol* is bogus and ill-conceived. But if converted to use ECMA-48 styling, then *that* is a possible storage (archival) format.

Teletext and optional subtitling (and news flashes)

Subtitles (in Teletext) are rarely, if ever (haven’t seen that), converted to web pages. Subtitling needs timing information (when to show the subtitles). This is not part of the Teletext protocol but must be part of the storage format for Teletext subtitling (again, I’m not familiar with the storage formats used, but they (or it) are likely proprietary for a particular system).

News flashes, a good idea but haven’t seen that used, are more one-off (show now, ..., stop showing), not related to running time of the TV program. But on the Teletext protocol side news flashes would use the same “box” controls as subtitling (plus that there are protocol bits saying whether the page a “normal” page, a subtitling page, or a news flash page).

Teletext and digital TV broadcast/multicast (DVB, IP-TV)

The original Teletext data was transmitted in the so-called synch lines of the analogue TV signal. For the digital TV broadcasts (and multicasts, a.k.a. IP-TV), the Teletext data is embedded in the digital TV signal, but not in the synch lines as there are no sync lines, indeed no scan lines at all. See reference [Teletext in DVB] for details of that embedding of Teletext in DVB (Digital Video Broadcasting).

Teletext references

Teletext is not dead, quite yet. But it is admittedly on the decline. There are few and fewer broadcasters that provide any Teletext news service. Teletext is still used for optional subtitling, for accessibility reasons. It may continue to be in use for optional subtitling for several years to come. In DVB there is a newer mechanisms for subtitling (apart from Teletext), using images, so any fonts need only be on the broadcaster side, not the TV set side. But that has still to catch on, it seems.

[Teletext]	<i>Enhanced Teletext specification</i> , ETSI EN 300 706 V1.2.1, https://www.etsi.org/deliver/etsi_en/300700_300799/300706/01.02.01_60/en_300706v010201p.pdf , 2003.
[Teletext in DVB]	Digital Video Broadcasting (DVB); <i>Specification for conveying ITU-R System B Teletext in DVB bitstreams</i> , ETSI EN 300 472 V1.4.1, https://www.etsi.org/deliver/etsi_en/300400_300499/300472/01.04.01_60/en_300472v010401p.pdf , 2017.
[TeletextWebList]	http://teletext.mb21.co.uk/live.shtml (many links are now defunct, due to the decline of use of Teletext for news pages).
[WikiTeletext]	https://en.wikipedia.org/wiki/Teletext , https://en.wikipedia.org/wiki/Teletext_character_set .

Some legacy encodings for C0/C1 according to INTERNATIONAL REGISTER OF CODED CHARACTER SETS TO BE USED WITH ESCAPE SEQUENCES

The register summary is given at: https://www.itscj-ipsj.jp/custom_contents/cms/linkfile/ISO-IR.pdf.

This does not cover all “old” control codes. For instance, “PETSCII”, “ATASCII” and EBCDIC are not covered (but see below). But for these, much of the special control codes have to do with colour changes and arrow key presses, which is all covered by standard control sequences in ECMA-48.

Here is a walk-through of the registered C0 and C1 variants for ISO/IEC 2022.

- **1 C0 Set of ISO 646** <https://www.itscj-ipsj.jp/ir/001.pdf> 4/0 4.2
A bit cryptic, but this should probably be seen as the “normal” C0, but with one major typo: the first ETX should be STX.
- **7 NATS, C0 Set** <https://www.itscj-ipsj.jp/ir/007.pdf> 4/1 4.9
Intended for “newspaper typography”. Seems to be mappable to ECMA-48 styling (plus word delete...).
- **26 IPTC, C0 Set** <https://www.itscj-ipsj.jp/ir/026.pdf> 4/3 4.14
“International press council”; has similarities with entry 7. Seems to be mappable to ECMA-48 styling (plus word delete...), plus maybe some APC...ST (the description in the registration is incomplete), which is private-use in ECMA-48.
- **36 C0 Set of ISO 646 with SS2 instead of IS4** <https://www.itscj-ipsj.jp/ir/036.pdf> 4/4 4.18
The title says it all. Easily mapped to ECMA-48 standard C0/C1. However, that is for SINGLE SHIFT, and thus not relevant for Unicode/10646, only for conversion (maybe to Unicode).
- **48 INIS, Control Set** <https://www.itscj-ipsj.jp/ir/048.pdf> 4/2 4.29
“International Nuclear Information System”/“Bibliographic data interchange”; Just declaring that they (supposedly) only use ESC, GS (IS3) and RS (IS2) out of the “normal” C0 codes.
- **74 C0 for Japanese standard JIS C 6225-1979** <https://www.itscj-ipsj.jp/ir/074.pdf> 4/6 4.48
“Normal” C0 except that IS4 is replaced by CEX (CONTROL EXTENSION). CEX should probably be mapped either to DLE or some ECMA-48 control sequence (depending on what CEX is...).
- **104 Minimum C0 Set for ISO 4873** <https://www.itscj-ipsj.jp/ir/104.pdf> 4/7 4.62
Ordinary C0 limited to ESC...
- **106 Teletex primary set of Control Functions CCITT Rec. T.61** <https://www.itscj-ipsj.jp/ir/106.pdf> 4/5 4.68 **(Not to be confused with Teletext.)**
C0 easily mapped to ECMA-48 standard C0/C1. However, that is for LOCKING SHIFTS, and thus not relevant for Unicode/10646, only for conversion (maybe to Unicode).
- **130 C0 Set of ISO 646 without SI and SO ASMO-662 and COMECON ST SEV 358** <https://www.itscj-ipsj.jp/ir/130.pdf> 4/8 4.80
Normal (ECMA-48) C0 but excluding SI and SO (LS1 and LS0). These two controls are basically always uninterpreted today (and must be uninterpreted for Unicode/10646) anyway.
- **132 Primary Control Set of Data Syntax I of CCITT Rec. T.101** <https://www.itscj-ipsj.jp/ir/132.pdf> 4/9 4.85
Easily mappable to ECMA-48 standard control codes and standard control sequences.
- **134 Primary Control Set of Data Syntax II of CCITT Rec. T.101** <https://www.itscj-ipsj.jp/ir/134.pdf> 4/10 4.93
Mostly easily mappable to standard ECMA-48 control codes/sequences, except for CURSOR ON/OFF and REPEAT, which may need some minor extension to ECMA-48 or the use of private use control codes/sequences.
- **135 Primary Control Set of Data Syntax III of CCITT Rec. T.101** <https://www.itscj-ipsj.jp/ir/135.pdf> 4/11 4.97
Easily mappable to ECMA-48 standard control codes and standard control sequences.
- **140 C0 Set of ISO 646 with EM replaced by SS2 - Czechoslovak 4/12 4.110 Standard CSN 369102** <https://www.itscj-ipsj.jp/ir/140.pdf> 4/12 4.110
The title says it all. Easily mapped to ECMA-48 standard C0/C1.

- **40 Additional Control Functions for Bibliographic Use according to DIN 31626**
<https://www.itscj-ipsj.jp/ir/040.pdf> 4/5 4.25
Has a number of control codes for affecting bibliographic sorting of strings. These can be mapped to private-use control sequences as per ECMA-48, or private-use Unicode characters.
- **56 Attribute Control set for UK Videotex British Telecom**
<https://www.itscj-ipsj.jp/ir/056.pdf> 4/0 4.32 (***Videotex was a precursor to Teletext.***)
Neither the 7-bit encoding (using escape sequences) nor the 8-bit encoding of that registration can be used in Teletext (and probably nor in Videotex): there are no escape sequences in Teletext, nor are there any 8-bit encodings so there is no C1 area. These “controls” are also inappropriate for anything outside of the Teletext protocol, since these (like the real C0 in Teletext) assume a particular code page setup particular to Teletext. In addition, the descriptions are wrong (or at least incomplete); these are spacing graphic characters (despite being called control characters), usually rendered as SPACE (this likely holds also for Videotex). And... there are no CSI nor ESC (despite the name) controls in Teletext (and likely likewise for Videotex). So this is just...very wrong.
- **67 Additional Control Functions for Bibliographic Use according to ISO 6630**
<https://www.itscj-ipsj.jp/ir/067.pdf> 4/2 4.38
See registration nr 124 below. Has a number of control codes for affecting bibliographic sorting of strings. These can be mapped to private-use control sequences as per ECMA-48, or private-use Unicode characters.
- **73 Attribute Control Set for Videotex CCITT** <https://www.itscj-ipsj.jp/ir/073.pdf> 4/1 4.42 (***Videotex was a precursor to Teletext.***) In addition to the problems listed with entry 56 (see above), this one is incorrect relative to the (Enhanced) Teletext standard w.r.t. the set of “controls”, even when ignoring the issues listed above (for entry 56). So this is just...very, very wrong (and likely wrong also w.r.t. Videotex).
- **77 C1 Control Set of ISO 6429-1983** <https://www.itscj-ipsj.jp/ir/077.pdf> 4/3 4.55
This seems to be the normal C1 controls.
- **105 Minimum C1 Set for ISO 4873** <https://www.itscj-ipsj.jp/ir/105.pdf> 4/7 4.65
“Normal” C1 limited to SS2 and SS3...
- **107 Teletex Supplementary Set of Control Functions CCITT Rec. T.61**
<https://www.itscj-ipsj.jp/ir/107.pdf> 4/8 4.72 (***(Not to be confused with Teletext.)***)
“Normal” C1 limited to PLU, PLD, and CSI.
- **124 Upward Compatible Version of ISO 6630 (Registration 67)**
<https://www.itscj-ipsj.jp/ir/124.pdf> 4/0 4/2 4.76
Has a number of control codes for affecting bibliographic sorting of strings. These can be mapped to private-use control sequences as per ECMA-48, or private-use Unicode characters.
- **133 Supplementary Control Set of Data Syntax I of CCITT Rec. T.101**
<https://www.itscj-ipsj.jp/ir/133.pdf> 4/4 4.89
(Another registration for Videotex, but different set of controls. See above. This one has likely never been used.)
- **136 Supplementary Control Set of Data Syntax III of CCITT Rec. T.101**
<https://www.itscj-ipsj.jp/ir/136.pdf> 4/6 4.102
(Yet another registration for Videotex, but different set of controls. See above. This one has likely never been used.)

Note that there is no registration for Teletext “controls”; and the registration for “normal” C0 has a major typo.

It is high time to scrap all this for the purposes of ISO/IEC 10646 (which mistakenly still references it). The C0/C1 need to be fixed and fixed to what is in ISO/IEC 6429. Again, this does not mean that implementations suddenly must interpret all or any of C0/C1 characters according to ISO/IEC 6429. All can stay as is in just about all modern implementations w.r.t. conformity to the Unicode standard and the ISO/IEC 10646 standard. Most of this is either completely outdated (since decades), or plain wrong (w.r.t. the relevant standard; Teletext in particular). Data using these registrations may still occur in archival data, but should be converted to follow ISO/IEC 6429 (ECMA-48) with suitable bit padding for UTF-16 and UTF-32, and appropriate bit encoding for C1 in UTF-8. If there is no direct mapping to ECMA-48 control codes or standard control sequences, there are a few control codes that are private use, a large number of control sequences that are private use (and one could use appropriate extensions to the standard ones as well, though that would be a matter of at least proposed standard extensions), and all control strings (APC...ST, OSC...ST, DCS...ST) are all private use. So there is no lack of appropriate ways of representing “odd” (and old) control codes in a standards compliant manner according to ISO/IEC 6429.

Bibliographic sorting controls

While most of the non-ISO/IEC 6429 registered control sets are either defective or misguided, one of the few that may actually have continued interest is the bibliographic sorting one, <https://www.itsci-ipsi.jp/ir/124.pdf>. (There are actually three bibliographic sorting control sets registered, but we will look only at the last one.)

While these are formulated as control codes in the registration, I can see three different reasonable ways of handling them in a Unicode and ECMA-48 context. (Apart from CUS to CGJ and preferring superscript/subscript styling to PLU/PLD.) 1) Use printable Unicode characters that are “extremely” unlikely to occur in a bibliographic registry entry (but one can use a character escaping scheme just in case). That way entry is made easier, one can pick from a small palette of characters, and they are visible through common system fonts. And the characters in question should be bidi neutral and bidi mirrored, so that they work in a bidi environment as well. 2) Use Unicode private use characters. These are “guaranteed” not occur in actual bibliographic entries, and can in some sense be “controls”, but can also be given glyphs in special fonts. But that is also a downside: you need special fonts to make them visible, and they might not work well in a bidi context. 3) Use control sequences. That is closer to the original (control codes), but are still invisible (if displayed properly), but are longer than just one character (so handling them for interpretation is a little bit harder). One can combine these approaches though. Use printable characters for input, then pre-process to private-use characters (and interpret character escapes), and then do the bibliographic sorting. Note that the character choices and control sequence choices here are for illustration of the approaches only, it is not a firm proposal.

Bibl. hex.	Mnem.	Name	printable	private-use	contr.seq.(p.u.)
87	CUS	Close-up for sorting	COMBINING GRAPHEME JOINER (U+034F)		
88	NSB	Non-sorting characters begin	⌞ (U+2770)	U+F80A	CSI 1:1s
89	NSE	Non-sorting characters end	⌟ (U+2771)	U+F80B	CSI 1:0s
8B	PLD	Partial line down	PLD (U+008B) (But better to use ECMA-48 superscript/subscript styling per prop. extension; CSI 56:1m, ..., CSI 56:0m; CSI 56:=1m, ..., CSI 56:0m)		
8C	PLU	Partial line up	PLU (U+008C) (But better to use ECMA-48 superscript/subscript styling per prop. extension; CSI 56:1m, ..., CSI 56:0m; CSI 56:=1m, ..., CSI 56:0m)		

Bibl. hex.	Mnem.	Name	printable	private-use	contr.seq.(p.u.)
91	EAB	Embedded annotation begin	⌈ (U+27EC)	U+F800	CSI 2:1s
92	EAE	Embedded annotation end	⌋ (U+27ED)	U+F801	CSI 2:0s
95	SIB	Sorting interpolation begin	{ (U+2983)	U+F802	CSI 3:1s
96	SIE	Sorting interpolation end	} (U+2984)	U+F803	CSI 3:0s
97	SSB	Secondary sorting value begin	(((U+2985)	U+F804	CSI 4:1s
98	SSE	Secondary sorting value end) (U+2986)	U+F805	CSI 4:0s
9C	KWB	Keyword begin	⌵ (U+2989)	U+F806	CSI 5:1s
9D	KWE	Keyword end	⌶ (U+298A)	U+F807	CSI 5:0s
9E	PSB	Permutation string begin	⌷ (U+2991)	U+F808	CSI 6:1s
9F	PSE	Permutation string end	⌸ (U+2992)	U+F809	CSI 6:0s

The registration form gives more information on how these controls (or higher-level protocol, if using printable characters) are to be used. The description in the registration is still incomplete and leaves much of the specification on how these controls affect the sorting to be found elsewhere (in particular ISO 6630:1986 *Documentation — Bibliographic control characters*).

EBCDIC control codes

Also EBCDIC control codes can be reasonably converted to ECMA-48 control codes and control sequences. But we need a few extensions to ECMA-48 for subscript/superscript styling and line indent. There are also three control codes with S/360 special meaning; below I've used private use control sequences for mapping those.

The table below is mostly based on <https://en.wikipedia.org/wiki/EBCDIC>. There appears to have been some variation (both in time and space), which are not covered here.

EBCDIC hex	Name	Mnem.	Description	EBCDIC to ECMA-48 (preliminary)	ICU
00	null	NUL		NULL (U+0000)	<U0000>
01	start of heading	SOH		SOH (U+0001)	<U0001>
02	start of text	STX		STX (U+0002)	<U0002>
03	end of text	ETX		ETX (U+0003)	<U0003>
04	Select	SEL	Device control character taking a single-byte parameter.	DLE (...) (U+0010, ...) (DLE is de facto deprecated since ECMA-24 is withdrawn)	<U009C>
05	horizontal tab.	HT		HT (U+0009)	<U0009>
06	Required New Line	RNL	Line-break resetting Indent Tab mode	CSI 75:0:0m, CSI 76:0:0m, NEL	<U0086>
07	delete	DEL		DEL (U+007F)	<U007F>

EBCDIC hex	Name	Mnem.	Description	EBCDIC to ECMA-48 (preliminary)	ICU
08	Graphic Escape	GE	Non-locking shift that changes the interpretation of the following byte.	SS2 (U+008E) (Must be uninterpreted in Unicode; though this control code should be interpreted by a converter to Unicode)	<U0097>
09	Superscript	SPS	Begin superscript or undo subscript. [toggles]	CSI 56:1m, CSI 56:0m (ECMA-48 as extended in my proposal to update the styling part of ECMA-48; normally ECMA-48 control sequences do not use toggling)	<U008D>
0A	Repeat	RPT	Switch to an operation mode repeating a print buffer	DLE (...) (U+0010, ...) (DLE is de facto deprecated since ECMA-24 is withdrawn)	<U008E>
0B	vertical tabulation	VT		VT (U+000B)	<U000B>
0C	form feed	FF		FF (U+000C)	<U000C>
0D	carriage return	CR		CR (U+000D)	<U000D>
0E	shift out	SO		SO (U+000E) (Must be uninterpreted in Unicode, though this control code should be interpreted by a converter to Unicode)	<U000E>
0F	shift in	SI		SI (U+000F) (Must be uninterpreted in Unicode, though this control code should be interpreted by a converter to Unicode)	<U000F>
10	data link escape	DLE		DLE (U+0010) (DLE is de facto deprecated since ECMA-24 is withdrawn)	<U0010>
11	device control 1	DC1		DC1/XON (U+0011)	<U0011>
12	device control 2	DC2		DC2 (U+0012)	<U0012>
13	device control 3	DC3		DC3/XOFF (U+0013)	<U0013>
14	Restore, Enable Presentation	RES/ENP	Resume output (after BYP/INP)	XON/DC1 (U+0011) (dup.)	<U009D>
15	New Line	NL	Line break.	NEL (or LF , or CRLF)	<U0085>
16	backspace	BS		BS (U+0008)	<U0008>
17	Program Operator Communication	POC	Followed by two one-byte operators that identify the specific function, for example a light or function key.	CSI n SP W (FNK) (The following two bytes need be converted to a suitable <i>n</i>)	<U0087>
18	cancel	CAN		CAN (U+0018)	<U0018>
19	end of medium	EM		EM (U+0019)	<U0019>
1A	Unit Backspace	UBS	A fractional backspace.	CSI y (a private-use control sequence)	<U0092>
1B	Customer Use One	CU1		PU1 (U+0091)	<U008F>
1C	Interchange File Separator	IFS		FS/IS4 (U+001C)	<U001C>

EBCDIC hex	Name	Mnem.	Description	EBCDIC to ECMA-48 (preliminary)	ICU
1D	Interchange Group Separator	IGS		GS/IS3 (U+001D)	<U001D>
1E	Interchange Record Separator	IRS		RS/IS2 (U+001E)	<U001E>
1F	Interchange Unit Separator	IUS		US/IS1 (U+001F)	<U001F>
20	Digit Select	DS	Used by S/360 CPU edit (ED) instruction	CSI 1z (a private-use control sequence)	<U0080>
21	Start of Significance	SoS	Used by S/360 CPU edit (ED) instruction. (Note: different from ISO/IEC 6429's SOS.)	CSI 2z (a private-use control sequence)	<U0081>
22	Field Separator	FDS	Used by S/360 CPU edit (ED) instruction.	CSI 3z (a private-use control sequence)	<U0082>
23	Word Underscore	WUS	Underscores the immediately preceding word.	CSI 4m, CSI 24m (surrounding the "immediately preceding word")	<U0083>
24	Bypass, Inhibit Presentation	BYP/INP	De-activates output, i.e., ignores all graphical characters and control characters besides transmission control codes and RES/ENP, until the next RES/ENP.	(Not sure the description is correct...) XOFF/DC3 (U+0013) (dup.) (or just throw away text: DC4 (? dup.))	<U0084>
25	line feed	LF		LF (U+000A)	<U000A>
26	end of transmission block	ETB		ETB (U+0017)	<U0017>
27	escape	ESC		ESC (U+001B)	<U001B>
28	Set Attribute	SA	Marks the beginning of a fixed-length device specific control sequence. Deprecated in favour of CSP.	CSI ... or DLE ...	<U0088>
29	Start Field Extended	SFE	Marks the beginning of a variable-length device specific control sequence. Deprecated in favour of CSP.	CSI ... or DLE ...	<U0089>
2A	Set Mode, Switch	SM/SW	Device specific control that sets a mode of operation, such as a buffer switch.	DLE (...) (U+0010, ...) (DLE is de facto deprecated since ECMA-24 is withdrawn)	<U008A>

EBCDIC hex	Name	Mnem.	Description	EBCDIC to ECMA-48 (preliminary)	ICU
2B	Control Sequence Prefix	CSP	Marks the beginning of a variable-length device specific control sequence. Followed by a class byte specifying a category of control function, a count byte giving the sequence length (including count and type bytes, but not the class byte or initial CSP), a type byte identifying a control function within that category, and zero or more parameter bytes.	CSI (U+009B) (The description (as per Wikipedia) says "device specific"; however, the description (in Wikipedia) does not give enough detail to specify exact conversions, or even which functions may be specified, so it is hard to say how to interpret the rest of the control sequence. A more detailed analysis, given more detailed information, is needed to give a full conversion specification)	<U008B>
2C	Modify Field Attribute	MFA	Marks the beginning of a variable-length device specific control sequence. Deprecated in favour of CSP.	CSI ... or DLE ...	<U008C>
2D	enquiry	ENQ		ENQ (U+0005)	<U0005>
2E	acknowledge	ACK		ACK (U+0006)	<U0006>
2F	bell	BELL	[The device control, not the bell symbol]	BELL (U+0007)	<U0007>
30	(reserved)		Reserved for future use by IBM	Conversion error: SUB	<U0090>
31	(reserved)		Reserved for future use by IBM	Conversion error: SUB	<U0091>
32		SYN	synchronous idle	SYN (U+0016)	<U0016>
33	Index Return	IR	Either move to start of next line (see also NL), or terminate an information unit (see also IUS/ITB).	IND (U+0084)?	<U0093>
34	Presentation Position	PP	Followed by two one-byte parameters (firstly function, secondly number of either column or line) to set the current position.	CSI <i>n;mH</i>	<U0094>
35	Transparent	TRN	Followed by one byte parameter that indicates the number of bytes of transparent data that follow.	CSI 8m ... CSI 28m (Enclosing the characters that should be transparent. Here I assume that "transparent" means "transparent colour". But it might mean "non-character data" in some way. In the latter case, one can use SOS...ST combined, if needed, with base-64 encoding of the non-character data.)	<U0095>
36	Numeric Backspace	NBS	Move backward the width of one digit.	BS (U+0008) dup., or CSI 1y	<U0096>
37	end of transmission	EOT		EOT (U+0004)	<U0004>

EBCDIC hex	Name	Mnem.	Description	EBCDIC to ECMA-48 (preliminary)	ICU
38	Subscript	SBS	Begin subscript or undo superscript. [toggles]	CSI 56:=1m, CSI 56:0m (ECMA-48 as extended in my proposal to update the styling part of ECMA-48; normally ECMA-48 control sequences do not use toggling)	<U0098>
39	Indent Tab	IT	Indents the current and all following lines, until RNL or RFF is encountered.	CSI 75:0:3m, CSI 76:0:3m (ECMA-48 as extended in my proposal to update the styling part of ECMA-48)	<U0099>
3A	Required Form Feed	RFF	Page-break resetting Indent Tab mode.	CSI 75:0:0m, CSI 76:0:0m, CR, FF	<U009A>
3B	Customer Use Three	CU3		PU2 (U+0092)	<U009B>
3C	device control 4	DC4		DC4 (U+0014)	<U0014>
3D	negative acknowledge	NAK		NAK (U+0015)	<U0015>
3E	(reserved)		Reserved for future use by IBM	Conversion error: SUB	<U009E>
3F	substitute	SUB		SUB (U+001A)	<U001A>
40	space	SP		SP (U+0020)	<U0020>
41	required space	RSP		NBSP (U+00A0)	missing
CA	Syllable hyphen	SHY		SHY (soft hyphen) (U+00AD)	missing
E1	Numeric space	NSP		FIGURE SPACE (U+2007)	missing
FF	Eight Ones	EO	Used as filler	DEL (U+007F); dup.	<U009F>

For ICU all mappings to C0 (for 21 years now) are OK w.r.t. ECMA-48. However, mappings to C1 are botched, with the exception of NEL. In addition, three “semi-control” codes are not mapped in ICU.

ATASCII, PETSCII, ISCII-ATR

The character encoding for many old computer systems do **not** follow the ISO/IEC 2022 architecture, and their control codes do not follow ISO/IEC 6429 and sometimes allocate control codes in what would otherwise be “graphic” areas. However, many of the “special” control codes (regardless of allocation) can mostly easily be mapped to standard control sequences in ISO/IEC 6429. A few may need to be mapped to some minor extensions or the use of private use control sequences.

Several of the ATASCII and PETSCII control codes are for arrow keys (there are ISO/IEC 6429 standard control sequences for that), or (foreground) colour changes (there are ISO/IEC 6429 standard control sequences for that). So much, if not all, of those legacy control codes can be mapped to ISO/IEC 6429 standard control sequences. (And mapped back, if needed.)

Below are mappings for ATASCII controls, PETSCII controls (varies per version), and ISCII styling controls (but not the C0/C1 controls in ISCII).

This should be enough to show the principle not only for these legacy encodings, but also for other “legacy” control codes, and how they should be handled (by interested parties, **not** to be handled or hosted by Unicode consortium).

The UTC, however, **must avoid supporting, and indeed deter**, disastrously destructive approaches such as:

“Control characters from microcomputer platforms and teletext were also determined to be out of scope for the UCS. These characters were located in what would today be considered the C0 control range (0x00–0x1F) or the C1 control range (0x7F–0x9F). Processes that need to interchange these codes should simply interchange the binary C0 or C1 value, extended to the UCS code space but without further mapping. Emulators should treat these control codes as appropriate for the targeted environment.”.

... and instead encourage doing mappings as below (and above). (Again, not suggesting that Unicode consortium actually define the mappings, nor host such mapping tables. Just to say, in short: “this is how to do it, details are left to interested parties”.)

ATASCII control characters to ECMA-48 mapping (and, in principle, backmapping)

ATASCII Hex.	Function	ECMA-48 (preliminary) based on https://en.wikipedia.org/wiki/ATASCII
00		(NULL)
1B	Escape key	ESC (U+001B) or, rather, INT - INTERRUPT (ESC a)
1C	Cursor Up	CUU - CURSOR UP (CSI A , using default arg.: 1)
1D	Cursor Down	CUD - CURSOR DOWN (CSI B , using default arg.: 1)
1E	Cursor Left	CUB - CURSOR LEFT (CSI D , using default arg.: 1)
1F	Cursor Right	CUF - CURSOR RIGHT (CSI C , using default arg.: 1)
7D	Clear Screen	RIS - RESET TO INITIAL STATE (ESC c)
7E	Delete	DEL (U+007F)
7F	Tab	HT (U+0009)
9B	End of line	CR (U+000D) or LF (U+000A) or even NEL (U+0085)
9C	Delete Line	DL - DELETE LINE (CSI M , using default arg.: 1)
9D	Insert Line	IL - INSERT LINE (CSI L , using default arg.: 1)
9E	Clear Tab stop	TBC – (CHARACTER) TABULATION CLEAR (CSI g)
9F	Set Tab stop	HTS - CHARACTER TABULATION SET (U+0088)
FD	Buzzer	BELL (U+0007)
FE	Delete Character	ECH - ERASE CHARACTER (CSI X , using default arg.: 1)
FF	Insert Character	ICH - INSERT CHARACTER (CSI (, using default arg.: 1)

Note that the three last ones are not in the C0/C1 ranges.

PETSCII control characters to ECMA-48 mapping, C0 part

PETSCII hex.	E-48	C esc.	PETSCII to ECMA-48 mapping, <i>per version</i> (preliminary) based on https://en.wikipedia.org/wiki/PETSCII		
00	NUL	\0	(NULL)		
01	SOH		---	---	
02	STX		---/ul.on	---	Ul.on. CSI 4m
03	ETX		Stop	ESC a (INT - INTERRUPT)	
04	EOT		---	---	
05	ENQ		White	CSI 37m	
06	ACK		---	---	
07	BEL	\a	BELL	---	BELL

PETSCII hex.	E-48	C esc.	PETSCII to ECMA-48 mapping, <i>per version</i> (preliminary) based on https://en.wikipedia.org/wiki/PETSCII			
08	BS	<code>\b</code>	Sh.dis./---	CSI =0:0 SP W (new)		---
09	HT	<code>\t</code>	Sh.en./HT	CSI =0:1 SP W (new)		HT
0A	LF	<code>\n</code>	---/LF	---		LF (dupl.) (or IND)
0B	VT	<code>\v</code>	Shift en.	---		CSI =0:1 SP W (new)
0C	FF	<code>\f</code>	Shift dis.	---		CSI =0:0 SP W (new)
0D	CR	<code>\r</code>	CR	CR		
0E	SO		Txt.mode	SO		
0F	SI		---/blink on	---		blink on: CSI 5m
10	DLE		---	---		
11	DC1/XON		Cur. down	CUD - CURSOR DOWN (CSI B, using default arg.: 1)		
12	DC2		Reverse on	CSI 7m		
13	DC3/XOFF		Home	CSI H (using default args. 1;1)		
14	DC4		Delete	DEL		
15	NAK		---	---		
16	SYN		---	---		
17	ETB		---	---		
18	CAN		---/ Tab.set/clr.	---	Toggle: HTS (U+0088), TBC: CSI 0g. ECMA-48 (fortunately) does not have toggling tab set/unset, so we need to use a private use control sequence, like CSI .	
19	EM		---	---		
1A	SUB		---	---		
1B	ESC	<code>\e</code>	ESC	---	(ESC)	
1C	FS/IS4		Red	CSI 91m		
1D	GS/IS3		Cur. right	CUF - CURSOR RIGHT (CSI C, using default arg.: 1)		
1E	RS/IS2		Green	CSI 92m		
1F	US/IS1		Blue	CSI 94m		

It is not clear what Shift-enable and Shift-disable did (Shift-lock/Caps-lock, or actually Shift?).

The PETSCII colours to SGR control sequences may need adjusting.

PETSCII control characters to ECMA-48 mapping, C1 part

PETSCII hex.	E-48	E-48 ESC+	PETSCII to ECMA-48 mapping, per version (preliminary) based on https://en.wikipedia.org/wiki/PETSCII			
80	PAD	@	---	---		
81	HOP	A	Ornge/d.prpl	CSI 99m (new)		CSI 35m
82	BPH	B	Blink on/ ul.off	---	blink on: CSI 5m	UI.off: CSI 24m
83	NBH	C	Run	OSC run ST		
84	IND	D	Blink off/---	---	blink off: CSI 25m	---
85	NEL	E	F1	CSI 1 SP W FNK (FUNCTION KEY)		

PETSCII hex.	E-48	E-48 ESC+	PETSCII to ECMA-48 mapping, per version (preliminary) based on https://en.wikipedia.org/wiki/PETSCII	
86	<u>SSA</u>	F	F2	CSI 2 SP W FNK (FUNCTION KEY)
87	<u>ESA</u>	G	F3	CSI 3 SP W FNK (FUNCTION KEY)
88	HTS	H	F4	CSI 4 SP W FNK (FUNCTION KEY)
89	HTJ	I	F5	CSI 5 SP W FNK (FUNCTION KEY)
8A	VTs	J	F6	CSI 6 SP W FNK (FUNCTION KEY)
8B	PLD	K	F7	CSI 7 SP W FNK (FUNCTION KEY)
8C	PLU	L	F8	CSI 8 SP W FNK (FUNCTION KEY)
8D	RI	M	LF	LF (U+000A, or IND or NEL)
8E	SS2	N	Graphics	SI
8F	SS3	O	---/Blink off	--- blink off: CSI 25m
90	DCS	P	Black	CSI 30m
91	<i>PU1</i>	Q	Curs. up	CUU - CURSOR UP (CSI A , using default arg.: 1)
92	<i>PU2</i>	R	Reverse off	CSI 27m
93	STS	S	CLR	ESC c (RIS - RESET TO INITIAL STATE)
94	CCH	T	Insert	CSI @ (ICH - INSERT CHARACTER)
95	MW	U	Brown/ dark yellow	CSI 38:2::165:42:42m, ...
96	SPA	V	Pink/ yllw-grn/ lght red	CSI 38:2::255:192:203m, ...
97	EPA	W	Dark grey/ pink/ dark cyan/ light grey	CSI 90m (new), ...
98	SOS	X	Med.grey/ blue-green	CSI 98m (new), ...
99	SGC	Y	Lght.green/ lght blue	CSI 38:2::152:251:152m, ...
9A	SCI	Z	Lght.blue/ dark blue	CSI 38:2::135:206:235m, ...
9B	CSI	[Lght.grey/ dark grey	CSI 97m (new), ...
9C	ST	\	Purple	CSI 95m
9D	OSC]	Curs. left	CUB - CURSOR LEFT (CSI D , using default arg.: 1)
9E	PM	^	Yellow	CSI 93m
9F	APC	_	Cyan	CSI 96m

The colour settings are for foreground colour. Background colour is set by other means.

The PETSCII colours to SGR control sequences may need adjusting to get the colours right. In addition several of the colours vary a bit too much per version (omitting the details here). Note that colours, underline, and arrow keys per ECMA-48 are commonly supported in modern terminal emulators; but ECMA-48 is not limited to terminal emulators by any means.

ISCII *styling* controls to ECMA-48 mapping (note: ATR (0xEF) is not in the C0/C1 ranges)

ISCII (all these ATR toggle)	ECMA-48 (as extended here; note: no toggling)
ATR BLD (= ATR 0) (bold)	CSI 1m or CSI 1:3m, CSI 22m
ATR ITA (= ATR 1) (italic)	CSI 3m, CSI 23m
ATR UL (= ATR 2) (underline)	CSI 4m, CSI 24m
ATR EXP (= ATR 3) (expand)	CSI 73:1:2m, CSI 73:1:1m
ATR HLT (= ATR 4) (highlight/bold)	CSI 1m or CSI 1:3m, CSI 22m
ATR OTL (= ATR 5) (outline)	CSI 80:6:5;38:6m (setting the all-around 0,05 em shadow colour to be the current text colour, and then setting the text colour to be the same as the (opaque) background colour), CSI 38:7;89m (if not combined with SHD)
ATR SHD (= ATR 6) (shadow)	CSI 81:45:6:10;80:6:5;38:6m (extra shadow to the ‘south-east’), CSI 38:7;89m (if not combined with OTL)
ATR TOP (= ATR 7) (top half of double height glyphs)	CSI 73:2:1;113:1m, CSI 73:1:1;113:0m
ATR LOW (= ATR 8) (low half of double height glyphs)	CSI 73:2:1;113:2m, CSI 73:1:1;113:0m
ATR DBL (= ATR 9) (double size)	CSI 73:2:2m (combined with ATR TOP and ATR LOW one can get CSI 73:4:2m and combined with ATR EXP also CSI 73:4:4m) or maybe CSI 72:0:2m (not clear from ISCII), CSI 73:1:1m or maybe CSI 72:0:0?5m
<i>line break char</i> (in ISCII all style settings are auto-reset at the beginning of each line)	<i>line break char</i> , CSI 0m

(C0/C1 not covered here, nor EXT (= SS2?), just the styling controls.) (“Here” refers to my proposal to update/extend styling controls for ECMA-48.)

Note again that ATR (= 0xEF) is outside of C0/C1.

Recommendations

The goal of this paper is not to preserve Teletext, even though quite a lot of is has been devoted to certain aspects of Teletext. But that has been in order to give a background for why it is a bad idea to regard C0/C1 areas as some kind of private use area(s), like suggested in L2/21-235 ‘Proposal to add further characters from legacy computers and teletext’ which suggests that:

“Control characters from microcomputer platforms and teletext were also determined to be out of scope for the UCS. These characters were located in what would today be considered the C0 control range (0x00–0x1F) or the C1 control range (0x7F–0x9F). Processes that need to interchange these codes should simply interchange the binary C0 or C1 value, extended to the UCS code space but without further mapping. Emulators should treat these control codes as appropriate for the targeted environment.”

This must not happen. It is a senseless idea for Teletext (as explained at length above) and is a very, very bad idea also for other, less extreme, settings. The C0/C1 needs to be fixed in semantics. That does not mean that Teletext and other “old” systems are left stranded. For Teletext it is perfectly possible to convert to another styling system that is valid outside of the Teletext protocol. For other odd-ball old control codes one should map each to the closest corresponding ISO/IEC 6429 control

code (there are even 4 private-use ones), escape sequence, control sequence (CSI.....; there are plenty of private use ones), or control string (often application program control string; APC....ST; all of which is private-use).

C0/C1 should be protected by conformity clauses, stability clauses, and a total rewrite of TUS section 23.1 (*Control Codes*; see below for a suggested rewrite) of the Unicode standard, and a rewrite of at least section 13.4 (Identification of control function set) of ISO/IEC 10646 (nobody, absolutely nobody, implements *those* escape sequences (for C0/C1) anyway).

C0/C1 needs to be as stable as the ASCII part of the Unicode and ISO/IEC 10646 standards. This does not mean that implementations must start supporting lots of control codes not supported today. Many of them will simply be uninterpreted, as they are today. Note that not even A-Z, a-z, 0-9 are required to be supported by conforming implementations. (All do, but that is for practical reasons, not conformity reasons.)

Requesters that want “new” control codes for handling some (presumably old) character sets when mapped to Unicode, should be directed to ISO/IEC 6429, which in addition to many standard control sequences, also has room for private-use control sequences and control strings. There is no need to add any “new” control code to Unicode for any “old” character set, while still having C0/C1 areas stable.

No new control codes in Unicode or ISO/IEC 10646

There is no need to encode any new control codes in Unicode. ISO/IEC 6429 provides plenty of private use possibilities: a few C0/C1 control codes (PUx, DCx) that are private-use, lots of private use control sequences, and all control strings are private-use. In addition, Unicode has private use areas. So, for control codes from legacy character sets that cannot be mapped closely to ISO/IEC 6429 C0/C1 controls codes or control sequences (many can be so mapped), there are plenty of private-use options. Some (in particular, box start/end in Teletext) may need extensions to the SGR functionality of ISO/IEC 6429 which are not formally private-use.

Nit: Teletext has no tabs, and no line feed (lines are handled by having line numbers in the Teletext *protocol*), and overrides the *entire* C0 (yes, C0, not C1, there is no C1 in Teletext) controls with graphic characters (SPACE or a “mosaic” character) that have control functions (often two each).

C0/C1 stability

C0/C1 needs to be stable (except for the handful of private-use control codes), not vary like a private use area, like some have, quite destructively (from a standards and standards use perspective), proposed. Below are some text changes to the Unicode standard text relevant to this as well as some text changes to ISO/IEC 10646 relevant to this. There may be more changes needed.

All Cc, Cf and Zp/Zl characters must be glyphless in normal circumstances (i.e., no font lookup). Note that SPACE is considered to be a glyph here. (Note that this makes Teletext “controls” (which are actually SPACE or a “mosaic” character) moved to C1 unsuitable also under the secondary alternative below). Of course, failure to parse ESC-sequences for C1 controls does not count, since such parsing should not be required. Nor does a “show invisibles” mode (a non-default mode, that still may be “on” by default for uninterpreted control codes or control sequences), which may show some generated glyph or a fallback glyph or special display glyph for Cc, Cf, and Zs/Zp/Zl characters.

Primary alternative

Change the approach for both C0 and C1 from a totally open “may” to a “shall”/“must” for both Unicode and ISO/IEC 10646 to reference ISO/IEC 6429 for the definition of the “control” characters in those areas. This way both C0 & C1 are stable in Unicode/10646 from a standards point of view.

Secondary alternative

Change the approach for C0 from a totally open “may” to a “shall”/“must” for both Unicode and ISO/IEC 10646 to reference ISO/IEC 6429 for the definition of the “control” characters in C0.

Change the approach for C1 from a totally open “may” to a “should” for both Unicode and ISO/IEC 10646 to reference ISO/IEC 6429 for the definition of the “control” characters in C1. However, the ESC-representation for C1 control codes always refer to the ISO/IEC 6429 C1 control codes.

This secondary alternative is passable, though not ideal, for two reasons: 1) C1 “controls” are often left uninterpreted, and 2) ISO/IEC 6429 (ECMA-48) has a mechanism for accessing the “true” C1 controls (according to ISO/IEC 6429) in case there is no C1 area, or, by extension, in case the C1 area is botched (from a ISO/IEC 6429 point of view). This is the default approach already taken in terminal emulators, since the existence of a plain C1 at all is not guaranteed in those contexts. This approach would allow the ICU mappings for EBCDIC and allow the bibliographic C1 control codes (but not the ESC-sequences for those).

Note that the ICU mapping of EBCDIC to C1 controls is **non**-conforming to ISO/IEC 10646 as it is currently written, since that C1 set is not registered.

Proposed text changes, and property changes for the Unicode standard

For the proposed changes below, the primary alternative above is assumed.

Old (<https://www.unicode.org/versions/Unicode14.0.0/UnicodeStandard-14.0.pdf>):

Control Codes. Sixty-five code points (U+0000..U+001F and U+007F..U+009F) are defined specifically as control codes, for compatibility with the C0 and C1 control codes of the ISO/IEC 2022 framework. A few of these control codes are given specific interpretations by the Unicode Standard. (See Section 23.1, Control Codes.)

Proposed (6429 instead of 2022):

Control Codes. Sixty-five code points (U+0000..U+001F and U+007F..U+009F) are defined specifically as control codes, for compatibility with the C0 and C1 control codes of ISO/IEC 6429:1992 *Information technology — Control functions for coded character sets*. A few of these control codes are given specific interpretations by the Unicode Standard. (See Section 23.1, Control Codes.)

Old (<https://www.unicode.org/versions/Unicode14.0.0/UnicodeStandard-14.0.pdf>):

Control Codes

In addition to the special characters defined in the Unicode Standard for a number of purposes, the standard incorporates the legacy control codes for compatibility with the ISO/IEC 2022 framework, ASCII, and the various protocols that make use of control codes. Rather than simply being defined as byte values, however, the legacy control codes are assigned to Unicode code points: U+0000..U+001F, U+007F..U+009F. Those code points for control

codes must be represented consistently with the various Unicode encoding forms when they are used with other Unicode characters. For more information on control codes, see *Section 23.1, Control Codes*.

Proposed (6429 instead of 2022):

Control Codes

In addition to the special characters defined in the Unicode Standard for a number of purposes, the standard incorporates the legacy control codes for compatibility with ISO/IEC 6429:1992 *Information technology — Control functions for coded character sets*, ASCII, and the various protocols that make use of control codes. Rather than simply being defined as byte values, however, the legacy control codes are assigned to Unicode code points: U+0000..U+001F, U+007F..U+009F. Those code points for control codes must be represented consistently with the various Unicode encoding forms when they are used with other Unicode characters. For more information on control codes, see *Section 23.1, Control Codes*.

Old (<https://www.unicode.org/versions/Unicode14.0.0/UnicodeStandard-14.0.pdf>):

23.1 Control Codes

There are 65 code points set aside in the Unicode Standard for compatibility with the C0 and C1 control codes defined in the ISO/IEC 2022 framework. The ranges of these code points are U+0000..U+001F, U+007F, and U+0080..U+009F, which correspond to the 8-bit controls 00_16 to 1F_16 (C0 controls), 7F_16 (delete), and 80_16 to 9F_16 (C1 controls), respectively. For example, the 8-bit legacy control code character tabulation (or tab) is the byte value 09_16; the Unicode Standard encodes the corresponding control code at U+0009.

The Unicode Standard provides for the intact interchange of these code points, neither adding to nor subtracting from their semantics. The semantics of the control codes are generally determined by the application with which they are used. However, in the absence of specific application uses, they may be interpreted according to the control function semantics specified in ISO/IEC 6429:1992.

In general, the use of control codes constitutes a higher-level protocol and is beyond the scope of the Unicode Standard. For example, the use of ISO/IEC 6429 control sequences for controlling bidirectional formatting would be a legitimate higher-level protocol layered on top of the plain text of the Unicode Standard. Higher-level protocols are not specified by the Unicode Standard; their existence cannot be assumed without a separate agreement between the parties interchanging such data.

Representing Control Sequences

There is a simple, one-to-one mapping between 7-bit (and 8-bit) control codes and the Unicode control codes: every 7-bit (or 8-bit) control code is numerically equal to its corresponding Unicode code point. For example, if the ASCII line feed control code (0A_16) is to be used for line break control, then the text “WX<LF>YZ” would be transmitted in

Unicode plain text as the following coded character sequence: <0057, 0058, 000A, 0059, 005A>.

Control sequences that are part of Unicode text must be represented in terms of the Unicode encoding forms. For example, suppose that an application allows embedded font information to be transmitted by means of markup using plain text and control codes. A font tag specified as “^ATimes^B”, where ^A refers to the C0 control code 01_16 and ^B refers to the C0 control code 02_16, would then be expressed by the following coded character sequence: <0001, 0054, 0069, 006D, 0065, 0073, 0002>. The representation of the control codes in the three Unicode encoding forms simply follows the rules for any other code points in the standard:

UTF-8: <01 54 69 6D 65 73 02>

UTF-16: <0001 0054 0069 006D 0065 0073 0002>

UTF-32: <00000001 00000054 00000069 0000006D 00000065 00000073 00000002>

Escape Sequences. Escape sequences are a particular type of protocol that consists of the use of some set of ASCII characters introduced by the *escape* control code, 1B_16, to convey extra-textual information. When converting escape sequences into and out of Unicode text, they should be converted on a character-by-character basis. For instance, “ESC-A” <1B 41> would be converted into the Unicode coded character sequence <001B, 0041>. Interpretation of U+0041 as part of the escape sequence, rather than as latin capital letter a, is the responsibility of the higher-level protocol that makes use of such escape sequences. This approach allows for low-level conversion processes to conformantly convert escape sequences into and out of the Unicode Standard without needing to actually recognize the escape sequences as such.

If a process uses escape sequences or other configurations of control code sequences to embed additional information about text (such as formatting attributes or structure), then such sequences constitute a higher-level protocol that is outside the scope of the Unicode Standard.

Specification of Control Code Semantics

Several control codes are commonly used in plain text, particularly those involved in line and paragraph formatting. The use of these control codes is widespread and important to interoperability. Therefore, the Unicode Standard specifies semantics for their use with the rest of the encoded characters in the standard. Table 23-1 lists those control codes.

Table 23-1. Control Codes Specified in the Unicode Standard		
Code Point	Abbreviation	ISO/IEC 6429 Name
U+0009	HT	character tabulation (tab)
U+000A	LF	line feed
U+000B	VT	line tabulation (vertical tab)
U+000C	FF	form feed
U+000D	CR	carriage return
U+001C	FS	information separator four
U+001D	GS	information separator three
U+001E	RS	information separator two
U+001F	US	information separator one
U+0085	NEL	next line

The control codes in Table 23-1 have the Bidi_Class property values of S, B, or WS, rather than the default of BN used for other control codes. (See Unicode Standard Annex #9, “Unicode Bidirectional Algorithm.”) In particular, U+001C..U+001E and U+001F have the Bidi_Class property values B and S, respectively, so that the Bidirectional Algorithm recognizes their separator semantics.

The control codes U+0009..U+000D and U+0085 have the White_Space property. They also have line breaking property values that differ from the default CM value for other control codes. (See Unicode Standard Annex #14, “Unicode Line Breaking Algorithm.”)

U+0000 null may be used as a Unicode string terminator, as in the C language. Such usage is outside the scope of the Unicode Standard, which does not require any particular formal language representation of a string or any particular usage of null.

Newline Function. In particular, one or more of the control codes U+000A line feed, U+000D carriage return, and the Unicode equivalent of the EBCDIC next line can encode a newline function. A newline function can act like a line separator or a paragraph separator, depending on the application. See Section 23.2, Layout Controls, for information on how to interpret a line or paragraph separator. The exact encoding of a newline function depends on the application domain. For information on how to identify a newline function, see Section 5.8, Newline Guidelines.

Proposed (incl.: HTJ should be treated the same as HT; bidi, lb, and some other prop. fixes):

23.1 Control Codes

There are 65 code points set aside in the Unicode Standard for compatibility with the C0 and C1 control codes defined in ISO/IEC 6429:1992 *Information technology — Control functions for coded character sets*. The ranges of these code points are U+0000..U+001F, U+007F, and U+0080..U+009F, which correspond to the 8-bit controls 00_16 to 1F_16 (C0 controls), 7F_16 (delete), and 80_16 to 9F_16 (C1 controls), respectively. For example, the 8-bit legacy control code *character tabulation* (or tab) is the byte value 09_16; the Unicode Standard encodes the corresponding control code at U+0009.

The Unicode Standard provides for the intact interchange of these code points, neither adding to nor subtracting from their semantics, except as noted below. The semantics of the control codes are generally determined by the application with which they are used. The control function semantics is specified in ISO/IEC 6429:1992. With some exceptions, especially the ones for NLF and tab, many applications do not interpret most of the control codes, and they are then usually rendered as zero-width and glyph-less (something that may cause security issues or strange errors).

In general, the use of control codes constitutes a higher-level protocol and is beyond the scope of the Unicode Standard, except for the control codes mentioned in section *Specification of Control Code Semantics* below.

For example, the use of ISO/IEC 6429 control sequences for controlling bidirectional formatting would be a legitimate higher-level protocol layered on top of the plain text of the Unicode Standard. Note though, that the use of ISO/IEC 6429 control sequences for

controlling bidirectional formatting is incompatible with the use of the Unicode bidi algorithm, a different higher-level protocol for the same purpose using control codes that are not in C0/C1. A more common use of ISO/IEC 6429 control sequences is to specify bold, italic, or underline style, as well as background and foreground colours for the text.

Even if a process does not interpret any of the ISO/IEC 6429 escape sequences, control sequences, or control strings, it may render them (in their entirety) as zero-width and glyph-less (if rendering is at all part of the process), or (in their entirety) as default-ignorable. However, that is not required by the Unicode standard.

Higher-level protocols, except for the Unicode bidi algorithm and certain cursive shaping, are not specified by the Unicode Standard; their existence cannot be assumed without a separate agreement between the parties interchanging such data (and likewise for the Unicode bidi algorithm).

Representing Escape Sequences, Control Sequences and Control Strings

There is a simple, one-to-one mapping between commonly used 7-bit (and 8-bit) control codes (those conforming to ISO/IEC 6429) and the Unicode control codes: every 7-bit (or 8-bit) control code is numerically equal to its corresponding Unicode code point. For example, if the ASCII *line feed* control code (0A_16) is to be used for line break control, then the text “WX<LF>YZ” would be transmitted in Unicode plain text as the following coded character sequence: <0057, 0058, 000A, 0059, 005A>.

Escape sequences and control sequences, as well as control strings, that are part of Unicode text must be represented in terms of the Unicode encoding forms. For example, suppose that an application allows embedded font information to be transmitted by means of markup using plain text and control codes. A font tag specified as “^ATimes^B”, where ^A refers to the C0 control code 01_16 and ^B refers to the C0 control code 02_16, would then be expressed by the following coded character sequence: <0001, 0054, 0069, 006D, 0065, 0073, 0002>. The representation of the control codes in the three Unicode encoding forms simply follows the rules for any other code points in the standard:

UTF-8: <01 54 69 6D 65 73 02>

UTF-16: <0001 0054 0069 006D 0065 0073 0002>

UTF-32: <00000001 00000054 00000069 0000006D 00000065 00000073 00000002>

(Note that this example is an example only; it has no known implementations, nor does it follow ISO/IEC 6429 syntax for control strings (something that does have some implementations). The appropriate control string syntax is APC...ST, where the ... is private use. It could be “**fnt=Times**”). [It would be preferable to use the latter, which is standards based and realistic, rather than something purely invented for this paragraph...i.e. APC **fnt=Times** ST.]

Escape sequences and control sequences. [That subtitle should probably be moved up, or deleted; there is already a title good enough above.] Escape sequences and control sequences are a particular type of protocol that consists of the use of some set of ASCII characters introduced by the *escape* control code, 1B_16, or the *control sequence introducer* control code, to convey extra-textual information. When converting escape sequences into and out of

Unicode text, they should be converted on a character-by-character basis. For instance, “ESC-A” <1B 41> would be converted into the Unicode coded character sequence <001B, 0041>. Interpretation of U+0041 as part of the escape sequence, rather than as latin capital letter A, is the responsibility of the higher-level protocol that makes use of such escape sequences. This approach allows for low-level conversion processes to conformantly convert escape sequences into and out of the Unicode encoding forms without needing to actually recognize the escape sequences as such, provided that the other encoding conforms to ISO/IEC 6429. ISO/IEC 6429 also specifies the general syntax for control strings but does not specify the content of such strings; they are entirely private-use. See example above.

If a process uses escape sequences or other configurations of control code sequences to embed additional information about text (such as formatting attributes or structure), then such sequences constitute a higher-level protocol that is outside the scope of the Unicode Standard.

Legacy charsets often contain control codes that do not directly match any of the control codes in C0/C1 of ISO/IEC 6429. However, often one can find matching standard control sequences in ISO/IEC 6429. For example, ATASCII's 1C_16 (cursor up) corresponds to ISO/IEC 6429's *CSI 1A*, (cursor up 1 step) and PETSCII's 1F_16 (blue) corresponds to ISO/IEC 6429's *CSI 34m* (blue-ish foreground) or maybe better *CSI 94m* (clear blue foreground, a popular extension to ISO/IEC 6429). For Teletext "control" codes, specified in the "Enhanced Teletext" standard, many of them serve three purposes in one, a graphic character, code page switching, and colour change. E.g. Teletext's 00_16 (alpha black), (usually) displays as a SPACE, switches to the current default "alphanumeric" 7-bit codepage, and changes to black foreground colour; despite not being true control characters, they are graphic characters with variable display, they reside in the C0 area of all of the 7-bit codepages of Teletext (aside: Teletext itself has no escape sequences, nor any control sequences, nor any C1 area). Generally, if there is no close match, one may consider using the private-use control codes, the private-use control sequences of ISO/IEC 6429 or even private-use Unicode code points. Unicode does not specify these mappings.

Specification of Control Code Semantics

Several control codes are commonly used in plain text, particularly those involved in line and paragraph formatting. The use of these control codes is widespread and important to interoperability. Therefore, the Unicode Standard specifies semantics for their use with the rest of the encoded characters in the standard. Table 23-1 lists those control codes.

Table 23-1. Control Codes Specified in the Unicode Standard

Code Point	Abbreviation	ISO/IEC 6429 Name
U+0009	HT	character tabulation (tab)
U+000A	LF	line feed
U+000B	VT	line tabulation (vertical tab)
U+000C	FF	form feed
U+000D	CR	carriage return
U+001C	FS	information separator four
U+001D	GS	information separator three
U+001E	RS	information separator two
U+001F	US	information separator one
U+0085	NEL	next line
U+0089	HTJ	character tabulation with justification

The control codes in Table 23-1 have the Bidi_Class property values of S or B, rather than the default of BN used for other control codes. (See Unicode Standard Annex #9, “Unicode Bidirectional Algorithm.”) In particular, U+001C..U+001E and U+001F have the Bidi_Class property values B and S, respectively, so that the Bidirectional Algorithm recognizes their separator semantics.

The control codes U+0009..U+000D and U+0085 have the White_Space property. They also have line breaking property values that differ from the default CM value for other control codes. (See Unicode Standard Annex #14, “Unicode Line Breaking Algorithm.”)

[Line break properties for U+001C..U+001F: Though these characters are basically obsolete, they do have “special” bidi properties. They should have line break properties consistent with the bidi properties. HTJ should have bidi and line break properties consistent with HT.]

U+0000 *null* may be used as a Unicode string terminator, as in C, C++ and several other programming languages. Such usage is outside the scope of the Unicode Standard, which does not require any particular formal programming language representation of a string or any particular usage of *null*.

Newline Function. In particular, one or more of the control codes U+000A *line feed*, U+000D *carriage return*, and the Unicode equivalent of the ISO/IEC 6429 *next line* (often used in EBCDIC based encodings) can encode a newline function. A newline function can act like a *line separator* or a *paragraph separator*, depending on the application. See Section 23.2, Layout Controls, for information on how to interpret a line or paragraph separator. The exact encoding of a newline function depends on the application domain. For information on how to identify a newline function, see Section 5.8, Newline Guidelines. Also U+000B (*line tabulation*) and U+000C (*form feed*) can encode new line functions, similar to that of *line separator*.

Proposed, LineBreak.txt, change to have:

001C.. 001E ; BK	# Cc	[3] <control-001E>.. <control-001F>	[cmp. bidi B]
001F ; BA	# Cc	<control-001F>	[cmp. bidi S]
0084 ; LF	# Cc	<control-0084>	[xterm still interprets this character]
0089 ; BA	# Cc	<control-0089>	[this is a character tabulation variant]

Proposed, UnicodeData.txt, change to have:

000C;<control>;Cc;0; S ;;N;FORM FEED (FF);;;
0084;<control>;Cc;0; S ;;N;; [xterm still interprets this character]
0089;<control>;Cc;0; S ;;N;CHARACTER TABULATION WITH JUSTIFICATION;;;

Proposed text changes for ISO/IEC 10646

Section 12 correctly states:

“Code extension control functions for the ISO/IEC 2022 code extension techniques (such as designation escape sequences, single shift, and locking shift) shall not be used with this coded character set.”

However, section 13.4 incorrectly contradicts that for control characters, harking back to ISO 2022 codepage switching which must not be used with the ISO/IEC 10646 encodings:

“13.4 Identification of control function set

When the escape sequences from ISO/IEC 2022 are used, the identification of each set of control functions (see Clause 12) of ISO/IEC 6429 to be used in conjunction with ISO/IEC 10646 shall be an identifier sequence of the type shown below.

ESC 02/01 04/00 identifies the full C0 set of ISO/IEC 6429

ESC 02/02 04/03 identifies the full C1 set of ISO/IEC 6429

For other C0 or C1 sets, the final octet F shall be obtained from the International Register of Coded Character Sets. The identifier sequences for these sets shall be

ESC 02/01 F identifies a C0 set

ESC 02/02 F identifies a C1 set

If such an escape sequence appears within a code unit sequence conforming to ISO/IEC 2022, it shall consist only of the sequences of bit combinations as shown above.

If such an escape sequence appears within a code unit sequence conforming to this document, it shall be padded in accordance with Clause 12.”

That text needs to be replaced. This leftover from before the synchronization with Unicode has made some think that the “shall not be used with this coded character set” from section 12 does not apply to the C0/C1 areas, which of course it does. Indeed, C0/C1 need be fixed to that of ISO/IEC 6429, by reference to that standard.

Proposed replacement:

13.4 Control function set

The control codes in C0 shall be as defined in ISO/IEC 6429.

The control codes in C1 [shall| should] be as defined in ISO/IEC 6429. The escape sequences designating control codes in C1 shall be referring to C1 control codes as defined in ISO/IEC 6429.

[only needed for the “should” option in the previous paragraph: If the control codes in C1 deviate from ISO/IEC 6429, that shall be clearly documented and such encodings shall not use any of the designations UTF-8, UTF-16BE, UTF-16LE, UTF-32BE or UTF-32LE, and the characters shall be single function control characters without any graphic component.]

As mentioned in section 12: Code extension control functions for the ISO/IEC 2022 code extension techniques (such as designation escape sequences) shall not be used with this coded character set.

There may be some other changes needed as well.

Conclusions

That ISO/IEC and Unicode has not stabilised the C0 and C1 areas, but left them as totally user-defined (a.k.a. private use) is a major standardisation failure. Fortunately, it has gone largely unnoticed (for 30 years!), and just about everyone assume that C0 and C1 are the ECMA-48 (ISO/IEC 6429) ones, and nothing else. Even if most handling of character strings ignores most of ISO/IEC 6429, it is always assumed that HT, VT, FF LF, CR and sometimes NEL are always exist and are at their ECMA-48 positions.

This standardisation flaw is especially embarrassing for Unicode and 10646 due to the high emphasis on stability that these standards feature. Even more embarrassing is the pre-synchronisation leftover in ISO/IEC 10646 that (contradicting a strong requirement not to use it) “sort of” allows for ISO/IEC 2022 style of swapping out C0 or C1 to something else. And Unicode only says “may use ISO/IEC 6429”, which is as far from a requirement to use it one can get, without saying “must not”.

Unfortunately, now, this catastrophic flaw has been brought forward in both character proposals as well as in the general forum, and people are explicitly asking to “use” this major catastrophic flaw. This is something, I wager, all existing software dealing with Unicode characters and encoding are totally unprepared for. Few standards and semistandards make provisions that the C0/C1 of ECMA-48 are used, and no other. And no software makes provisions for any kind of changeability of C0/C1 in any way whatsoever. So allowing C0/C1 to effectively be all private-use (in addition to the explicit private use control codes) is a major disaster for these two standards.

Fortunately, it is easy to get out this problem. Stabilize the C0/C1 areas to be the ECMA-48 (ISO/IEC 6429) ones. That is what (almost) everyone has assumed anyway. That requires a few text changes to both standards, and a proposal to some of those changes are give above. This does not mean that implementations need support more of ECMA-48 than they do already. But it also means that implementations do not need to prepare for, or be able to explicitly reject, other C0/C1 definitions. The latter capabilities are basically required as the standards are written today. And no-one should do such changes to their implementations, not even a little bit.

But what about character sets that do have other control codes (not counting the ones that just do restrictions, which in themselves are catastrophic enough, or minor movements of control codes)? Most of these “other” control codes can easily be mapped to ECMA-48 control codes or control sequences. Some need to use an extension or two, or even, when appropriate, use private use control sequences. Example mappings are given above, for Teletext, EBCDIC, PETSCII, ATASCII, ISCII styling as well as bibliographic sorting control codes.

We have given detailed motivation for why the “looseness” w.r.t. C0/C1, basically making them private-use areas is a really bad idea, standardswise. We have also given, in fair detail, mappings of “odd” control code sets to ECMA-48 (ISO/IEC 6429), showing that “odd” controls in existing character sets can be mapped to either 1) existing control codes or control sequences in ECMA-48, or 2) control sequences that can be extensions to existing functionality in ECMA-48, or (less common) use private use control sequences or even (all private use) control strings as defined in ECMA-48.

Thus, there is no need to ever define any more control codes in Unicode/10646, except possibly for very script specific ones.

Unicode and 10646 need to stabilise the C0 and C1 areas, so that they are never perceived as private-use areas, and indeed have (or at least should have for C1) a stable semantics. Something that is assumed in all modern software that deal with characters and character strings. This stabilisation should have been done from start, butt better late than never. But this is a stability that has been generally assumed (at least for C0) anyway, so formally stabilising is not a breaking change. It just reaffirms what has already been generally assumed.

On the other hand, not stabilising C0/C1, but instead openly “admit” that C0/C1 are private-use areas (by whatever name) is highly destructive and breaks absolutely all modern software that does anything with characters. and likely breaks a lot of standards and semistandards as well (for programming languages, data structure languages, and more), since few have explicit guards saying that C0 (and C1) must be the ECMA-48 ones. So, no action on this issue is a **non**-option. A highly destructive non-option.
