



Universal Multiple-Octet Coded Character Set
International Organization for Standardization
Международная организация по стандартизации

ISO/IEC JTC1/SC2/WG2
Universal Multiple-Octet Coded Character Set (UCS) – ISO/IEC 10646
Secretariat: ANSI

Doc Type: Expert Contribution
Title: Bengali encoding in ISO/IEC 10646 and BSD 1520
Source: Michael Everson, Everson Gunn Teoranta (IE)
Project: JTC1.02.18.01
Status: Discussion
Action: FYI
Due Date: —
Distribution: WG2
Medium: Paper and PDF
No. of Pages: 6

This contribution is in response to N1634, a contribution from Dr Golam Mowlah, Director General of Bangladesh Standards and Testing Institution (BSTI), dated 1997-06-29. In that document, BSTI requests the “incorporation” of the code table of BSD 1520:1995 “Bangla Coded Character Set” into ISO/IEC 10646. The rationale given for this incorporation is that the “Bengali coded character which is available in ISO/IEC publication 10646:1993 edition in row 9, locating 0980–09FF (hex) is not applicable for Bangladesh as that coded characters are not similar to [the Bangladeshi] national coded character”.

WG2 has made a thorough review of the code table in N1634 and has compared it to the Bengali code table in ISO/IEC 10646 and has found that BSD 1520 is mappable to 10646 as it currently stands, that no change to 10646 is required, and that with an adequate table-lookup it will be possible for data coded in BSD 1520 to be transformed into 10646 coding. This determination is in accordance with the agreed WG2 and UTC principle that character codes (or character names) *shall not be changed* once they have been allocated.

The analysis given below is imperfect because, unfortunately, N1634 does not give name identifications for BSD 1520 characters 83–FF, and so this analysis may have errors deriving from our imperfect knowledge of Bengali conjunct formation. However, it is clear from the analysis that even those conjuncts which we have not identified can be represented with the current 10646 coding conventions – it is simply a question of decomposing the conjunct forms. We were unable to identify 4 of the 125 characters in BSD 1520 from the range 83–FF (99, C6, F2, F8).

Almost all of the characters from that range are glyph representations of an underlying coding compatible with 10646 coding for this script. The character set appears (because of the characters CA, and D0–D6) to be based on an Apple Macintosh implementation for Bengali. Numerous precomposed conjuncts (e.g. 8D–92) are given, as well as a number of on-screen composable variants (such as the positional variants for following -RA, intended to combine with characters of differing widths and heights (DE–E0, ঞ ঞ ঞ). If the implementations in Bangladesh are fully

WorldScript compatible, the table-lookup for conversion to 10646 will be relatively straightforward. If the implementations are glyph-based, the lookup-table mapping will be more complex, but still definable and simple enough to implement.

Four characters appearing in UCS do not appear in BSD 1520: U+09C4, U+09E1, U+09E2, U+09E3. One character, BENGALI ABAGRAHA, does not appear in either ISO/IEC 10646 or BSD 1520 (cf. the LOC Romanization tables). Was this unified with DEVANAGARI AVAGRAHA or is a unique BENGALI ABAGRAHA required? The Unicode Standard should note, informatively, that ৭ BENGALI DIGIT SEVEN is used to denote the Hindu god Ganesha.

The table which follows gives the BSD 1520 character set, with its hexadecimal notation, and with strings of 10646 characters to which they can be mapped if the implementation in Bangladesh is glyph-based. Many scripts require mapping tables to convert between 10646 and local standards, and most vendors' technology already has facilities to support such mapping. However, full information on local standards is important to the correct formation of those tables.

BSTI should also be aware that SC2/WG2 and the Unicode Consortium take its contribution in N1634 seriously, and wish to ensure them that implementors of 10646 and Unicode intend to support mapping between BSD 1520 and 10646. However, because UCS is *not* a "registry" for national standards, BSD 1520:1995 cannot be adopted to *replace* the existing code table in Row 9 of the Standard. What BSTI can do, to facilitate mapping between BST 1520 and 10646, in order to produce comprehensive lookup tables for mapping between UCS and the Bangladeshi Standard, is provide us with a full description of the implementation (Worldscript-based or glyph-based) and, if glyph based, an exhaustive list of Bangla syllables coded according to the principles of BSD 1520:1995.

If careful analysis of BSD 1520 shows that one or more characters cannot be mapped directly (or with reasonable, local, context analysis) to 10646, then those characters may be candidates to be added to the standard.

Char.	1520	UCS			
SP	20	0020	^	3C	003C
!	21	0021		3D	003D
"	22	2033	V	3E	003E
#	23	0023	?	3F	003F
£	24	09F3	©	40	0040
%	25	0025	ঐ	41	0985
x	26	00D7	ঊ	42	0986
'	27	0027	ঋ	43	0987
(28	0028	৐	44	0988
)	29	0029	৑	45	0989
*	2A	002A	৒	46	098A
+	2B	002B	৓	47	098B
,	2C	0026	৔	48	098F
-	2D	002D	৕	49	0990
.	2E	002E	৖	4A	0993
/	2F	002F	ৗ	4B	0994
0	30	09E6	৙	4C	0995
১	31	09E7	৚	4D	0996
২	32	09E8	৛	4E	0997
৩	33	09E9	ড়	4F	0998
৪	34	09EA	ঢ়	50	0999
৫	35	09EB	৞	51	099A
৬	36	09EC	য়	52	099B
৭	37	09ED	ৠ	53	099C
৮	38	09EE	ৡ	54	099D
৯	39	09EF	ৢ	55	099E
:	3A	003A	ৣ	56	099F
;	3B	003B	৤	57	09A0
			৥	58	09A1

କ	93	0995 09CD 09AE	୧	B0	09CD 09B9
କ	94	09CD 0096	୧	B1	09A4 09CD 09A5
କ	95	09CD 0997	୧	B2	09A6 09CD
କ	96	0997 09C1	୧	B3	09A6 09CD
କ	97	0997 09CD 09A6	୧	B4	09A6 09CD 09A6
କ	98	0997 09CD 09A6	୧	B5	09A6 09CD 09A7
କ	99		୧	B6	09A6 09CD 09AC
କ	9A	0999 09CD 0995	୧	B7	09A8 09CD
କ	9B	0950 09CD 094E	୧	B8	09A8 09CD
କ	9C	09CD 099A	୧	B9	09CD 09A8
କ	9D	099A 09CD 099E	୧	BA	09CD 09A8
କ	9E	099C 09CD 099C	୧	BB	09A8 09CD 09A0
କ	9F	099C 09CD 099D	୧	BC	09A8 09CD 09A1
କ	A0	099C 09CD 099E	୧	BD	09A8 09CD 09A7
କ	A1	09F9	୧	BE	09A8 09CD 09B6
କ	A2	099E 09CD 099A	୧	BF	09B6 09CD
କ	A3	099E 09CD 099B	୧	C0	09CD 09B6
କ	A4	099E 09CD 099C	୧	C1	09B6 09CD 099F
କ	A5	099E 09CD 099D	୧	C2	09B6 09CD 09B6
କ	A6	0956 09C1	୧	C3	09B6 09C1
କ	A7	09A1 09CD 09A1	୧	C4	09CD 09B7
କ	A8	09A3 09CD 0956	୧	C5	09F0
କ	A9	09A3 09CD 09A0	୧	C6	
କ	AA	09A3 09C1	୧	C7	09CD 09AC
କ	AB	09CD 0961	୧	C8	09CD 09AC
କ	AC	09CD 0961 09CD 0961	୧	C9	09CD 09AC
କ	AD	0961 09CD 0961	୧	CA	00CA
କ	AE	0961 09CD 09B0	୧	CB	09AC 09CD 099C
କ	AF	09CD 0961 09CD 09B0	୧	CC	09AC 09CD 09A6

NBSP

ଅ	CD	09AC 09CD 09A7	କ	EA	09B7 09CD 099F
୧	CE	09CD 09AD	ଖ	EB	09B7 09CD 09A0
୨	CF	09AD 09CD 09B0	ଗ	EC	09B8 09CD
୩	D0	2013	ଘ	ED	___ 09CD
୪	D1	2014	ଙ	EE	09B9 09C1
୫	D2	201C	ଚ	EF	09B9 09C3
୬	D3	201D	ଟ	F0	09B9 09CD 09A3
୭	D4	2018	ଠ	F1	09B9 09CD 09AE
୮	D5	2019	ଡ	F2	
୯	D6	00D6	ଣ	F3	09C1
୧୦	D7	09CD 09AD 09CD 09B0	ତ	F4	09CD 09A4
୧୧	D8	096A 09CD	ଥ	F5	09C2
୧୨	D9	09CD 096A	ଦ	F6	09C3
୧୩	DA	09CD 096A	ଧ	F7	09C3
୧୪	DB	09CD 09AF	ନ	F8	09C8?
୧୫	DC	09B0 09CD	ପ	F9	09CD
୧୬	DD		ଫ	FA	09F4
୧୭	DE	09CD 09B0	ବ	FB	09F5
୧୮	DF	09CD 09B0	ଭ	FC	09F6
୧୯	E0	09CD 09B0	ୱ	FD	09F7
୨୦	E1	09B2 09CD	୲	FE	09F8
୨୧	E2	09B2 09CD	୳	FF	09BC
୨୨	E3	09CD 09B2			
୨୩	E4	09CD 09B2			
୨୪	E5	09B2 09CD 09A1			
୨୫	E6	09B6 09C1			
୨୬	E7	09B6 09CD			
୨୭	E8	09B7 09CD			
୨୮	E9	09B7 09CD 099E			