

Standardization of Sanskrit for Electronic Data Transfer and Screen Representation

Dominik Wujastyk

9 September 1990

Text Encoding Guidelines

During the 8th World Sanskrit Conference, Vienna 1990, a panel was held to discuss the standardization of Sanskrit for electronic data transfer. Participants were encouraged to acquire and study the *ACH-ACL-ALLC Guidelines for the Encoding and Interchange of Machine-readable Texts*, edited by Lou BURNARD and C. M. SPERBERG-MCQUEEN (Chicago and Oxford, 1990). These *Guidelines* are available free of charge in Europe from L. Burnard, Oxford University Computing Service, 13 Banbury Road, Oxford OX2 6NN, England, or in the USA from C. M. Sperberg-McQueen, Computer Center (MIC 135), University of Illinois at Chicago, Box 6998, Chicago, IL 60680, USA.

7-bit coding for file transfer

Professor H. Falk presented a program called CONVERT that conveniently converts any coding scheme used in a data file to any other coding scheme. This program was generously made available at no cost, together with Turbo Pascal source code. Prof. Falk also presented a very useful 7-bit, multi-byte “mediation code” which will be of general use for file exchange.

8-bit character set for text display

Finally, although the above two provisions cover all essential needs, the panel still felt that a standard assignment of graphic codes for the display of Sanskrit transliteration would be helpful. An ad hoc committee of interested parties was formed, and two 8-bit ‘code pages’ were designed. One, *Classical Sanskrit* (CS), for standard use and another, *Classical Sanskrit Extended* (CSX), which included the former, but also provided for Vedic, MIA, Tamil and some special usages. These code pages take as their point of departure IBM’s code page 437, the default set of character codes built into the IBM PC and clones. The characters listed below are replacements for the characters in code page 437 which have the same numerical code. E.g., character number 224 in code page 437 is a Greek letter alpha (α); CS redefines it to be a with a macron (\bar{a}). All codes not specified below are assumed to be as code page 437. E.g., character number 130 is e acute (\acute{e}).

The codes assigned were as follow:

Classical Sanskrit (CS)

166 l tilde	\tilde{l}	240 N overdot	\dot{N}
167 m overdot	\dot{m}	241 t underdot	$\underset{\cdot}{t}$
224 a macron	\bar{a}	242 T underdot	$\underset{\cdot}{T}$
225 not used (normally German <i>eszett</i> , β)		243 d underdot	$\underset{\cdot}{d}$
226 A macron	\bar{A}	244 D underdot	$\underset{\cdot}{D}$
227 i macron	\bar{i}	245 n underdot	$\underset{\cdot}{n}$
228 I macron	\bar{I}	246 N underdot	$\underset{\cdot}{N}$
229 u macron	\bar{u}	247 s acute	\acute{s}
230 U macron	\bar{U}	248 S acute	\acute{S}
231 r underdot	$\underset{\cdot}{r}$	249 s underdot	$\underset{\cdot}{s}$
232 R underdot	$\underset{\cdot}{R}$	250 S underdot	$\underset{\cdot}{S}$
233 r underdot macron	$\underset{\cdot}{\bar{r}}$	251 not used (normally the root sign $\sqrt{\quad}$)	
234 R underdot macron	$\underset{\cdot}{\bar{R}}$	252 m underdot	$\underset{\cdot}{m}$
235 l underdot	$\underset{\cdot}{l}$	253 M underdot	$\underset{\cdot}{M}$
236 L underdot	$\underset{\cdot}{L}$	254 h underdot	$\underset{\cdot}{h}$
237 l underdot macron	$\underset{\cdot}{\bar{l}}$	255 H underdot	$\underset{\cdot}{H}$
238 L underdot macron	$\underset{\cdot}{\bar{L}}$		
239 n overdot	\dot{n}		

Classical Sanskrit Extended (CSX) additions

The following definitions are added to the above Classical Sanskrit character set.

159 r underbar	ṛ	199 r underdot grave	ṛ̇
168 a macron breve	ā̆	207 r underdot macron acute	ṛ̇́
169 i macron breve	ī̆	208 a tilde	ã
170 u macron breve	ū̆	209 i tilde	ĩ
173 n underbar	ṅ	210 u tilde	ũ
181 a macron acute	á	211 e tilde	ẽ
182 a macron grave	à	212 o tilde	õ
183 i macron acute	í	213 e breve	ĕ
184 i macron grave	ì	214 o breve	ŏ
189 u macron acute	ú	215 l underbar	ḷ
190 u macron grave	ù		
198 r underdot acute	ṛ́		

These codes were chosen to have minimal impact on the standard IBM PC extended ASCII character set, but they are intended for general use in displaying Indological texts on any machine with an 8-bit (or greater) character set.

Dr. D. Wujastyk will be making available small programs that load the above character sets into the EGA or VGA display adaptors, for IBM PC users.

The above character codings have been approved by R. E. Emmerick, H. Falk, R. Lariviere, G. J. Meulenbeld, H. Nakatani, M. Tokunaga, D. Wujastyk, P. Schreiner and M. Yano.

These character codings are primarily intended for use in situations when the screen display of these characters is required, such as in word processing. They may, of course, be used for data transfer, where, however, a 7-bit code (perhaps with multi-byte character codes) is still preferable. One such 7-bit scheme is provided by H. Falk (see 2. above).

These character codings are currently open for discussion and comments may be directed to Dr. D. Wujastyk at

Wellcome Institute,
183 Euston Road,
London NW1 2BN, England,

or by email at

Bitnet/Earn: `dow@harvunxw` or

Janet: `D.Wujastyk@uk.ac.ucl`.

After a suitable lapse of time, the character sets will be sent to ECMA and ISO for registration. They will also be sent to the Text Encoding Initiative for registration, probably with H. Falk's 7-bit coding scheme.

Such registration in no way enforces these schemes; it merely makes them available centrally for reference. Other schemes may also be registered in the future.