

**ISO/IEC JTC 1/SC 2/WG 2
PROPOSAL SUMMARY FORM TO ACCOMPANY SUBMISSIONS
FOR ADDITIONS TO THE REPERTOIRE OF ISO/IEC 10646¹**

Please fill all the sections A, B and C below.

(Please read Principles and Procedures Document for guidelines and details before filling this form.)

See <http://www.dkuug.dk/JTC1/SC2/WG2/docs/summaryform.html> for latest *Form*.

See <http://www.dkuug.dk/JTC1/SC2/WG2/docs/principles.html> for latest *Principles and Procedures* document.

See <http://www.dkuug.dk/JTC1/SC2/WG2/docs/roadmaps.html> for latest roadmaps.

A. Administrative

1. Title:	Proposal to Encode Kharoṣṭhī in Plane 1 of ISO/IEC 10646
2. Requester's name:	Andrew Glass, Stefan Baums, Richard Salomon, UTC
3. Requester type (Member body/Liaison/Individual contribution):	Individual contribution
4. Submission date:	18 September 2003
5. Requester's reference (if applicable):	_____
6. (Choose one of the following:)	
This is a complete proposal:	Yes
or, More information will be provided later:	_____

B. Technical - General

1. (Choose one of the following:)	
a. This proposal is for a new script (set of characters):	Yes
Proposed name of script:	Kharoṣṭhī / KHAROSTHI
b. The proposal is for addition of character(s) to an existing block:	_____
Name of the existing block:	_____
2. Number of characters in proposal:	65
3. Proposed category (see section II, Character Categories):	C
4. Proposed Level of Implementation (1, 2 or 3) (see clause 14, ISO/IEC 10646-1: 2000):	Level 3
Is a rationale provided for the choice?	Yes
If Yes, reference: Combining marks used.	
5. Is a repertoire including character names provided?	Yes
a. If YES, are the names in accordance with the 'character naming guidelines in Annex L of ISO/IEC 10646-1: 2000?	Yes
b. Are the character shapes attached in a legible form suitable for review?	Yes
6. Who will provide the appropriate computerized font (ordered preference: True Type, or PostScript format) for publishing the standard? Andrew Glass (True Type)	
If available now, identify source(s) for the font (include address, e-mail, ftp-site, etc.) and indicate the tools used:	Not yet available.
7. References:	
a. Are references (to other character sets, dictionaries, descriptive texts etc.) provided?	Yes
b. Are published examples of use (such as samples from newspapers, magazines, or other sources) of proposed characters attached?	Yes
8. Special encoding issues:	
Does the proposal address other aspects of character data processing (if applicable) such as input, presentation, sorting, searching, indexing, transliteration etc. (if yes please enclose information)?	
Yes. It covers Kharoṣṭhī bidirectional behavior and gives normative rules required for rendering the script.	
9. Additional Information:	
Submitters are invited to provide any additional information about Properties of the proposed Character(s) or Script that will assist in correct understanding of and correct linguistic processing of the proposed character(s) or script. Examples of such properties are: Casing information, Numeric information, Currency information, Display behaviour information such as line breaks, widths etc., Combining behaviour, Spacing behaviour, Directional behaviour, Default Collation behaviour, relevance in Mark Up contexts, Compatibility equivalence and other Unicode normalization related information. See the Unicode standard at http://www.unicode.org for such information on other scripts. Also see http://www.unicode.org/Public/UNIDATA/UnicodeCharacterDatabase.html and associated Unicode Technical Reports for information needed for consideration by the Unicode Technical Committee for inclusion in the Unicode Standard.	

¹ Form number: N2352-F (Original 1994-10-14; Revised 1995-01, 1995-04, 1996-04, 1996-08, 1999-03, 2001-05, 2001-09)

C. Technical - Justification

1. Has this proposal for addition of character(s) been submitted before? If YES explain _____	No
2. Has contact been made to members of the user community (for example: National Body, user groups of the script or characters, other experts, etc.)? If YES, with whom? Richard Salomon, Andrew Glass If YES, available relevant documents: Kharosthī Manuscript Paleography	Yes
3. Information on the user community for the proposed characters (for example: size, demographics, information technology use, or publishing use) is included? Reference: _____	Scholars
4. The context of use for the proposed characters (type of use; common or rare) Reference: _____	Scholarly; Rare
5. Are the proposed characters in current use by the user community? If YES, where? Reference: Scholars worldwide	Yes
6. After giving due considerations to the principles in <i>Principles and Procedures document</i> (a WG 2 standing document) must the proposed characters be entirely in the BMP? If YES, is a rationale provided? _____ If YES, reference: _____	No
7. Should the proposed characters be kept together in a contiguous range (rather than being scattered)?	Yes
8. Can any of the proposed characters be considered a presentation form of an existing character or character sequence? If YES, is a rationale for its inclusion provided? _____ If YES, reference: _____	No
9. Can any of the proposed characters be encoded using a composed character sequence of either existing characters or other proposed characters? If YES, is a rationale for its inclusion provided? _____ If YES, reference: _____	No
10. Can any of the proposed character(s) be considered to be similar (in appearance or function) to an existing character? If YES, is a rationale for its inclusion provided? _____ If YES, reference: _____	Yes Yes See below
11. Does the proposal include use of combining characters and/or use of composite sequences (see clauses 4.12 and 4.14 in ISO/IEC 10646-1: 2000)? If YES, is a rationale for such use provided? _____ If YES, reference: See below; and Kharosthī Manuscript Paleography Is a list of composite sequences and their corresponding glyph images (graphic symbols) provided? _____ If YES, reference: See below; and Kharosthī Manuscript Paleography	Yes Yes Yes
12. Does the proposal contain characters with any special properties such as control function or similar semantics? If YES, describe in detail (include attachment if necessary)	Yes Virāma (10A3F)
13. Does the proposal contain any Ideographic compatibility character(s)? If YES, is the equivalent corresponding unified ideographic character(s) identified? _____ If YES, reference: _____	No

Proposal for Kharoṣṭhī script

This is a proposed assignment for Kharoṣṭhī characters. The Kharoṣṭhī script was used to write Gāndhārī and Sanskrit as well as various mixed dialects termed ‘Gāndhārī Hybrid Sanskrit’ (see Salomon 2001). The characters in this proposal are derived from sources in the Kharoṣṭhī script from across the whole range of known manuscripts and inscriptions. The intention is to provide a standard method for writing Kharoṣṭhī, and also a common means for the electronic storage of manuscript data. The Unicode Consortium has not previously published a proposal for Kharoṣṭhī.

Brief History of the Kharoṣṭhī script

The Kharoṣṭhī script is one of the two ancient writing systems of India in the historical period. Unlike the pan-Indian Brāhmī script, Kharoṣṭhī was confined to the northwest of India centered on the region of *Gandhāra* (modern northern Pakistan and eastern Afghanistan; see map). The exact details of its origin remain obscure despite the attention of several generations of scholars, but it is almost certainly related to Aramaic, stemming from the time of the Achaemenid conquest and occupation of that region from 559–336 BCE (Salomon 1998: 51–4). The Kharoṣṭhī script first appears in a fully developed form in the Aśokan inscriptions at Shāhbāzgarhī and Mānsehrā which have been dated to around 250 BCE (Hultsch 1925: xxxv). The script continued to be used in Gandhāra and neighboring regions, sometimes alongside Brāhmī, until around the third century CE, when it disappeared from its homeland (Salomon 1996: 375). The Kharoṣṭhī script was also used for official documents and epigraphs in the Central Asian cities of Khotan and Niya in the third and fourth centuries CE, and appears to have survived in Kucha and neighboring areas along of the Northern Silk Road until the seventh century. This form of the script has been termed Formal Kharoṣṭhī in recent publications (Sander 1999: 72, Lin 2003: 1).



Map: Geographical extent of the Kharoṣṭhī script

The Kharoṣṭhī script was initially deciphered around the middle of the nineteenth century by James Prinsep and others who worked from the short bicult inscriptions (Greek and Kharoṣṭhī) on the coins of the Indo-Greek and Indo-Scythian kings. The decipherment has been refined over

the last 150 years as more material has come to light. We now have several examples of Sanskrit, or Sanskritized Gāndhārī, written in Kharoṣṭhī script. The current proposal makes provision for encoding the level of Sanskrit found in the known documents (see Salomon 2001).

Formal Kharoṣṭhī

A distinct form of the Kharoṣṭhī script used to write both Gāndhārī and Tocharian B is found in a small number of documents from sites along the Northern Silk Route. Some work on the decipherment of this form of the Kharoṣṭhī script has been published recently (see Lin 2003). It would be possible to render the texts in Formal Kharoṣṭhī in this publication with the code points proposed here, however, it is likely that when this form of the writing system is fully understood, it may be desirable to introduce additional characters. Furthermore, as new items in the Kharoṣṭhī script continue to be discovered in Afghanistan, Pakistan and Central Asia, other signs may be needed at some point in the future. Therefore, we have reserved some code points to accommodate any additional needs for both standard and Formal Kharoṣṭhī.

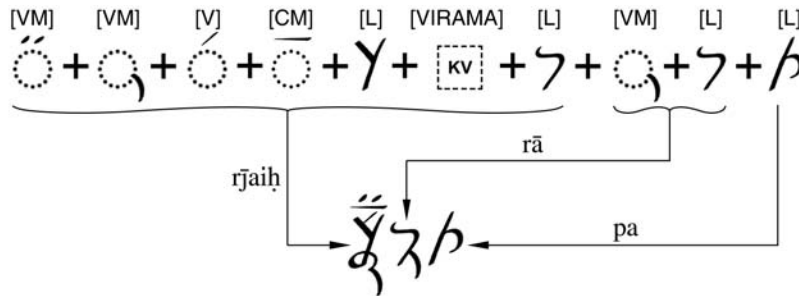
The Writing System

The Kharoṣṭhī script is a member of the Indic script family and conforms to the alphasyllabic or abugida script type. However, unlike the other scripts of this group, it is written from right to left. Kharoṣṭhī letters do not have positional variants as in Arabic and Hebrew.

Unicode Bidirectional Algorithm. Kharoṣṭhī can be implemented using the rules of the Unicode Bidirectional Algorithm as they apply to Arabic and Hebrew, with the exception that in Kharoṣṭhī both letters and digits are written from right to left.

Convention. In what follows, we have followed the unicode naming conventions for the Indic scripts (see <http://www.unicode.org/charts/PDF/U0900.pdf>), with slight adaptations based on current scholarly conventions for naming Kharoṣṭhī letters (see Glass 2000: 33–113).

Diacritic Marks/Vowels. All vowels other than *a* are written with diacritic marks in Kharoṣṭhī. In addition, there are six vowel modifiers and three consonant modifiers which are written with combining diacritics. Some letters may take more than one such diacritical mark. In these cases the correct sequence should be: Letter (L) + [Consonant Modifier (CM)] + [Vowel (V)] + [Vowel Modifier (VM)]. For example the Sanskrit word *parārdhyaiḥ* might be rendered in Kharoṣṭhī script as **parārjaiḥ* (written from right to left):



Numerical Signs. Kharoṣṭhī employs a set of numeral signs unique to the script. These have been included in this proposal. The numerals, like the letters, are written from right to left. Numbers in Kharoṣṭhī are based on an additive system. There is no zero, nor separate signs for the numbers

5–9. The number 1996, for example, would appear as: 1000 4 4 1 100 20 20 20 20 10 4 2 (see Glass 2000: 139–43).

୫ ୪ ୪ ୩ ୩ ୩ ୩ ୩ ୧ ୧ ୪ ୪ ୧

Punctuation. Nine different punctuation marks are used in Kharoṣṭhī manuscripts and inscriptions. These have been included in this proposal (see Glass 2000: 144–7).

Minimum Rendering Requirements. Rendering requirements for Kharoṣṭhī are similar to those used for Devanāgarī. The remainder of this section specifies a minimum set of rules that provide legible Kharoṣṭhī diacritic and ligature substitution behavior.

Combining Classes. The various combining diacritics attach to the full characters in different ways. A number of classes have been determined on the basis of their standard positions.

VOWEL SIGNS:

Combining *-i*:

Horizontal: example $a + -i \rightarrow i$

𑀓 + 𑀓̇ → 𑀓̇

members of this class: a, na, ha .

Diagonal: example $ka + -i \rightarrow ki$

𑀓𑀕 + 𑀓̇ → 𑀓̇𑀕

members of this class: $ka, k̄a, kha, ga, gha, ca, cha, ja, ña, ṭa, ṭha, ḥa, ḍa, ḍha, ṇa, ta, da, dha, ba, bha, ya, ra, va, ṣa, sa, za$.

Vertical: example $tha + -i \rightarrow thi$

𑀓𑀕 + 𑀓̇ → 𑀓̇𑀕

members of this class: $tha, pa, pha, ma, la, śa$.

Combining *-u*:

Attached: example $a + -u \rightarrow u$

𑀓 + 𑀓̈ → 𑀓̈

members of this class: $a, ka, k̄a, kha, ga, gha, ca, cha, ja, ña, ṭa, ṭha, ḍa, ḍha, ṇa, ta, tha, da, dha, na, pa, pha, ba, bha, ya, ra, la, va, śa, ṣa, sa, za$.

Independent: example $ha + -u \rightarrow hu$

𑀓 + 𑀓̈ → 𑀓̈

members of this class: $ṭa, ha$.

Ligatured: example $ma + -u \rightarrow mu$

𑀓 + 𑀓̈ → 𑀓̈

members of this class: ma .

Combining *-r*:

Attached: example $a + -r \rightarrow r$

𑀓 + 𑀓̇ → 𑀓̇

members of this class: *a, ka, ká, kha, ga, gha, ca, cha, ja, ta, da, dha, na, pa, pha, ba, bha, va, śa, sa.*

Independent: example *ma + -ṛ → mṛ*

U + ṛ → Ṛ

members of this class: *ma, ha.*

Combining -e:

Horizontal: example *a + -e → e*

Ṛ + ē → Ṝ

members of this class: *a, na, ha.*

Diagonal: example *ka + -e → ke*

Ṛ + ē → Ṝ

members of this class: *ka, ká, kha, ga, gha, ca, cha, ja, ña, ṭa, ṭha, ṭha, ḍa, ḍha, ṇa, ta, dha, ba, bha, ya, ra, va, śa, sa, za.*

Vertical: example *tha + -e → the*

Ṛ + ē → Ṝ

members of this class: *tha, pa, pha, la, śa.*

Ligatured: example *da + -e → de*

Ṛ + ē → Ṝ

members of this class: *da, ma.*

Combining -o:

Diagonal: example *a + -o → o*

Ṛ + ō → Ṝ

members of this class: *a, ka, ká, kha, ga, gha, ca, cha, ja, ña, ṭa, ṭha, ṭha, ḍa, ḍha, ṇa, ta, tha, da, dha, na, ba, bha, ma, ra, la, va, śa, sa, za, ha.*

Vertical: example *pa + .o → po*

Ṛ + ō → Ṝ

members of this class: *pa, pha, ya, śa.*

VOWEL MODIFIERS:

Combining VOWEL LENGTH MARK:

This sign may be used with *-a, -i, -u, -ṛ*, to indicate the equivalent long vowel *-ā, -ī, -ū, -ṝ*. In combination with *-e* and *-o* it indicates the diphthongs *-ai* and *-au*.

Example *ma + ¯ → mā*

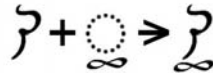
U + ¯ → Ū

combines with: *-a, -i, -ṛ, -u, -e, -o.*

Combining DOUBLE RING BELOW:

This sign appears in some of the Central Asian documents. Its precise phonetic value has not yet been established.

Example *sa* + ◌ → *sq*



combines with: *-a, -u.*

Combining ANUSVARA:

This sign indicates nasalization of the vowel or a nasal segment following the vowel.

Example *a* + *-ṃ* → *aṃ*



combines with: *-a, -i, -u, -ṛ, -e, -o.*

Combining VISARGA:

This sign is generally used to indicate unvoiced syllable-final [h]. A secondary usage is as a vowel length marker.

Example *ka* + *-ḥ* → *kaḥ*



combines with: *-a, -i, -u, -ṛ, -e, -o.*

CONSONANT MODIFIERS:

Combining BAR ABOVE:

This sign is used to indicate various modified pronunciations depending on the consonants involved, such as nasalization or aspiration.

Example *ja* + *-̄* → *jā*



combines with: *kṣa ga, ca, ja, na, ma, śa, ṣa, sa, ha.*

Combining CAUDA:

This sign is used to indicate various modified pronunciations of the consonants involved, particularly fricativization.

Example *ga* + *·* → *gá*



combines with: *ga, ja, ḍa, ta, da, pa, ya, va, śa, sa.*

Combining DOT BELOW:

The precise value of this sign has not yet been determined.

Example *ma* + *·* → *ṃa*



combines with: *ma, ha.*

COMBINING VIRAMA:

This is a control character. When not followed by a consonant it causes the preceding consonant to be written as subscript to the left of the letter before it. If followed by another consonant, it will trigger a combined form consisting of two or more consonants. The resulting form may also be subject to combinations with the above Combining Diacritics.

Examples:

Pure VIRAMA:

$dha + i + k + [\text{VIRAMA}] \rightarrow dhik$

३ + ि + क + [KV] ⇒ धि३क

Ligatures:

$ka + [\text{VIRAMA}] + \text{ṣa} \rightarrow kṣa$

क + [KV] + ष ⇒ क्ष

$ma + [\text{VIRAMA}] + ra \rightarrow mra$

म + [KV] + र ⇒ म्र

$va + [\text{VIRAMA}] + ha \rightarrow vha$

व + [KV] + ह ⇒ व्ह

$sa + [\text{VIRAMA}] + ta \rightarrow sta$

स + [KV] + त ⇒ स्त

members of this class: $kṣV$, tsV , mrV , vhV , stV .

Consonants with special combining forms:

$sa + [\text{VIRAMA}] + \text{ṣa} \rightarrow sṣa$

स + [KV] + ष ⇒ स्ष

$ṛa + [\text{VIRAMA}] + ta \rightarrow rta$

ॠ + [KV] + त ⇒ र्त

$ta + [\text{VIRAMA}] + \text{ṛa} \rightarrow tra$

त + [KV] + ॠ ⇒ त्र

$ḷa + [\text{VIRAMA}] + pa \rightarrow ḷpa$

ॡ + [KV] + प ⇒ ष्ट

$pa + [\text{VIRAMA}] + \text{ḷa} \rightarrow pla$

प + [KV] + ॡ ⇒ ष्ट

$ka + [\text{VIRAMA}] + \underline{\underline{la}} \rightarrow kla$

$\text{𑌕} + \boxed{\text{KV}} + \text{𑌗} \Rightarrow \text{𑌕}$

$ta + [\text{VIRAMA}] + \underline{\underline{va}} \rightarrow tva$

$\text{𑌖} + \boxed{\text{KV}} + \text{𑌘} \Rightarrow \text{𑌖}$

members of this class: CyV, rCV, CrV, lCV, ClV, CvV.

Consonants with full combined forms:

$\underline{\underline{ka}} + [\text{VIRAMA}] + \underline{\underline{ta}} \rightarrow kta$

$\text{𑌕} + \boxed{\text{KV}} + \text{𑌖} \Rightarrow \text{𑌕}$

$\underline{\underline{kha}} + [\text{VIRAMA}] + ka + [\text{VIRAMA}] + \underline{\underline{sa}} \rightarrow khk\u0308sa$

$\text{𑌕} + \boxed{\text{KV}} + \text{𑌕} + \boxed{\text{KV}} + \text{𑌖} \Rightarrow \text{𑌕}$

members of this class: $k, kh, g, \acute{g}, c, j, \tilde{n}, \acute{t}, \acute{th}, \acute{d}, \acute{dh}, \acute{n}, t, th, d, dh, n, p, b, bh, m, y$ (in ryV), l (in lmV), v (in vrV), $\acute{s}, \acute{\acute{s}}, s, z, h$.

Kharoṣṭhī
















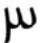




















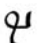























Range: 10A00 to 10A5F

These charts contain only proposed assignments and should not be considered valid until such time as the Unicode Consortium formally accepts them.

Andrew Glass created the fonts used in these charts.














Code chart

The code chart characters are normalized forms based on manuscripts of the *first* century CE

















	10A0	10A1	10A2	10A3	10A4	10A5
0	 10A00	 10A10	 10A20	 10A30	 10A40	 10A50
1	 10A01	 10A11	 10A21	 10A31	 10A41	 10A51
2	 10A02	 10A12	 10A22	 10A32	 10A42	 10A52
3	 10A03	 10A13	 10A23	 10A33	 10A43	 10A53
4			 10A24		 10A44	 10A54
5	 10A05	 10A15	 10A25		 10A45	 10A55
6	 10A06	 10A16	 10A26		 10A46	 10A56
7		 10A17	 10A27		 10A47	 10A57
8			 10A28	 10A38		 10A58
9		 10A19	 10A29	 10A39		
A		 10A1A	 10A2A	 10A3A		
B		 10A1B	 10A2B			
C	 10A0C	 10A1C	 10A2C			
D	 10A0D	 10A1D	 10A2D			
E	 10A0E	 10A1E	 10A2E			
F	 10A0F	 10A1F	 10A2F	 10A3F		

Name chart







The name chart characters are normalized forms based on manuscripts of the *first* century CE. Additional information about individual characters in this block can be found in [Appendix 1](#).

Glyph	Unicode code point	Name	Transcription
	10A00	KHAROSTHI LETTER A	a
	10A01	KHAROSTHI VOWEL SIGN I	i
	10A02	KHAROSTHI VOWEL SIGN U	u
	10A03	KHAROSTHI VOWEL SIGN VOCALIC R	r̥
	10A04	(This position shall not be used)	
	10A05	KHAROSTHI VOWEL SIGN E	e
	10A06	KHAROSTHI VOWEL SIGN O	o
	10A07	(This position shall not be used)	
	10A08	(This position shall not be used)	
	10A09	(This position shall not be used)	
	10A0A	(This position shall not be used)	
	10A0B	(This position shall not be used)	
	10A0C	KHAROSTHI VOWEL LENGTH MARK	-
	10A0D	KHAROSTHI SIGN DOUBLE RING BELOW	◌̣̣
	10A0E	KHAROSTHI SIGN ANUSVARA	◌̣̣̣
	10A0F	KHAROSTHI SIGN VISARGA	◌̣̣̣̣
	10A10	KHAROSTHI LETTER KA	ka
	10A11	KHAROSTHI LETTER KHA	kha
	10A12	KHAROSTHI LETTER GA	ga

Glyph	Unicode code point	Name	Transcription
𑍆	10A13	KHAROSTHI LETTER GHA	gha
	10A14	(This position shall not be used)	
𑍇	10A15	KHAROSTHI LETTER CA	ca
𑍈	10A16	KHAROSTHI LETTER CHA	cha
𑍉	10A17	KHAROSTHI LETTER JA	ja
	10A18	(This position shall not be used)	
𑍊	10A19	KHAROSTHI LETTER NYA	ña
𑍋	10A1A	KHAROSTHI LETTER TTA	ṭa
𑍌	10A1B	KHAROSTHI LETTER TTHA	ṭha
𑍍	10A1C	KHAROSTHI LETTER DDA	ḍa
𑍎	10A1D	KHAROSTHI LETTER DDHA	ḍha
𑍏	10A1E	KHAROSTHI LETTER NNA	ṇa
𑍐	10A1F	KHAROSTHI LETTER TA	ta
𑍑	10A20	KHAROSTHI LETTER THA	tha
𑍒	10A21	KHAROSTHI LETTER DA	da
𑍓	10A22	KHAROSTHI LETTER DHA	dha
𑍔	10A23	KHAROSTHI LETTER NA	na
𑍕	10A24	KHAROSTHI LETTER PA	pa
𑍖	10A25	KHAROSTHI LETTER PHA	pha
𑍗	10A26	KHAROSTHI LETTER BA	ba

Glyph	Unicode code point	Name	Transcription
	10A27	KHAROSTHI LETTER BHA	bha
	10A28	KHAROSTHI LETTER MA	ma
	10A29	KHAROSTHI LETTER YA	ya
	10A2A	KHAROSTHI LETTER RA	ra
	10A2B	KHAROSTHI LETTER LA	la
	10A2C	KHAROSTHI LETTER VA	va
	10A2D	KHAROSTHI LETTER SHA	śa
	10A2E	KHAROSTHI LETTER SSA	ṣa
	10A2F	KHAROSTHI LETTER SA	sa
	10A30	KHAROSTHI LETTER ZA	za
	10A31	KHAROSTHI LETTER HA	ha
	10A32	KHAROSTHI LETTER KKA	ka
	10A33	KHAROSTHI LETTER TTTHA	ṭha
	10A34	(This position shall not be used)	
	10A35	(This position shall not be used)	
	10A36	(This position shall not be used)	
	10A37	(This position shall not be used)	
	10A38	KHAROSTHI SIGN BAR ABOVE	-
	10A39	KHAROSTHI SIGN CAUDA	´ or _ see Appendix 1
	10A3A	KHAROSTHI SIGN DOT BELOW	.
	10A3B	(This position shall not be used)	

Glyph	Unicode code point	Name	Transcription
	10A3C	(This position shall not be used)	
	10A3D	(This position shall not be used)	
	10A3E	(This position shall not be used)	
KV	10A3F	KHAROSTHI VIRAMA = halant · suppresses inherent vowel	see VIRAMA
१	10A40	KHAROSTHI DIGIT ONE	1
२	10A41	KHAROSTHI DIGIT TWO	2
३	10A42	KHAROSTHI DIGIT THREE	3
४	10A43	KHAROSTHI DIGIT FOUR	4
१०	10A44	KHAROSTHI NUMBER TEN	10
२०	10A45	KHAROSTHI NUMBER TWENTY	20
१००	10A46	KHAROSTHI NUMBER HUNDRED	100
१०००	10A47	KHAROSTHI NUMBER THOUSAND	1000
	10A48	(This position shall not be used)	
	10A49	(This position shall not be used)	
	10A4A	(This position shall not be used)	
	10A4B	(This position shall not be used)	
	10A4C	(This position shall not be used)	
	10A4D	(This position shall not be used)	
	10A4E	(This position shall not be used)	
	10A4F	(This position shall not be used)	
.	10A50	KHAROSTHI PUNCTUATION DOT	.
◦	10A51	KHAROSTHI PUNCTUATION SMALL CIRCLE	◦
○	10A52	KHAROSTHI PUNCTUATION CIRCLE	○

Glyph	Unicode code point	Name	Transcription
	10A53	KHAROSTHI PUNCTUATION CRESCENT BAR	€
	10A54	KHAROSTHI PUNCTUATION MANGALAM	⊕
	10A55	KHAROSTHI PUNCTUATION LOTUS	☼
	10A56	KHAROSTHI PUNCTUATION DANDA	
	10A57	KHAROSTHI PUNCTUATION DOUBLE DANDA	
	10A58	KHAROSTHI PUNCTUATION LINES	≈
	10A59	(This position shall not be used)	
	10A5A	(This position shall not be used)	
	10A5B	(This position shall not be used)	
	10A5C	(This position shall not be used)	
	10A5D	(This position shall not be used)	
	10A5E	(This position shall not be used)	
	10A5F	(This position shall not be used)	

Text Samples



Figure 1: Aśokan inscription at Shāhbāzgarhī, ca. 250 BCE (Hultsch 1925).



Figure 2: Relic vase inscription of Theodoros, ca. 50 BCE (Konow 1929: Plate 1).



Figure 3: Coin of King Azes with legend in Greek and Kharoṣṭhī, ca. 50 BCE. (The Royal Collection of Coins and Medals, National Museum, Denmark. Photographs by Stefan Baums and Helle Horsnæs. Inventory Number B.P. 917.)



Figure 4: Detail from British Library Kharoṣṭhī Fragment 5B, ca. 50 CE (Salomon 2000: Plate 2).

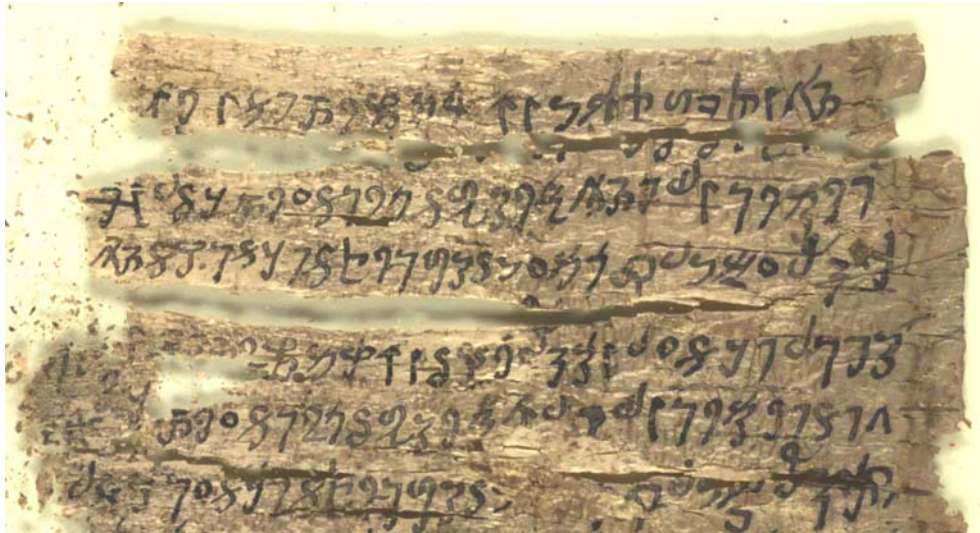


Figure 5: Detail from British Library Kharoṣṭhī Fragment 14, ca. 50 CE (Allon 2001: Plate 7).

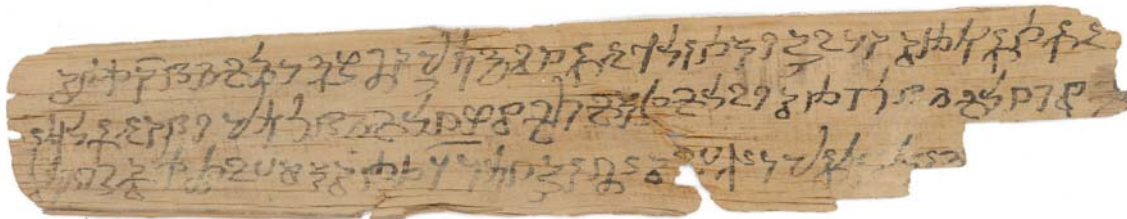


Figure 6: Fragment 44 from the Schøyen Collection, ca. 150 CE (Braarvig 2000: Plate 10.2).

budhabayaṇo (𑀧𑀲𑀭𑀯𑀢) in line 29, the scribe originally omitted the *ya* and wrote *budhabaṇa* (𑀧𑀲𑀢𑀢) before realizing his mistake and making the right leg of the *ya* from the tip of the *ṇa*, that is, *ya* 𑀲 (29.6). The *khu* 𑀧 in 38.16 was first written as a plain *kha*, and the *u*-vowel loop was added subsequently. Other examples where the scribe modified an incorrect character include *śa* 𑀱 (10.12), where the scribe corrected what he probably originally wrote as *va*. Similarly, the *ji* in 51.18 seems to have been corrected from an original *ci*. In *tva* 𑀲 (44.11), the scribe wrote a normal *ta* 𑀲 before separately adding a postconsonantal *v*. Likewise in *pra* 𑀢

Figure 7: Sample text including Kharoṣṭhī characters from a recent publication (Allon 2001: 66).















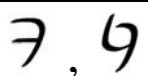




- 10A3F. This is the Kharoṣṭhī *virāma*. It is used to indicate the suppression of the inherent vowel. It not a mark or sign in itself, but a control character that causes the consonant which it follows to appear as a subscript to the preceding akṣara. When followed immediately by another consonant it triggers a conjunct form representing both consonants. See [Combining with VIRAMA](#) above. It can only follow a consonant, or a consonant modifier. It cannot follow a space, a vowel, a vowel modifier, a number, a punctuation sign, or another VIRAMA.
- 10A40 – 10A47. These are the Kharoṣṭhī numerals. They are written from right to left like the letters. The Kharoṣṭhī number system is additive/multiplicative, there is no zero, and no decimal point.
- 10A50 – 10A57. These are the Kharoṣṭhī punctuation signs. Nine punctuation signs have been identified from across the range of Kharoṣṭhī sources. Some of these punctuation signs could be considered similar (in appearance or function) to existing characters. However, we feel that independent code points should be assigned to the Kharoṣṭhī punctuation signs so that Kharoṣṭhī documents posted on the internet may be readable and searchable for those who do not have specialized Kharoṣṭhī fonts installed. For example, such documents should be readable and searchable using a future version of Arial Unicode or any other single, fallback Unicode font.

Appendix 2: Sort Order

There is an ancient abecedary connected with the Kharoṣṭhī script called Arapacana, named after its first five akṣaras. There is, however, no evidence that words were ever sorted in this order. A further complication is that there is no record in Kharoshti of the complete Arapacana sequence, while Sanskrit records are not in total agreement about the inventory and order of the letters. Therefore, we do not propose using the Arapacana as the basis for sorting.

In modern scholarly practice, Gāndhārī is sorted in much the same order as Sanskrit. Vowel length, however, even when marked, is ignored when sorting Kharoṣṭhī. In the following table, when two signs are given in a single row, they should be treated as equivalent in the sorting algorithm, the first sign having priority in tie-resolving situations, for example, *ka*, *ka*, *ki*.

Character	Unicode code point	Transcription
॑	10A00	a
॒	10A01	i
॒	10A02	u
॒	10A03	ṛ
॒	10A05	e





Character	Unicode code point	Transcription
	10A06	o
	10A0E	m̐ (preceding Ø, y-h)
	10A0F	ḥ
	10A3F	see VIRAMA
	10A10, 10A32	k, k̐
	10A11	kh
	10A12	g
	10A13	gh
	10A0E	m̐ (preceding k-gh) see note below
	10A15	c
	10A16	ch
	10A17	j
	10A19, 10A0E	ñ, m̐ (preceding c-ñ) see note below
	10A1A	ṭ
	10A1B, 10A33	ṭh, ṭh̐
	10A1C	ḍ
	10A1D	ḍh
	10A1E, 10A0E	ṇ, m̐ (preceding ṭ-ṇ) see note below
	10A1F	t

Character	Unicode code point	Transcription
†	10A20	th
ſ	10A21	d
ʒ	10A22	dh
ſ, ʒ	10A23, 10A0E	n, ɱ (preceding t-n) see note below
þ	10A24	p
ƥ	10A25	ph
ʒ	10A26	b
ʒ	10A27	bh
ʘ, ʒ	10A28, 10A0E	m, ɱ (preceding p-m) see note below
ʌ	10A29	y
ʒ	10A2A	r
†	10A2B	l
ʒ	10A2C	v
ɱ	10A2D	ś
ɱ	10A2E	ş
ʒ	10A2F	s
ʒ	10A30	z
2	10A31	h
l	10A40	l

Character	Unicode code point	Transcription
μ	10A41	2
μ	10A42	3
χ	10A43	4
∩	10A44	10
3	10A45	20
ι	10A46	100
ϣ	10A47	1000
.	10A50	.
◦	10A51	◦
○	10A52	○
∈	10A53	∈
⊕	10A54	⊕
⊙	10A55	⊙
	10A56	
	10A57	
≈	10A58	≈

The following characters, omitted in the above table, should be transparent to the sorting algorithm:

Character	Unicode code point	Transcription
⊙	10A0C	-

Character	Unicode code point	Transcription
	10A0D	o
	10A38	-
	10A39	´ or _
	10A3A	.

The sort value of ANUSVARA (10A0E) is context dependent:

- When followed by a space, the letters *y–h* (10A29 – 10A31), a number (10A40 – 10A47), a punctuation mark (10A50 – 10A57), or any non-Kharoṣṭhī character, it is considered to be a ‘true’ *anusvāra* and follows *o* (10A07) in the sort order.
- When followed by the letters *k–gh*, or *k* (10A10 – 10A13, or 10A32), it is considered to be a velar nasal and follows *gh* (10A13) in the sort order.
- When followed by the letters *c–ñ*, (10A15 – 10A19), it is functionally equivalent to *ñ* (10A19), and follows *j* (10A17) in the sort order.
- When followed by the letters *t–n*, or *ḥ* (10A1A – 10A1E, or 10A33), it is functionally equivalent to *ṇ* (10A1E), and follows *dh* (10A1D) in the sort order.
- When followed by the letters *t–n*, (10A1F – 10A23), it is functionally equivalent to *n* (10A23), and follows *dh* (10A22) in the sort order.
- When followed by a vowel or the letters *p–m*, (10A00 or 10A24 – 10A28), it is functionally equivalent to *m* (10A28), and follows *bh* (10A27) in the sort order.

The sort values of the Kharoṣṭhī digits will not produce a correct sorting of Kharoṣṭhī numerals, because of the multiplicative element in the Kharoṣṭhī numeral system. If possible, the Kharoṣṭhī numerals should be sorted according to their numeric values.

Appendix 3: Word Breaks, Line Breaks and Hyphenation

Most Kharoṣṭhī manuscripts are written as continuous text with no indication of word boundaries. Only a few examples are known where spaces have been used to separate words or verse quarters. Most scribes have tried to finish a word before starting a new line. There are no examples of anything akin to hyphenation in Kharoṣṭhī manuscripts. In cases where a word would not completely fit into a line, its continuation simply appears at the beginning of the next line. Modern scholarly practice will in most cases make use of spaces and hyphenation. When necessary, hyphenation should be applied on the model of Sanskrit.

References

- Allon, Mark. 2001. *Three Ekottarikāgama-Type Sūtras: British Library Kharoṣṭhī Fragments 12 and 14*. Gandhāran Buddhist Texts 2. Seattle: University of Washington Press.
- Boyer, A. M., E. J. Rapson, and E. Senart. 1920–9. *Kharoṣṭhī Inscriptions Discovered by Sir Aurel Stein in Chinese Turkestan*. 3 pts. (pt. 3 by Rapson and P. S. Noble). Oxford: Clarendon Press.

- Braarvig, Jens, ed. 2000. *Manuscripts in the Schøyen Collection I: Buddhist Manuscripts*, vol. 1. Oslo: Hermes Publishing.
- Glass, Andrew. 2000. “A Preliminary Study of Kharoṣṭhī Manuscript Paleography.” Master’s thesis, Department of Asian Languages and Literature, University of Washington. [http://depts.washington.edu/ebmp/downloads/Glass_2000.pdf]
- Hultzsch, E. 1925. *The Inscriptions of Aśoka*. Second edition. Corpus Inscriptionum Indicarum 1. Oxford: Clarendon Press.
- Konow, Sten, ed. 1929. *Kharoṣṭhī Inscriptions with the Exception of Those of Aśoka*. Corpus Inscriptionum Indicarum 2.1. Calcutta: Government of India. Plate 1.
- Lin, Meicun. 2003. “Five Gāndhārī Documents from Kizil in the Le Coq Collection.” *Kodai Bunka* 55.3: 1–22.
- Salomon, Richard. 1996. “Brahmi and Kharoṣṭhī” in Daniels and Bright, eds. *The World’s Writing Systems*. New York: Oxford University Press.
- . 1998. *Indian Epigraphy: A Guide to the Study of Inscriptions in Sanskrit, Prakrit, and Other Indo-Aryan Languages*. New York: Oxford University Press.
- . 2000. *A Gāndhārī Version of the Rhinoceros Sūtra: British Library Kharoṣṭhī Fragment 5B*. Gandhāran Buddhist Texts 2. Seattle: University of Washington Press.
- . 2001. “‘Gāndhārī Hybrid Sanskrit’: New Sources for the Study of the Sanskritization of Buddhist Literature.” *Indo-Iranian Journal* 44: 241–252.
- Sander, Lore. 1999. “Early Prakrit and Sanskrit Manuscripts from Xinjiang (second to fifth/sixth Centuries C.E.): Paleography and Literary Evidence and Their Relation to Buddhist Schools.” In J. McRae and J. Nattier eds. *Collection of Essays 1993: Buddhism across Boundaries: Chinese Buddhism and the Western Regions*. Taipei: Foguang Cultural Enterprise Co. Ltd.: 61–106.

Comments or Discussion

Please send any responses to this proposal to Andrew Glass (email: asg@u.washington.edu). Please also CC to Richard Salomon (email: rsalomon@u.washington.edu) and Stefan Baums (email: baums@u.washington.edu).