2016-12-31

# Proposal to encode the Old Sogdian script in Unicode

Anshuman Pandey pandey@umich.edu

December 31, 2016

## 1 Introduction

This is a proposal to encode the 'Old Sogdian' script in Unicode. It is a significant revision of the following document:

• L2/15-089 "Preliminary Proposal to Encode the Old Sogdian Script in Unicode"

An ISO proposal summary form is attached. This proposal addresses comments made on previous versions in the following reports:

- L2/16-037 "Recommendations to UTC #146 January 2016 on Script Proposals"
- L2/17-037 "Recommendations to UTC #150 January 2017 on Script Proposals"

A proposed Unicode encoding for the later 'Sogdian' script has been presented in:

• L2/16-371R "Revised proposal to encode the Sogdian script in Unicode"

The present proposal has been reviewed by Nicholas Sims-Williams and Yutaka Yoshida, who are leading scholars of Sogdian studies.

# 2 Background

The proposed Unicode encoding for 'Old Sogdian' encompasses a group of related scripts used in the following records for representing Sogdian (ISO 639: sog), an ancient Eastern Iranian language:

• Kultobe inscriptions The oldest Sogdian records are stone inscriptions found at Kultobe, hereafter 'K', in modern Kazakhstan (see Sims-Williams & Grenet 1998; Grenet, et al 2007). Fourteen inscriptions have been discovered and studied (see figures 26, 27). They have not been concisely dated, but the archaic features of the script and language indicate that they precede the 'Ancient Letters'.

- 'Ancient Letters' The earliest attested Sogdian manuscripts are known as the 'Ancient Letters' (see figures 28–35), hereafter 'AL'. These paper documents were found in 1907 by Aurel Stein in Dunhuang, western China. Based upon internal evidence the 'AL' may have been written during 312–314 CE (Sims-Williams 1985; Grenet, et al 1998).
- *Upper Indus inscriptions* Sogdian text appears on more than 600 rock carvings at Shatial and other sites in the Gilgit region of Pakistan (see figures 36, 37). These 'Upper Indus inscriptions', hereafter 'UII', have been dated to the 4th–7th centuries CE (Sims-Williams 1989, 2000), and some more precisely to the latter half of the 5th century (Yoshida 2013).
- Short inscriptions on coins and vessels A script resembling that used in AL and UII is used for inscriptions on coins and vessels from the ancient principality of Chach, situated around modern Tashkent, Uzbekistan, and surrounding areas (see figure 39).

The scripts of these records are derived from Imperial Aramaic and exhibit the following features:

- Repertoire Of the 22 letters of the Aramaic alphabet, 20 are attested in the repertoires of these scripts. Analogues for *teth* and *qoph* do not exist. Of these 20, 17 have distinctive representations, while 3 share a resemblance. In AL and UII, the shapes of *daleth* and *ayin* are in general identical to *resh*, but may be distinctive in K. The letter *taw* has a unique final form in K. All 20 letters are exhibited in K 4 and occur collectively in AL. The AL contain additional letters that do not occur in K, such as distinctive final forms of *aleph*, *beth*, *nun*, *sadhe*, *taw*; special forms of *ayin*; and a new form of *he* (see § 3.1). Numerical signs are attested in AL and UII, but not in K.
- Letterforms The shapes of letters in AL and UII are nearly identical. The letterforms of K are more archaic and reflect constraints imposed by the method and medium of inscription. The shapes of gimel, he, yodh, lamedh, shin in the three varieties differ from the Aramaic originals and corresponding letters in related Iranian scripts. They may be considered characteristically 'Sogdian'. The special forms of ayin in AL do not occur in K or UII, or in any other script. A comparison of letters in related scripts is shown in table 1 and figure 42.
- Structure Each variety is a non-joining abjad, similar to Hebrew. Letters retain their shapes within a word. Some letters have distinctive word-final forms, but there are no formal conventions for their usage. The strokes of adjacent letters of a word may connect or overlap as the result of cursive writing. This type of conjunction differs from that of later 'formal' and 'cursive' Sogdian scripts, which possess intrinsic conjoining behaviors similar to Arabic, as shown below:

		Old Sogdian	Later Sogdian
swyδyk	'Sogdian'	מכאלגע	ويمسلاس
sm>rknδc	'of Samarkand'	מאצגגלרב	ويويعديك

• *Directionality* These old Sogdian varieties are written from right to left in lines that advance from top to bottom. Some UII are written vertically with letters rotated 90° counter-clockwise with lines that advance from left to right (see § 4.5).

These scripts may be considered typologically identical on the basis of their graphical and structural features. For purposes of character encoding they may be unified within a single Unicode script block. Using this approach texts would be represented using the same character set, but the display would be managed through the selection of fonts designed specifically for the K, AL, and UII varieties.

The proposed Unicode block is named 'Old Sogdian'. This identifier has been selected because proper names do not exist for individual script varieties or for the family. The script of AL has been referred to as "Sogdian Aramaic" (Skjærvø 1996), which may be applied applied equally to the other two varieties. However, the descriptor 'Aramaic' is not used in Unicode names for other scripts descended from Aramaic. The bare name 'Sogdian' is used in the catalogue of the International Dunhuang Project for referring to both early and later script varieties. It is, however, practical to reserve this name for a Unicode block for the more well-known 'formal' and 'cursive' styles, which have been proposed for encoding in a unified 'Sogdian' block (see L2/16-371). The designation 'Old Sogdian' suitably identifies these early varieties while emphasizing their genetic relationship with later 'Sogdian' script styles.

# 3 Character Repertoire

The proposed repertoire contains 40 characters: 29 letters, 10 numbers, 1 heterogram. Names for letters correspond to those of the 'Imperial Aramaic' block. Representative glyphs are based upon forms in the AL unless specified below. The encoded set may differ from traditional and scholarly inventories of script varieties that occur in written and inscriptional sources. Such differences naturally arise from the requirements for digitally representing a script in plain text and for preserving the semantics of characters.

In this document, names in italics refer to scholarly names for graphemes while names in small capitals refer to proposed Unicode characters, eg.  $\checkmark$  is aleph and OLD SOGDIAN LETTER ALEPH. For sake of brevity, the descriptor 'OLD SOGDIAN' is dropped when refering to Old Sogdian characters, eg. OLD SOGDIAN LETTER ALEPH is referred to as ALEPH. Characters of other scripts are designated by their full Unicode names. Latin transliteration of Old Sogdian letters follows the scholarly convention. Aramaic heterograms are transliterated using the corresponding uppercase letters, with some exceptions as shown in the table below.

#### 3.1 Letters

Glyph	Character name	Latin
×	OLD SOGDIAN LETTER ALEPH	)
	OLD SOGDIAN LETTER FINAL ALEPH	)
3	OLD SOGDIAN LETTER BETH	$\beta$ ; B
د	OLD SOGDIAN LETTER FINAL BETH	$\beta$ ; B
и	OLD SOGDIAN LETTER GIMEL	$\gamma;G$
я	OLD SOGDIAN LETTER HE	h
ے	OLD SOGDIAN LETTER FINAL HE	h

2	OLD SOGDIAN LETTER WAW	W
J	OLD SOGDIAN LETTER ZAYIN	z
N	OLD SOGDIAN LETTER HETH	х; Ӊ
5	OLD SOGDIAN LETTER YODH	у
У	OLD SOGDIAN LETTER KAPH	k
7	OLD SOGDIAN LETTER LAMEDH	$\delta;L$
*	OLD SOGDIAN LETTER MEM	m
J	OLD SOGDIAN LETTER NUN	n
٦	OLD SOGDIAN LETTER FINAL NUN	n
1	OLD SOGDIAN LETTER FINAL NUN WITH VERTICAL TAIL	n
n	OLD SOGDIAN LETTER SAMEKH	S
5	OLD SOGDIAN LETTER AYIN	C
<b>५</b> ७	OLD SOGDIAN LETTER ALTERNATE AYIN	C
9	OLD SOGDIAN LETTER PE	p
٠	OLD SOGDIAN LETTER SADHE	c
ے	OLD SOGDIAN LETTER FINAL SADHE	c ; Ṣ
٢	OLD SOGDIAN LETTER FINAL SADHE WITH VERTICAL TAIL	c
У	OLD SOGDIAN LETTER RESH-DALETH-AYIN	r, d, <sup>c</sup>
n	OLD SOGDIAN LETTER SHIN	š
מ	OLD SOGDIAN LETTER TAW	t
ىر	OLD SOGDIAN LETTER FINAL TAW	t
ק	OLD SOGDIAN LETTER FINAL TAW WITH VERTICAL TAIL	t

# 3.1.1 Notes on letters

**aleph** In word-final positions in AL, *aleph* is written as  $\perp$  FINAL ALEPH, in which the horizontal stroke at the baseline is elongated. The letter  $\perp$  ALEPH has the shape  $\bowtie$  in K. This form is a glyphic variant. See figure 1 for attestations.

**beth** In word-final positions in AL, **S** BETH is written as **S** FINAL BETH, in which the horizontal stroke at the baseline is elongated. See figure 2 for attestations.

*gimel* See figure 3 for attestations of **▶** GIMEL.

waw See figure 6 for attestations of **3** waw.

zayin See figures 7 and 8 for attestations of J ZAYIN.

**heth** See figure 9 for attestations of N HETH.

teth An Old Sogdian analogue for Aramaic teth does not exist. In K, the teth in Aramaic heterograms is represented using א TAW: QTLt is written as אמלפן KTLt (K 3.3).

*yodh* See figure 10 for attestations of **5** YODH.

**kaph** See figure 11 for attestations of **y** KAPH.

**lamedh** The letter LAMEDH has the shape  $\Delta$  in K and  $\Delta$  in AL (see figure 12). The AL form is the representative glyph. In AL 5, *lamedh* appears as as  $\Delta$ . Differences between  $\Delta$ ,  $\Delta$ ,  $\Delta$  are stylistic, not semantic. The forms  $\Delta$  and  $\Delta$  are to be treated as glyphic variants of  $\Delta$ .

mem See figure 13 for attestations of ≯ MEM.

nun Occurrences of nun are represented using J Nun, J Final Nun,  $\uparrow$  Final Nun with vertical tail (see figure 14). The representative glyph J for Nun is derived from K. The final forms occur only in AL. While nun has the distinctive shape J in K, it has the shape J in AL when non-final, which is identical to J ZAYIN, eg. ZNH (K 4.1) and ZNH (AL 2.10). When word-final in AL, nun is written as both J and J, eg. JNM (AL 2.2) and JNM (AL 2.6). The regular and final forms are contrastive in AL (see figure 8). They are not glyphic variants. All three characters are required for fully representing nun in plain text.

**samekh** The letter >> samekh occurs as the two-part form of in K 4. This archaic form is to be treated as a glyphic variant. See figure 15 for attestations.

ayin The letter ayin occurs only in Aramaic heterograms. It has the regular shape ש and the special shapes and so (see figure 16). The regular ש ayin occurs in both K and AL, eg. שלוע 'BDt (K 4.1), אלוע

LZK (AL 2.12), L (AL 6.6). In AL, the shape of regular ayin is identical to resh (and daleth). In K, there is a possibility that ayin might be a distinctive letter. The similarity between ayin and resh is inherited from Aramaic, compare Aramaic letter ayin and Aramaic letter resh. However, there is insufficient information for determining whether or not the differences between ayin and resh in K are semantically significant. Therefore, a separate character for regular ayin is not proposed at present. It is to be represented using Resh-ayin-daleth. The letters Ayin and Alternate ayin occur only in AL for writing the heterogram D, eg. La (AL 2.1), La (AL 3 verso), La (AL 3.1), La (AL 5.1). Although Alternate ayin, it is appropriate to define two characters on account of their graphical structures. The is a glyphic variant of with an ornate tail; the is a variant with a truncated tail. These three forms are unified as Alternate ayin, which may be used for representing these special forms in plain text. See figure 25 and § 3.3 for attestations.

pe The letter **9** PE is has the variant 'open' shape **9**, which is a glyphic variant (see figure 17).

sadhe This letter is represented using באבול SADHE, באבול SADHE, and באבול SADHE WITH VERTICAL TAIL (see figure 18). The final forms occur only in AL. In AL 2, sadhe has the shape whenever it occurs at the margin, eg. באבול אבול אבול אבול SADHE, eg. באבול אבול SADHE, eg. באבול SADHE, eg. באבול

qoph An Old Sogdian analogue for Aramaic qoph does not exist. In K, the qoph in Aramaic heterograms is represented using א KAPH: QTLt is written as א KTLt (K 3.3). It used to be believed that א qoph was retained in AL as p and reassigned for the number 100. This p is now identified as the fraction ½ (Grenet, et al 1998).

resh In AL, the letter y is used for resh, daleth, and ayin (see figure 19). According to the Unicode character-glyph model, letters with identical glyphic representations are considered variants and are unified as a single character. As the sound [r] represented by resh is phonemic in Sogdian, and those represented by ayin and daleth are not, the letter y is used ubiquitously for resh. Accordingly, daleth and ayin are unified with resh as y resh-ayin-daleth. This approach follows the Unicode model for Inscriptional Pahlavi, in which waw, ayin, resh are represented using 2 U+10B65 INSCRIPTIONAL PAHLAVI LETTER WAW-AYIN-RESH; and similarly, mem and qoph using y U+10B6C INSCRIPTIONAL PAHLAVI LETTER MEM-QOPH. Despite occurring after daleth and ayin in the alphabetical order, resh is ordered first in the name resh-ayin-daleth because it occurs more frequently in the sources; daleth is ordered before ayin for the same reason.

shin See figure 20 for attestations of > SHIN.

This letter is represented using the taw, it is final taw as it is final taw as it is final taw. It is final taw as it is final taw is often written using a glyphic variant with a curved tail in AL. All three characters are required for fully representing taw in plain text.

#### 3.1.2 Note on final forms

Distinctive final forms of *aleph*, *beth*, *nun*, *sadhe*, *taw* are included in the repertoire as separate characters. These final forms differ from the nominal forms in the shape of their terminals, which are elongated horizontally or which descend vertically. An analysis of AL indicates that final forms are regularly used at the

end of words, and that some final forms are used specifically at the end of line. The analysis also suggests an intentional differentiation between nominal and final forms of only these five letters. For instance, the elongated baseline in the final form  $\bot$  of  $\bot$  aleph and the final form  $\bot$  of  $\bot$  beth may be interpreted as a natural flourish made by the scribe at the end of a word. But, such strokes occur consistently with final forms of these two letters across the AL corpus, and are not simply stylistic. On the other hand, commonly occuring letters such as  $\bigcirc$  waw,  $\bigcirc$  yodh,  $\bigcirc$  resh have curved terminal strokes that present a natural opportunity for stylistic elongation at end of a word. However, there appears to be deliberate avoidance of such flourishes when writing these letters in final position. These letters, in turn, may be compared to  $\bigcirc$  kaph and  $\bigcirc$  pe, whose shapes inherently possess an elongated tail that is often extended in final position, and which may be considered a natural stylistic flourish.

In addition to illustrating distinctive final forms for aleph, beth, nun, sadhe, taw, the available sources also point to the existence of two types of final forms for nun, sadhe, taw. These three letters occur in word-final position with either an elongated horizontal stroke or with a descending vertical tail. There is some evidence to suggest contrastive contextual usage of the two types. For instance, word-final radhe is written in AL 2 with a vertical tail radhe whenever it occurs at the margin, and with a horizontal tail in other positions within a line (see figure 18). Throughout the AL corpus, I nun is written using both in and if at the ends of words and lines. The same applies to the usage of in and if a taw. These final forms of nun and taw appear to be used interchangably and occur on the same line or in close proximity. It may be possible that one form was intentionally selected over the other based upon spacing requirements along a line. For instance, a scribe may have chosen the form with a horizontal tail to fill space, or the form with vertical tail to compensate for lack of space. The usage of both forms within the same source suggests that scribes perceived of a semantic distinction between the horizontal and vertical final forms for nun, sadhe, taw.

It is difficult to ascertain the nature of Sogdian scribal conventions that were in vogue in the early 4th century CE, when the AL were written. There are no sources that provide descriptions of orthographic rules or explanations for the existence of final forms for only five letters of the repertoire. There are no clues that offer insights into the development of two final forms for *nun*, *sadhe*, *taw*; or, that specify the rationale for their usage or the criteria for a scribe's preference of one form over the other in a given context. The available sources simply show that both final forms are used for these three letters.

For this reason, the two final forms for *nun*, *sadhe*, *taw* have been included as separate characters in the proposed repertoire. Without knowledge of the conventions for usage of the two forms, it is impractical to exclude one set from the repertoire. Moreover, given that there is some evidence to suggest scribal preferences for a particular form in a given context, it is improper to consider the forms as stylistic variants instead of semantic alternates. Furthermore, in terms of the Unicode character-glyph model it is difficult to specify which of final 1/nun, 1/nun, 1/nun, 1/nun, and an an antiferral sequence of the sequence of the

When developing a Unicode encoding for an ancient script such as Old Sogdian, it is most practical to permit the extant sources to guide the process. This is especially important when there is an absence of knowledge regarding a particular orthographic convention in such a script. It is from this perspective that the proposed repertoire for Old Sogdian has been defined. The proposal author has discussed the issue with scholars, who will be the primary users of the encoding. These experts have expressed a requirement for representing in plain text both vertical and final forms of final *nun*, *sadhe*, and *taw* as they occur in the sources in order to accurately and completely digitize Sogdian records.

## 3.2 Numbers

The repertoire contains 10 numerical characters. These occur in AL and UII, but not in the extant K sources. See figures 22–24 for attestations.

Glyph	Character name	Numeric value
J	OLD SOGDIAN NUMBER ONE	1
ມ	OLD SOGDIAN NUMBER TWO	2
m	OLD SOGDIAN NUMBER THREE	3
m	OLD SOGDIAN NUMBER FOUR	4
mm	OLD SOGDIAN NUMBER FIVE	5
٥	OLD SOGDIAN NUMBER TEN	10
3	OLD SOGDIAN NUMBER TWENTY	20
j	OLD SOGDIAN NUMBER THIRTY	30
ھ	OLD SOGDIAN NUMBER ONE HUNDRED	100
P	OLD SOGDIAN FRACTION ONE HALF	1/2

Primary units The primary units are expressed using joined repetitions of the sign J that are generally grouped in sets of three or four and separated by spaces, eg. J for 2, J for 3, J for 8. As the script is non-conjoining, no simple method exists for representing the ligated repetitions of J. For that reason, the numbers J one .. J four are encoded atomically. This model for one .. Four follows the Unicode encoding for Inscriptional Parthian, eg. J U+10B58 INSCRIPTIONAL PARTHIAN NUMBER ONE .. J U+10B5B INSCRIPTIONAL PARTHIAN NUMBER FOUR. The numbers 5–9 are written using sequences of one .. Four arranged in groups separated by spaces. The number 5 may also be represented using the character J FIVE, which is attested as a single unit in AL 7 (see figure 23). Representations of all primary numbers are shown in the table below.

Tens The 3 TEN resembles a vertically compressed 3 LAMEDH. The shapes for 3 TWENTY and 3 THIRTY are formed from vertical stacks of 3 TEN. Multiples of 10 greater than 20 are produced using appropriate repetitions and groupings of TEN and TWENTY. The character THIRTY is not commonly used in compound numbers. The number 30 may be also represented as 33, which is a compound of TWENTY and TEN.

Hundreds The number 100 is written using ONE HUNDRED. The glyph resembles the letter of GIMEL above a serpentine form, but it is an atomic character. The ONE HUNDRED also functions as a unit mark for the hundreds. Multiples of hundred are represented by prefixing the appropriate groupings of ONE .. FOUR before ONE HUNDRED.

Thousands The number 1000 is expressed using the Aramaic heterogram  $\mathcal{L}P$ , which is represented using the sequence  $\mathcal{L}P$  one,  $\mathcal{L}P$  LAMEDH,  $\mathcal{L}P$  PE>. The sequence  $\mathcal{L}P$  also functions as a unit mark for the thousands. The  $\mathcal{L}P$  one is an inherent part of the  $\mathcal{L}P$  unit. Multiples are expressed by prefixing primary numbers before the unit, eg. 2000 is  $\mathcal{L}P$ , 3000 is  $\mathcal{L}P$ .

Ten thousands The number 10000 is expressed using the Sogdian word צעאכע  $\beta rywr$ . There is no distinctive numerical sign for this value.

**Fraction** The **/** FRACTION ONE HALF is placed after another numerical character.

# 3.2.1 Notation system

The ordering of numbers follows the right-to-left directionality of the script. The expression of numbers is additive. Compounds of different units are produced by placing larger units first. However, in some inscriptions on silver coins the units precede the tens (see Livshits 2015: 234), which follows the order of spoken numbers. Spaces are used for separating groups of primary numbers.

Value	Number	Input string →
4½	hm	< <b>₩</b> FOUR, <b>/</b> FRACTION ONE HALF>
5	m m	THREE, SP SPACE, U TWO>
5	mm	<wu>mu FIVE&gt;</wu>
6	m m	<m [sp]="" m="" space,="" three="" three,=""></m>
7	m m	<pre><pre>FOUR, [sp] SPACE, m THREE&gt;</pre></pre>
$7\frac{1}{2}$	<b>/</b> ա ա	$<$ <b>uu</b> four, $[sp]$ space, <b>u</b> three, $\nearrow$ fraction one half $>$
8	ım ım	<pre><pre>FOUR, [sp] SPACE, pu FOUR&gt;</pre></pre>
9	ա ա ա	$<$ <b><math>m{u}</math></b> THREE, $[sp]$ SPACE, $m{u}$ THREE, $[sp]$ SPACE, $m{u}$ THREE $>$
13	m2	< <b>&gt;</b> TEN, <b>m</b> THREE>
15	Zum 2	< <b>&gt;</b> TEN, <b>J</b> FIVE>
30	3	< <b>3</b> THIRTY>
30	23	< <b>3</b> TWENTY, <b>5</b> TEN>
32	εςμ	< <b>3</b> TWENTY, <b>5</b> TEN, <b>11</b> TWO>
100	4	< <b>~</b> ONE HUNDRED>

200	רבה ח	<ul> <li>TWO, [SP] SPACE, CONE HUNDRED&gt;</li> </ul>
500	رهاس س	$<$ <b><math>\mathbf{u} THREE, <math>[\mathbf{s}P]</math> SPACE, <math>\mathbf{u}</math> TWO, <math>\mathbf{u}</math> ONE HUNDRED<math>&gt;</math></math></b>
1000	ىكو	$<$ <b>J</b> ONE, $\checkmark$ LAMEDH, $\checkmark$ PE $>$
2000	ىر درو	< <b>u</b> two, $[sp]$ space, <b>J</b> one, $2$ lamedh, <b>9</b> pe $>$
10000	בעאכע	<ul><li>S BETH, Y RESH, 5 YODH, 7 WAW, Y RESH&gt;</li></ul>

Attestations for the above numbers are shown in figures 22–24. The repertoire provides for the presentation of any numerical value, even if not attested. For example, the number 2453 could be represented as:

Value	Number	Input string →
2453	π ۱ر6 سπک ξεπ	Y TWO, SP SPACE, J ONE, LAMEDH, PE, SP SPACE, LAMEDH, PE, SP SPACE, SP SPACE, SP SPACE, ST THIRTY, S TWENTY, W THREE
2453	u دره سری ۶۶دس	<u \(="" \)="" \mathfrak{9}="" j="" lamedh,="" one,="" p="" pe,="" sp="" space,="" space,<="" two,=""> um four, \( \mathfrak{\mathfrak{M}} \) one hundred, sp space, 3 twenty, \( 3 \) twenty, \( 3 \) ten, \( m \) three&gt;</u>

#### 3.3 Heterogram

The repertoire contains 1 heterogram.

Glyph	Character name	Value
70	OLD SOGDIAN HETEROGRAM AYIN-DALETH	Ъ

Aramaic heterograms are represented as words spelled using conventional letters, eg. 'HRZY is written  $\leq$  ALEPH, M HETH, Y RESH, J ZAYIN, 5 YODH>. The heterogram 'D is the sole exception. Meaning "to", 'D occurs in the address and salutation of a letter, eg. 'P by w xwt'w  $\beta$ 'rkkw" (to lord master Barak"). Morphologically, it is comprised of ayin and daleth. Yet, instead of the expected spelling \*Y  $\leq$  RESH-AYIN-DALETH, Y RESH-AYIN-DALETH>, the ayin is written using special forms: Y, N, Y, Y, (see figure 25). An explanation for this curious orthography may be that ayin and daleth had disappeared from the script by the time of AL, and the original phonetic values of these letters never existed in Sogdian. Therefore, scribes were unaware of these letters and of the original spelling of the Aramaic word, so they stylized the writing of 'D (Sims-Williams, personal correspondence, 2016).

There are two ways to analyze these representations of 'D. First, as a conventional word comprised of the letters *ayin* and *daleth*. These forms of *ayin*, which occur only in this heterogram, are included in the repertoire as  $\longrightarrow$  AYIN and  $\Longrightarrow$  ALTERNATE AYIN; the  $\smile$  and  $\Longrightarrow$  could be considered glyphic variants of ALTERNATE

AYIN. Accordingly, 'D may be represented as <AYIN | ALTERNATE AYIN, RESH-AYIN-DALETH>. Secondly, 'D is a logographic unit comprised of a ligature or a set of two letters. This unit may be treated as an atomic character, eg. TO OLD SOGDIAN HETEROGRAM AYIN-DALETH. These approaches are not mutually exclusive and both are practical for character encoding. Depending upon the context, 'D may be spelled using a sequence of letters or represented using an atomic character.

The case of 'D is similar to the Latin '&' ampersand. The '&' represents the Latin word et "and". Morphologically, it is a ligation of the Latin letters 'e' and 't', eg. e7, & The base letters began to be obscured as the ligature became more stylized, eg. & The logographic nature of '&' is apparent in the abbreviation "&c" for Latin et cetera "and so forth", where it masks 'et'. Latin et can be represented both using the sequence <e, t> and atomic characters, such as e7 U+1F670 SCRIPT LIGATURE ET ORNAMENT.

The character name for HETEROGRAM AYIN-DALETH is derived from the conventional transliteration 'D of the heterogram. The representative glyph **To** is derived from AL 3 and has been selected because it is structurally a ligature. Variant forms may be managed through fonts.

# 4 Script Details

#### 4.1 Bidirectional model

Old Sogdian may be implemented using the Unicode Bidirectional Algorithm. There are no requirements for shaping.

#### 4.2 Punctuation

Punctuation marks are not attested. Words are separated using spaces in K and AL. Inter-word spacing is inconsistent in the UII.

## 4.3 Line-breaking

There are no rules for line-breaking. The available sources show line-breaks after the end of a word. Word are not split across lines. Consequently, hyphens or other continuation marks are not attested. In digital layouts, line-breaks may occur after any character.

#### 4.4 Collation

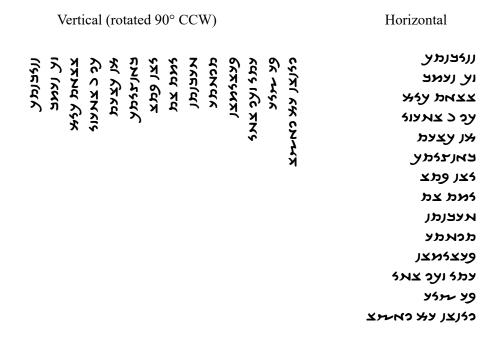
The sort order for Old Sogdian is as follows:

```
ightharpoonup resh-ayin-daleth < 
ightharpoonup shin < 
ightharpoonup taw << 
ightharpoonup final taw with vertical tail
```

#### 4.5 Vertical text

The majority of Old Sogdian records have horizontal orientations. Some UII records are inscribed vertically (Yoshida 2013). There are no formal conventions for text orientation. However, in vertical environments, Old Sogdian text is oriented from top to bottom with lines that advance from left to right. Letters are rotated 90° counter-clockwise from their regular upright shapes.

By default, Old Sogdian may be oriented horizontally in plain text representations. However, support for vertical orientations of the script is required for accurately displaying Old Sogdian text that is natively vertical. Below is a vertical text from Shatial rock 36:38 (see figure 38) and its horizontal representation:



nnyβntk / ZK nrsβ / '''yt kym / kw 10 'ḤRZY / MN k 'rt / βγncytk / y 'n pt '[-] / [-]yst 't / xrβntn / twxtr / pr'ys'n / rty ZKw 'ḤY / pr šyr / wyn 'n 'M wyš'

"(I), Nanai-vandak the (son of) Narisaf have come (here) in/on the (day/year) ten and asked a boon from the spirit of the sacred place Kârt (that) I may arrive at Kharvandan (= Tashkurgan) very quickly and see (my) brother in good (health) with joy." (Yoshida 2013: 379–380).

The "Unicode Technical Report #50: Unicode Vertical Text Layout" describes the <code>vertical\_orientation</code> (vo) property for specifying the orientation of characters in vertical environments. For Old Sogdian, this property would be defined as: <code>vertical\_orientation=R</code> or <code>vo=R</code>, where the value 'R' indicates that the glyphs are rotated in vertical layout. The rotation is 90° counter-clockwise.

## 4.6 Character Data

# 4.6.1 Character properties

In the format of UnicodeData.txt:

```
10F00;OLD SOGDIAN LETTER ALEPH;Lo;0;R;;;;N;;;;
10F01; OLD SOGDIAN LETTER FINAL ALEPH; Lo; 0; R;;;;; N;;;;;
10F01; OLD SOGDIAN LETTER BETH; Lo; 0; R;;;;; N;;;;;
10F03;OLD SOGDIAN LETTER FINAL BETH;Lo;0;R;;;;;N;;;;;
10F04; OLD SOGDIAN LETTER GIMEL; Lo; 0; R;;;;; N;;;;;
10F05; OLD SOGDIAN LETTER HE; Lo; 0; R;;;;; N;;;;;
10F06; OLD SOGDIAN LETTER FINAL HE; Lo; 0; R;;;; N;;;;
10F07; OLD SOGDIAN LETTER WAW; Lo; 0; R;;;;; N;;;;;
10F08; OLD SOGDIAN LETTER ZAYIN; Lo; 0; R;;;;; N;;;;;
10F09; OLD SOGDIAN LETTER HETH; Lo; 0; R;;;;; N;;;;;
10F0A; OLD SOGDIAN LETTER YODH; Lo; 0; R;;;;; N;;;;;
10F0B; OLD SOGDIAN LETTER KAPH; Lo; 0; R;;;;; N;;;;;
10F0C; OLD SOGDIAN LETTER LAMEDH; Lo; 0; R;;;;; N;;;;
10F0D; OLD SOGDIAN LETTER MEM; Lo; 0; R;;;;; N;;;;
10F0E; OLD SOGDIAN LETTER NUN; Lo; 0; R;;;;; N;;;;;
10F0F;OLD SOGDIAN LETTER FINAL NUN; Lo; 0; R;;;;; N;;;;;
10F10; OLD SOGDIAN LETTER FINAL NUN WITH VERTICAL TAIL; Lo; 0; R;;;;; N;;;;;
10F11; OLD SOGDIAN LETTER SAMEKH; Lo; 0; R;;;;; N;;;;;
10F12; OLD SOGDIAN LETTER AYIN; Lo; 0; R;;;;; N;;;;
10F13;OLD SOGDIAN LETTER ALTERNATE AYIN;Lo;0;R;;;;;N;;;;;
10F14;OLD SOGDIAN LETTER PE;Lo;0;R;;;;N;;;;
10F15;OLD SOGDIAN LETTER SADHE;Lo;0;R;;;;;N;;;;;
10F16; OLD SOGDIAN LETTER FINAL SADHE; Lo; 0; R;;;;; N;;;;;
10F17; OLD SOGDIAN LETTER FINAL SADHE WITH VERTICAL TAIL; Lo; 0; R;;;;; N;;;;;
10F18; OLD SOGDIAN LETTER RESH-AYIN-DALETH; Lo; 0; R;;;;; N;;;;
10F19; OLD SOGDIAN LETTER SHIN; Lo; 0; R;;;;; N;;;;;
10F1A; OLD SOGDIAN LETTER TAW; Lo; 0; R;;;;; N;;;;
10F1B;OLD SOGDIAN LETTER FINAL TAW;Lo;0;R;;;;;N;;;;
10F1C;OLD SOGDIAN LETTER FINAL TAW WITH VERTICAL TAIL; Lo; 0; R;;;;; N;;;;;
10F1D; OLD SOGDIAN NUMBER ONE; No; 0; R;;;; 1; N;;;;
10F1E; OLD SOGDIAN NUMBER TWO; No; 0; R;;;; 2; N;;;;;
10F1F; OLD SOGDIAN NUMBER THREE; No; 0; R;;;; 3; N;;;;;
10F20; OLD SOGDIAN NUMBER FOUR; No; 0; R;;;; 4; N;;;;;
10F21;OLD SOGDIAN NUMBER FIVE; No; 0; R;;;; 5; N;;;;;
10F22; OLD SOGDIAN NUMBER TEN; No; 0; R;;;; 10; N;;;;;
10F23;OLD SOGDIAN NUMBER TWENTY; No; 0; R;;;; 20; N;;;;;
10F24; OLD SOGDIAN NUMBER THIRTY; No; 0; R;;;; 30; N;;;;;
10F25;OLD SOGDIAN NUMBER ONE HUNDRED; No;0;R;;;;100;N;;;;;
10F26; OLD SOGDIAN FRACTION ONE HALF;; No; 0; R;;;; 1/2; N;;;;
10F27; OLD SOGDIAN HETEROGRAM AYIN-DALETH; Lo; 0; R; ;; ;; N; ;; ;;
```

## 4.6.2 Linebreaking

In the format of LineBreak.txt:

```
10F00..10F1C;AL  # Lo [29] OLD SOGDIAN LETTER ALEPH..
OLD SOGDIAN LETTER FINAL TAW WITH VERTICAL TAIL
10F1D..10F26;AL  # No [10] OLD SOGDIAN NUMBER ONE..OLD SOGDIAN FRACTION ONE HALF
10F27;AL  # Lo OLD SOGDIAN HETEROGRAM AYIN-DALETH
```

## 5 References

- Anderson, Deborah; et. al. 2016. "Recommendations to UTC #146 January 2016 on Script Proposals". L2/16-037. http://www.unicode.org/L2/L2016/16037-script-rec.pdf
- ——. 2017. "Recommendations to UTC #150 January 2017 on Script Proposals". L2/17-037. http://www.unicode.org/L2/L2017/17037-script-ad-hoc.pdf
- Grenet, Frantz; Nicholas Sims-Williams; Étienne de La Vaissière. 1998. "The Sogdian Ancient Letter V". *Bulletin of the Asia Institute*, Alexander's Legacy in the East: Studies in Honor of Paul Bernard, New series, vol. 12, ed. Osmund Bopearachchi, Carol Altman Bromberg, and Frantz Grenet, pp. 91–104.
- Grenet, Frantz; Nicholas Sims-Williams; Aleksandr Podushkin. 2007. "Les plus anciens monuments de la langue sogdienne: les inscriptions de Kultobe au Kazakhstan". *Comptes rendus des séances de l'Académie des Inscriptions et Belles-Lettres*, 151° année, N. 2, 2007, pp. 1005–1034.
- Pandey, Anshuman. 2016. "Revised proposal to encode the Sogdian script in Unicode" (L2/16-371). http://www.unicode.org/L2/L2016/16371-sogdian.pdf
- Reichelt, Hans. 1928–31. *Die soghdischen Handschriftenreste des Britischen Museums*. Heidelberg: C. Winter's Universitätsbuchhandlung.
- Sims-Williams, Nicholas. 1975. "Notes on Sogdian Palaeography". *Bulletin of the School of Oriental and African Studies, University of London*, vol. 38, no. 1 (1975), pp. 132–139.
- ——. 1981a. "The Sogdian sound-system and the origins of the Uyghur script". *Journal Asiatique*, pp. 347–360.
- ——. 1981b. "Remarks on the Sogdian letters γ and x (with special reference to the orthography of the Sogdian version of the Manichean church-history)", Appendix in W. Sundermann, *Mitteliranische manichäische Texte kirchen-geschichtlichen Inhalts*, Berliner Turfantexte, XI, Berlin: Akademie-Verlag, pp. 194–198.
- ——. 1985. "Ancient Letters". *Encyclopædia Iranica*, vol. II, fasc. 1, pp. 7–9. http://www.iranicaonline.org/articles/ancient-letters
- . 1989. Sogdian and Other Iranian Inscriptions of the Upper Indus. Corpus Inscriptionum Iranicarum, pt. II (Inscriptions of the Seleucid and Parthian Periods and of Eastern Iran and Central Asia), v. III (Sogdian), no. I. London: Published on behalf of Corpus Inscriptionum Iranicarum by School of Oriental and African Studies.
- ——. 2000. "The Iranian Inscriptions of Shatial". *Indologica Taurinensia*, v. 23–24 (Professor Gregory M. Bongard-Levin Felicitation Volume), pp. 523–541.
- Sims-Williams, Nicholas; Frantz Grenet. 2007. "The Sogdian Inscriptions of Kultobe". *Shygys*, 2006, vol. 1 pp. 95–111.
- Skjærvø, Prods Oktor. 1996. "Aramaic Scripts for Iranian Languages." *The World's Writing Systems*, edited by Peter T. Daniels and W. Bright, pp. 515–535. New York and Oxford: Oxford University Press.

Waugh, Daniel C. [comp]. 2004. "The Sogdian Ancient Letters", translated by Nicholas Sims-Williams. https://depts.washington.edu/silkroad/texts/sogdlet.html

Yoshida, Yutaka. 2002. "In search of traces of the Sogdians 'Phoenicians of the Silk Road". *Berlin-Brandenburgische Akademie der Wissenschaften*. Berichte und Abhandlungen, Band 9, p. 185–200. Berlin.

——. 2013. "When Did Sogdians Begin to Write Vertically?". *Tokyo University Linguistic Papers*, vol. 33, pp. 375–394.

# 6 Acknowledgments

I express my deep gratitude to Nicholas Sims-Williams (SOAS, University of London) for providing detailed comments on earlier versions of this proposal and for informative discussions regarding all facets of the script. I am also thankful to Yutaka Yoshida (University of Kyoto) for reviewing versions of this proposal and for providing valuable feedback. I thank them both for their patient responses to my numerous inquiries and for overlooking my ignorance of the script. I am grateful to Roozbeh Pournader (Google, San Francisco) for discussing Unicode encodings for Iranian scripts and for his feedback on the earliest draft proposal.

The present proposal was funded in part by the Adopt-A-Character Program of the Unicode Consortium. A previous version was made possible in part through a Google Research Award, granted to Deborah Anderson for the Script Encoding Initiative, which funded a post-doctoral research position for me in the Department of Linguistics, University of California, Berkeley during 2015–2016. Preliminary research was made possible through the Script Encoding Initiative at Berkeley. Any views, findings, conclusions or recommendations expressed in this publication do not necessarily reflect those of the Unicode Consortium; the University of California, Berkeley; or Google.

	10F0	10F1	10F2
0	<b>4</b>	10F10	<b>10F20</b>
1	10F01	<b>30</b>	<b>10F21</b>
2	<b>5</b>	<b>1</b> 0F12	<b>)</b>
3	<b>1</b> 0F03	<b>59</b>	<b>3</b>
4	10F04	<b>9</b>	<b>3</b>
5	<b>1</b> 0F05	<b>5</b>	10F25
6	10F06	<b>1</b> 0F16	<b>1</b> 0F26
7	<b>1</b> 0F07	<b>1</b> 0F17	<b>10F27</b>
8	<b>J</b>	<b>У</b>	
9	<b>)</b>	<b>)</b>	
Α	<b>5</b>	<b>)</b>	
В	<b>9</b>	10F1B	
С	<b>1</b> 0F0C	<b>P</b>	
D	10F0D	<b>J</b>	
E	<b>J</b>	<b>1</b> 0F1E	
F	10F0F	10F1F	

This block unifies the scripts used in the Ancient Letters and the Kultobe and Upper Indus inscriptions.

## Letters

10F00 

■ OLD SOGDIAN LETTER ALEPH 10F01 OLD SOGDIAN LETTER FINAL ALEPH 10F02 SOLD SOGDIAN LETTER BETH 10F03 🛥 OLD SOGDIAN LETTER FINAL BETH 10F04 • OLD SOGDIAN LETTER GIMEL 10F05 × OLD SOGDIAN LETTER HE 10F06 🗗 OLD SOGDIAN LETTER FINAL HE 10F07 2 OLD SOGDIAN LETTER WAW 10F08 J OLD SOGDIAN LETTER ZAYIN 10F09 N OLD SOGDIAN LETTER HETH 10F0A 4 OLD SOGDIAN LETTER YODH 10F0B y OLD SOGDIAN LETTER KAPH 10F0C 2 OLD SOGDIAN CETTER KAPH OLD SOGDIAN LETTER LAMEDH 10F0D 🤧 OLD SOGDIAN LETTER MEM 10F0E , OLD SOGDIAN LETTER NUN 10F0F → OLD SOGDIAN LETTER FINAL NUN 10F10 OLD SOGDIAN LETTER FINAL NUN WITH VERTICAL TAIL 10F11 > OLD SOGDIAN LETTER SAMEKH 10F12 - OLD SOGDIAN LETTER AYIN • used only in the Aramaic heterogram 'D • resh-ayin-daleth is used in other heterograms 10F13 so OLD SOGDIAN LETTER ALTERNATE AYIN • used only in the Aramaic heterogram 'D • resh-ayin-daleth is used in other heterograms 10F14 **9** OLD SOGDIAN LETTER PE 10F15 🗸 OLD SOGDIAN LETTER SADHE 10F16 - OLD SOGDIAN LETTER FINAL SADHE 10F17 OLD SOGDIAN LETTER FINAL SADHE WITH VERTICAL TAIL ע 10F18 OLD SOGDIAN LETTER RESH-AYIN-DALETH 10F19 - OLD SOGDIAN LETTER SHIN 10F1A DOLD SOGDIAN LETTER TAW OLD SOGDIAN LETTER FINAL TAW 10F1C p OLD SOGDIAN LETTER FINAL TAW WITH VERTICAL TAIL

#### Numbers

10F1D J OLD SOGDIAN NUMBER ONE
10F1E J OLD SOGDIAN NUMBER TWO
10F1F J OLD SOGDIAN NUMBER THREE
10F20 J OLD SOGDIAN NUMBER FOUR
10F21 J OLD SOGDIAN NUMBER FIVE
10F22 J OLD SOGDIAN NUMBER TEN
10F23 J OLD SOGDIAN NUMBER TWENTY
10F24 J OLD SOGDIAN NUMBER THIRTY
10F25 L OLD SOGDIAN NUMBER ONE HUNDRED
10F26 OLD SOGDIAN FRACTION ONE HALF

# Heterogram

10F27 TO OLD SOGDIAN HETEROGRAM AYIN-DALETH
• ligature of the Aramaic heterogram `D

	Old Sogdian	Inscriptional Pahlavi	Inscriptional Parthian	Imperial Aramaic
aleph	∡, _∡	П		*
beth	۳, ی	<b>ل</b>	ح	>
gimel	и	7	J	4
daleth	<b>(Y</b> )	3	کِ	,
he	거, 스	ಆ	$\mathscr{H}$	7)
waw	2	2	2	•
zayin	J	s	ſ	1
heth	N	s.	N	"
teth	_	2	לל	Ø
yodh	5	2	J	٨
kaph	y	1	9	y
lamedh	7	}	5	L
mem	<b>%</b>	ঠ	Я	<b>3</b>
nun	ا, ب, ا	1	<b>ـ</b> ـ	5
samekh	n	n	D	,
ayin	J-, 500, (Y)	(2)	خ	v
pe	9	4	<i>&gt;</i>	,
sadhe	۲, ۲, ۲	٤	_^	77
qoph	_	( <b>&amp;</b> )	מ	ア
resh	У	(2)	9	,
shin	<b>71</b>	22	¥	v
taw	ק,ת, מ	r	ゔ	٢

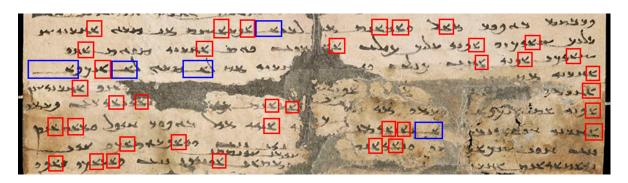
Table 1: Comparison of Old Sogdian letters with those in Unicode blocks for related Iranian scripts and Aramaic. Parenthesis indicate that a letter has been unified with another in the respective encoding. In Inscriptional Pahlavi, *ayin* and *resh* are unified with *waw*, and *qoph* with *mem*. For Old Sogdian, *daleth* and regular *ayin* are unified with *resh*.

	Old Sogdian	Inscriptional Pahlavi	Inscriptional Parthian	Imperial Aramaic
ONE	J	1	J	1
TWO	n	n	IJ	V
THREE	m	m	JJJ	\//
FOUR	mn	m	וווו	_
FIVE	mn	_	_	_
TEN	2	٦	٧	$\neg$
TWENTY	3	3	9	3,
THIRTY	ŧ	_	_	_
ONE HUNDRED	4	ķ	<u>ح</u>	47
ONE THOUSAND	_	ન	ځ	X
TEN THOUSAND	_	_	_	**
ONE HALF	P	_	_	_

Table 2: Comparison of Old Sogdian numerical signs with those in Unicode blocks for related Iranian scripts and Aramaic.



Inscriptional, archaic form **x** of **x** ALEPH (K 4.1–4).

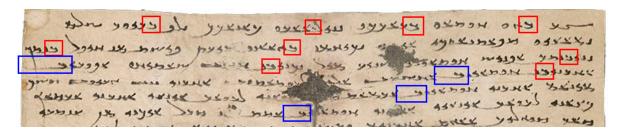


Written forms of  $\checkmark$  Aleph (red) and  $\checkmark$  Final Aleph (blue) (AL 2.1–6).

Figure 1: Specimens of aleph.



Inscriptional forms of **೨** BETH (K 4.1−2).



Written forms of S BETH (red) and S FINAL BETH (blue) (AL 2.1–6).

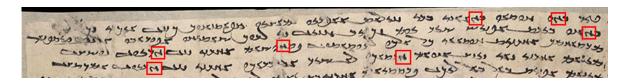
Figure 2: Specimens of beth.



Inscriptional form of  $\kappa$  GIMEL (K 4.6).



Written forms of  $\aleph$  GIMEL (AL 2.7–12).



Written forms of  $\kappa$  GIMEL (AL 3.1–4).

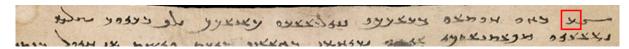
Figure 3: Specimens of gimel



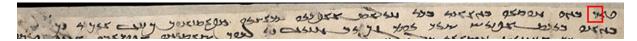
Inscriptional form of daleth in עבעק 'BDt written as  $\mathbf{y}$  (= RESH-AYIN-DALETH) (K 4.1).



Written form of daleth in yso 'D written as y (= RESH-AYIN-DALETH) (AL 1.1).



Written form of *daleth* in South Written as South (= RESH-AYIN-DALETH) (AL 2.1).



Written form of *daleth* in YSO 'D written as Y (= RESH-AYIN-DALETH) (AL 3.1).



The letter *daleth* in yso 'D written as y (= resh-ayin-daleth) (AL 3 verso).



Usage of  $\mathbf{y}$  (= RESH-AYIN-DALETH) for representing *daleth* (blue), *ayin* (green), and *resh* (red) (AL 2.1–12).

Figure 4: Specimens of daleth.

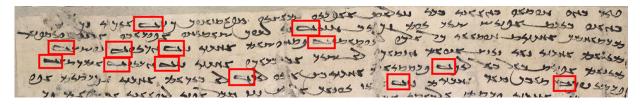


Inscriptional form of א HE in אכעא SWRH and אס(א) (H)WH (K 2.3–4).





Written forms of ح FINAL HE in لاح ZNH and VINAL VINAL



Ubiquitous usage of → FINAL HE in AL 3.1–6.

Figure 5: Specimens of he.



Inscriptional forms of **5** waw (K 4.2–4).



Written forms of **3** waw (AL 2.1–5).

Figure 6: Specimens of waw.



Inscriptional form of J ZAYIN (K 4).



Written form of J ZAYIN (AL 2.34–36).

Figure 7: Specimens of zayin. See also figure 8.



Inscriptional forms of J ZAYIN and J NUN in K 4: JNUN in

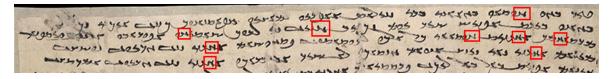


Written forms of J zayin (magenta) and nun at the end of word (AL 2.33–41). Final nun is represented using both J FINAL NUN (green) and J FINAL NUN WITH VERTICAL TAIL (blue).

Figure 8: Comparison of zayin and nun. See also figure 14.



Inscriptional forms of **▶** HETH (K 4.3–7).

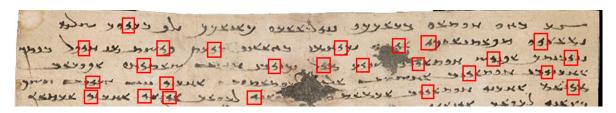


Written forms of N HETH (AL 3.1–4).

Figure 9: Specimens of heth.



Inscriptional forms of 5 YODH (K 4.1–3).

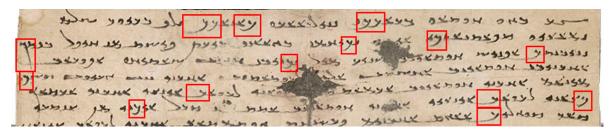


Written forms of 5 YODH (AL 2.1–5).

Figure 10: Specimens of yodh.



Inscriptional forms of **y** KAPH (K 4.1–3).

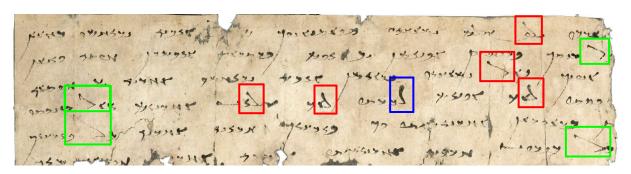


Written forms of y KAPH (AL 2.1–4).

Figure 11: Specimens of *kaph*.



Inscriptional, archaic form  $\mathbf{S}$  of  $\mathbf{S}$  Lamedh (K 4.1).

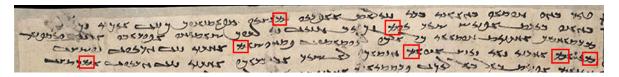


Written forms of  $\Delta$  LAMEDH (red) and its variant forms  $\Delta$  (green) and  $\Delta$  (blue) (AL 6.1–8).

Figure 12: Specimens of lamedh.



Inscriptional forms of → MEM (K 4.1–3).

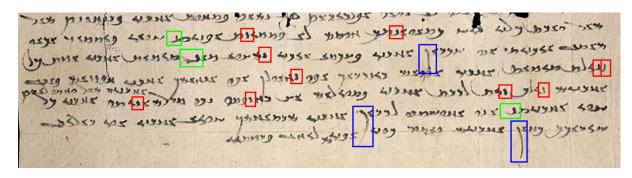


Written forms of > MEM (AL 3.1-4).

Figure 13: Specimens of mem.



Inscriptional form of J NUN (K 4).



Written forms of J NUN (red), → FINAL NUN (green), FINAL NUN WITH VERTICAL TAIL (blue) (AL 1.7–12).



Usage of  $\square$  Final nun (red) and | Final nun with vertical tail (blue) in the word MN:  $\square \bowtie$  and  $|\bowtie$  (AL 2.2–7).

Figure 14: Specimens of *nun*. See also figure 8.



Archaic form Jo of > SAMEKH (K 4.1-4).

معر معدد عدد عدد عدد مرد المدر المد

Written forms of > SAMEKH (AL 1.7–12).

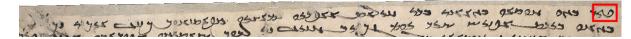
Figure 15: Specimens of samekh.



The letter ayin in Jo Written using Jo AYIN (AL 2.1).



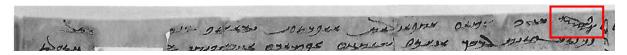
The ayin in you 'D written using so alternate ayin (AL 1.1).



The ayin in yso 'D written using so alternate ayin (AL 3.1).



The ayin in yso 'D written using the glyphic variant of of so alternate ayin (AL 3 verso).



The ayin in Y 'D' written using the glyphic variant of of so alternate ayin (AL 5.1).



The letter ayin in עבעק 'BDt inscribed as ש (= RESH-AYIN-DALETH) (K 4.1).



The letter ayin in עלוע "LZK written using צ (= RESH-AYIN-DALETH) (AL 2.12).



The letter ayin in  $\Delta \mathbf{v}$  Written using  $\mathbf{v}$  (= RESH-AYIN-DALETH) (AL 6.6).

Figure 16: Specimens of ayin.



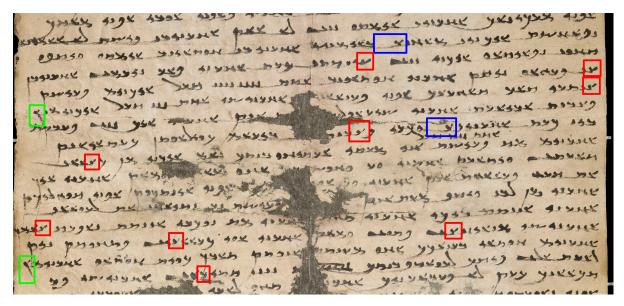
Glyphic variant  $\mathbf{9}$  of  $\mathbf{9}$  PE (K 4.1–6).

Written forms of **9** PE (AL 1.6–12).

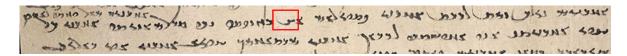
Figure 17: Specimens of pe.



Inscriptional form of **S** SADHE (K 4.1–6).



Written forms of  $\mathbf{r}$  sadhe (red),  $\mathbf{r}$  final sadhe (blue), and  $\mathbf{r}$  final sadhe with vertical tail (green) in (AL 2)



Curved variant of Final sadhe (AL 1.10).

Figure 18: Specimens of sadhe.



Usage of  $\mathbf{y}$  for representing *daleth* (blue), *ayin* (green), and *resh* (red) (AL 2.1–12). As shown,  $\mathbf{y}$  is most commonly used for *resh*. The letter  $\mathbf{y}$  is proposed for encoding as the unified character RESH-AYIN-DALETH.

Figure 19: Comparison of daleth, ayin, and resh.



Inscriptional forms of > SHIN (K 4.1–3).



Written forms of > SHIN (AL 2.1–4).

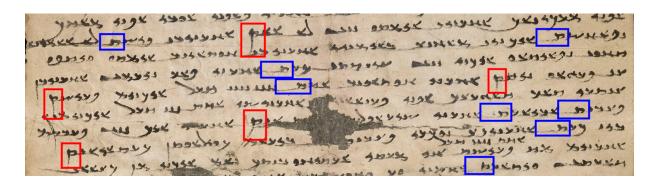


Written forms of > SHIN (AL 3.1–3).

Figure 20: Specimens of shin.



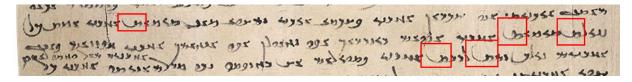
Inscriptional, archaic forms **J5** (red) and **[5** (blue) of **5** TAW and **[5** FINAL TAW WITH VERTICAL TAIL (K 4.1–2). The distinction is apparent in **[5** \*\*Swtt (line 2), which contains both nominal and final forms.



Written forms of Final Taw (blue) and p Final Taw with vertical Tail (red) at the end of word (AL 2.28–36).

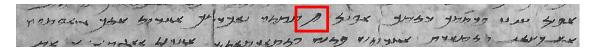


Contrasive usage of לב FINAL TAW and ק FINAL TAW WITH VERTICAL TAIL in two instances of the word prnxwnt: פצואכולם and פצואכולם (AL 1.5–6).



Curved variant → of → FINAL TAW (AL 1.8–10).

Figure 21: Specimens of taw.



The fraction  $\frac{1}{2}$  / (AL 5.10).



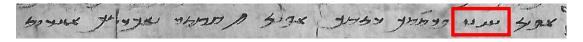
The number  $3 \mu$  (AL 2.32).



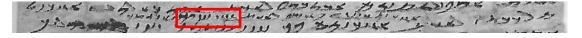
The number 4 **J JJ** (AL 5.26).



The number  $4\frac{1}{2}$  /**L** (AL 5.24).



The number 5 **u u** (AL 5.10).

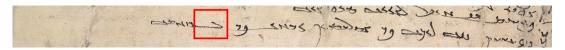


The number 7½ **/w .w.** (AL 5.26).



The number 8 **....** (AL 2.31).

Figure 22: Examples of numbers in the 'Ancient Letters'. See also figures 23 and 24.



The number 10 **>** (AL 3.26).



The number 13 **LL** (AL 2.62).



The number 15 **CAL** 7.8).



The number 20 **3** (AL 5.21).



The number 30 **§** (AL 5.32).



The number 32 **J3** (AL 2.62).

Figure 23: Additional examples of numbers in the 'Ancient Letters'. See also figures 22 and 24.



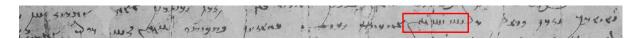
The number 100 🕰 (AL 2.19).



The number 200 **u u** (AL 7.3).



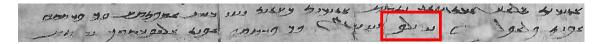
The number 500 **w w** (AL 5.9).



The number 800 **(AL** 4.3).



The number 1000 ملو (AL 2.1).



The number 2000 **ي** لاكو (AL 5.9).



The number 10000 represented using the word צעאכע βrywr (AL 2.1).

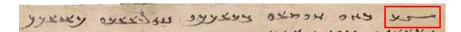
Figure 24: Further examples of numbers in the 'Ancient Letters'. See also figures 22 and 23.



The heterogram 'D written as yso <so alternate ayın, y resh-ayın-daleth> (AL 1.1).



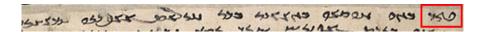
The heterogram 'D written as (y) so < so alternate ayın, (y resh-ayın-daleth) > (AL 1 verso).



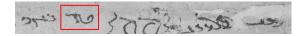
The heterogram 'D written as אבי < אווא, א RESH-AYIN-DALETH> (AL 2.1).



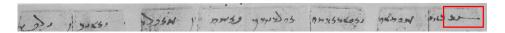
The heterogram 'D written as ys < s AYIN, y RESH-AYIN-DALETH> (AL 2 verso).



The heterogram 'D written as yso <so alternate ayın, y resh-ayın-daleth> (AL 3.1).



The heterogram 'D written as the ligature TO HETEROGRAM AYIN-DALETH (AL 3 verso).



The heterogram 'D written as YS < AYIN, Y RESH-AYIN-DALETH> (AL 4.1).



The heterogram 'D written as yso <50 alternate ayın, y resh-ayın-daleth> using the glyphic variant of alternate ayın (AL 5.1).

Figure 25: Specimens of the heterogram 'D.



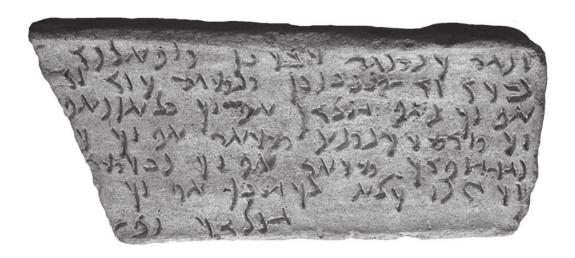


Figure 26: Two images of Kultobe inscription 4 (KII 26859/1). Top from Sims-Williams 2007; bottom from Grenet, et al 2007.



Figure 27: Kultobe inscriptions 2, 1, 3, 5, 10 (from Grenet, et al 2007).

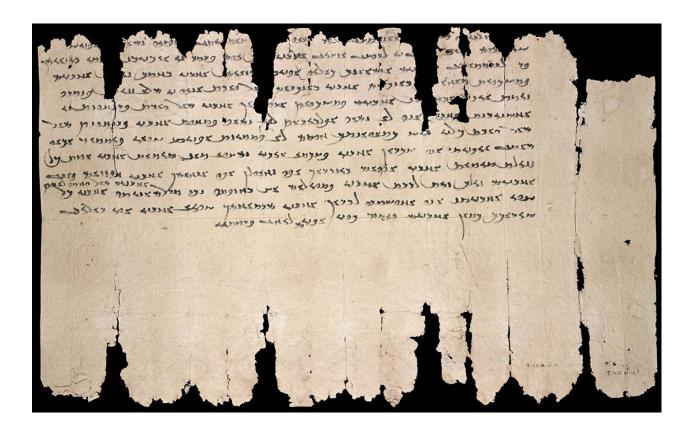


Figure 28: The 'Ancient Letter 1' (British Library, International Dunhuang Project: Or. 8212/92.1 recto 1). "From her daughter, the free-woman Miwnay, to her d[ear] mother [Chatis]." (translation by Sims-Williams in Waugh 2004).

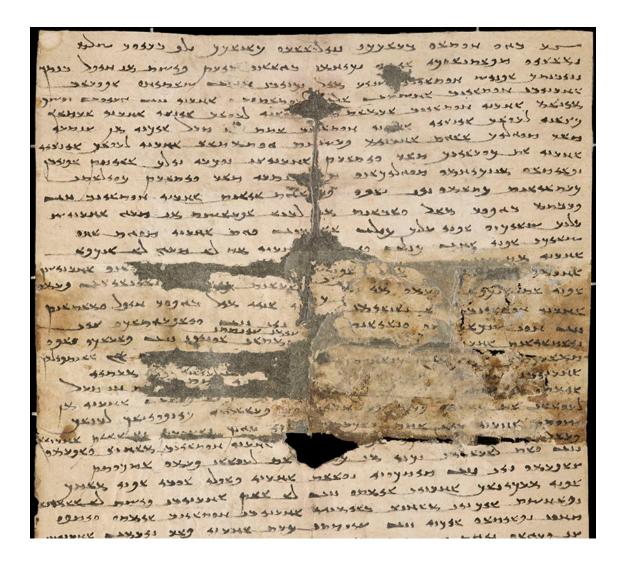


Figure 29: Top portion of 'Ancient Letter 2' (British Library, International Dunhuang Project: Or. 8212/95 side a). "To the noble lord Varzakk (son of) Nanai-thvar (of the family) Kanakk. Sent [by] his servant Nanai-vandak." (translation by Sims-Williams in Waugh 2004). Continued in figure 30.

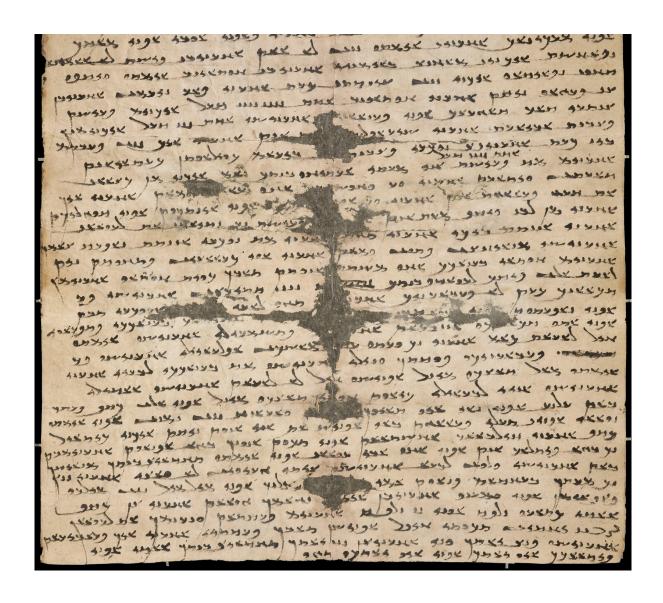


Figure 30: Bottom portion of 'Ancient Letter 2' (British Library, International Dunhuang Project: Or. 8212/95 side a). Continued from figure 29.

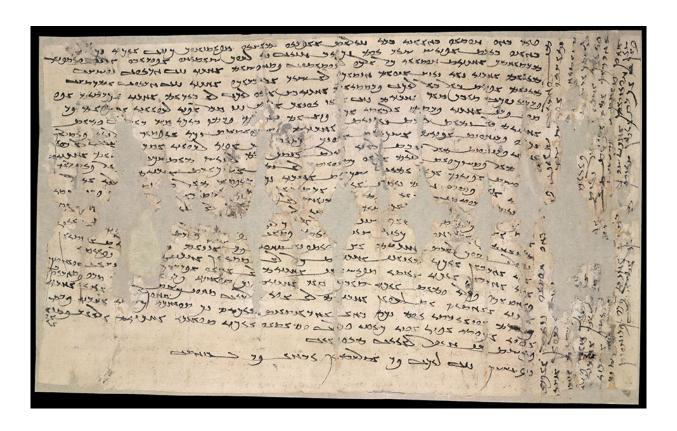


Figure 31: The 'Ancient Letter 3' (British Library, International Dunhuang Project: Or. 8212/98 recto 1). "To (my) noble lord (and) husband Nanai-dhat." (translation by Sims-Williams in Waugh 2004).

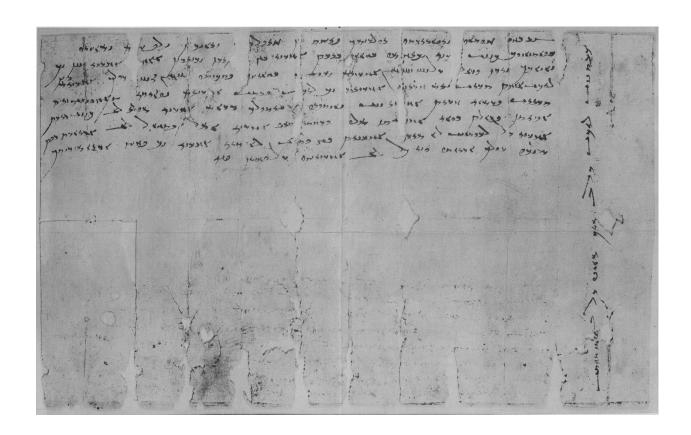


Figure 32: The 'Ancient Letter 4' (British Library: Or. 8212/93 recto; reproduced in Reichelt 1928: plate IV).



Figure 33: The 'Ancient Letter 5' (from Grenet, et al. 1998: 94). "To the noble lord, the chief merchant Aspandhāt. [Sent] by your servant [Frī-khwatāw]."



Figure 34: The 'Ancient Letter 6' (British Library, International Dunhuang Project: Or. 8212/97).

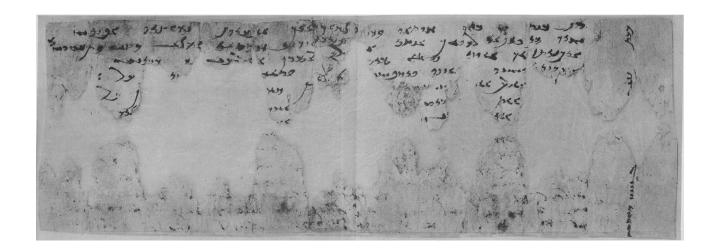


Figure 35: The 'Ancient Letter 7' (British Library: Or. 8212/96 recto; reproduced in Reichelt 1928: plate VII).



Figure 36: Sogdian rock inscription from Shatial (from Sims-Williams 1989: plate 10b) The inscription reads או איי אוווא אוויא אוווא אוויא אווויא אווויא אווויא אווויא אווויא אווויא אווויא אוויא אוויא אווויא אוויא או



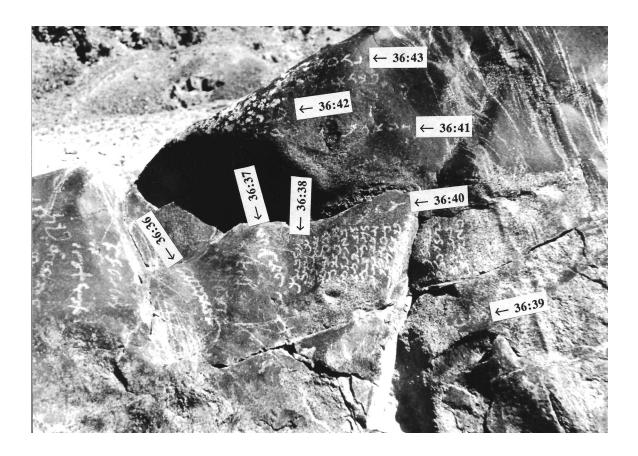
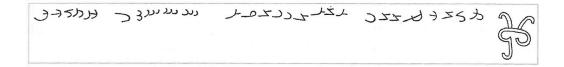


Figure 38: Rock at Shatial containing horizontal and vertical inscriptions in the Old Sogdian script (from Sims-Williams 1989: plate 109b). The text of 36:38 is shown in section 4.5.



Figure 39: Silver coin from Chach bearing an Old Sogdian inscription, 3rd—4th c. CE (reproduced in Grenet 2007: 1023). Reverse: profile of human head. Obverse: tamgha in the center with the text מצאבענצפע כוכואפע c'c'nn'pc wnwnxwr.



## Sogdian script

In the Sogdian script used in the "Ancient Letters" (TABLE 48.2), most of the letters are distinct and do not change shape when joined. In the "formal" and "Uyghur" Sogdian scripts, most of the letters are joined and, owing to the use of a broad pen, are frequently difficult to distinguish. In the earlier form, 'is still distinguished from n; but in the later,  $\dot{} = n$ ,  $\dot{} = n$ . Some scribes distinguish z from n by not connecting z to the preceding letter, but others make no distinction. In the later, increasingly cursive, form, other letters tend to become indistinguishable as well:  $\gamma/x/s/\tilde{s}$ ,  $r/\beta/y$ . Some letters are distinguished only in final position (by some scribes), e.g.,  $n \sim z$ ,  $x \sim \gamma$ .

z is sometimes distinguished from n or z from  $\check{z}$  by a diacritical point  $\underline{\cdot}$ , and the foreign sound b was noted as  $\underline{\dot{\circ}}$   $\dot{\mathbf{p}}$ .

### SAMPLES OF SOGDIAN

#### ANCIENT LETTERS

	_		wr <sup>&gt;&gt;</sup> βδynn				DO←
oιsse wn <sup>33</sup> γ		رويدوو XyKZ	کبور <b>پ</b> YZKYA	و <b>ده</b> بعدوره ykwn <sup>°</sup> zt <sup>°</sup> J			אלנאיני rwyrβ
				תיים וגא nn ktnβ	_	_	

1. Translitera	tion: OD	βγw	$xwt^{3}w$	βr³kk	nny	δβ <sup>&gt;&gt;</sup> rw	k <sup>3</sup> n <sup>3</sup> kk
2. Normalization:		βαγυ	xutāw	βarak	nan	ıē-θβār	kanak
3. Gloss:	to	lord.ACC	master	Barak	k Na	na's-gift	Kanak
ı. ıLP	βrywr	ŠLM	nm³c	yw	sp <sup>3</sup> 1	z'nwky	AYKZY
2. (ēw-)zār	βrēwar	*āfrīwa	n namā	icyu	spā	tzānūk	kaδ-uti
3. thousand ten.thousand greeting(?) reverence.ACC bended.knee when-that.and							
ı. ZKyXMw	βγ"νης	βyrt	pyšt	MN	хурθ	βntk	nnyβntk
2. wēšanu	βaγān(u)	βyart	pišt	con	xēpθ	βantē	nanē-βantē
3. them.obl	lords.0BL	received	written	from	own	servant	Nana's-servant

'To the Divine Master Barak(?) Nanethvar Kanak a thousand, ten thousand greetings, reverently with bended knees when received by their divinities. Written by his own servant Nanevante.'

-From the Old Sogdian "Ancient Letters" found in a mailbag in the Great Wall (AL II, Reichelt 1931: 12 and pl. 2).

Figure 41: Description of the Sogdian script of the 'Ancient Letters' (from Skjærvø 1996: 529).

TABLE 48.2: Main East Iranian Scripts Developed from Aramaic

Aramaic	Sogdian Ancient Letters	Sogdian sutra script	Manichean Sogdian	Christian Sogdian	Principal Phonetic Values (Sogdian)
>	\$	۵, 4	N	س 2	a, ā
b	<b>5</b>	۵, ٥	<u> </u>	-	b, ß
(β)			<u>ت</u>		β
g	**	*	4	1	g, y
(γ)			ž	_	γ
d	y		٠،٢	•	$d,\delta$
h ( <u>h</u> )		E	×	<b>a</b>	a, Ø
w	•	<b>4</b> , •	•	•	$w,\breve{\bar{o}},\breve{\bar{u}}$
z		J	<	•	z
(j)			7		ž
(ž)		٠.	Ë	٧	ž
ḥ (h)	<b>H</b> u	۱۵, ۵	અ	**	γ, x, h
ţ			e	Ŋ	t
у	4	۵, ۰	•	J	y, ĕ, ĭ
k	ל	و,ما	_	•	k
(x)			ف	ડ	x
l (δ)	>	1,0	22	9 7	δ
m	٧,	<b>\$</b> , <b>\$</b>	z z	<b>79 29</b>	m
n		L, .	۲.	<b>\</b> 1	n
S	وو	», »	<u> </u>	•	s
C	5	<b>©</b>	_	_	Ø
p	•	ی	_	S	p
(f)		•	خـ	4	f
ș (c)	سو	•	ىبى	S	č, j
q			ह्य ह्य	<b>45</b>	k
r	,	<b>4</b> , 3	રં નં	j	r
š	مهو	», »	ယ	<b>y x</b>	š
t	در د	6, 6	V	1	t, θ

Figure 42: Table showing various scripts for writing Sogdian (from Skjærvø 1996: 519).

# 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 (P & P) from <a href="http://std.dkuug.dk/JTC1/SC2/WG2/docs/principles.html">http://std.dkuug.dk/JTC1/SC2/WG2/docs/principles.html</a> for guidelines and details before filling this form.

Please ensure you are using the latest Form from <a href="http://std.dkuug.dk/JTC1/SC2/WG2/docs/summaryform.html">http://std.dkuug.dk/JTC1/SC2/WG2/docs/summaryform.html</a>.

Please ensure you are using the latest Form from <a href="http://std.dkuug.dk/JTC1/SC2/WG2/docs/summaryform.html">http://std.dkuug.dk/JTC1/SC2/WG2/docs/summaryform.html</a>. See also <a href="http://std.dkuug.dk/JTC1/SC2/WG2/docs/roadmaps.html">http://std.dkuug.dk/JTC1/SC2/WG2/docs/roadmaps.html</a> for latest Roadmaps.

### A. Administrative

7417444114414							
		Old Sogdian script in Unico	ode				
2. Requester's name:	Anshuman Pa	ndey <pandey @umich.edu=""></pandey>					
3. Requester type (Member body/Liaison/							
4. Submission date:	2016-12-3	31					
5. Requester's reference (if applicable):							
<ol><li>Choose one of the following:</li></ol>							
This is a complete proposal:			Yes				
(or) More information will be prov	ided later:						
B. Technical – General							
<ol> <li>Choose one of the following:</li> </ol>							
<ul> <li>a. This proposal is for a new script (</li> </ul>	set of characters):		Yes				
Proposed name of script:		Old Sogdian					
<ul> <li>b. The proposal is for addition of ch</li> </ul>	aracter(s) to an existin	g block:					
Name of the existing block:							
2. Number of characters in proposal:			40				
3. Proposed category (select one from be							
A-Contemporary B.1-Specializ	ed (small collection)	B.2-Specialized (large c	ollection)				
C-Major extinct X D-Attested ex	ktinct	E-Minor extinct					
F-Archaic Hieroglyphic or Ideographic		G-Obscure or questionable usa	ge symbols				
4. Is a repertoire including character nam	es provided?		Yes				
a. If YES, are the names in accorda	nce with the "characte	r naming guidelines"					
in Annex L of P&P document	?		Yes				
<ul> <li>b. Are the character shapes attached</li> </ul>	ed in a legible form suit	able for review?	Yes				
5. Fonts related:							
a. Who will provide the appropriate	computerized font to th	e Project Editor of 10646 for pub	olishing the				
standard?		,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	<b>3</b> · ·				
	Anshuman F						
b. Identify the party granting a license for use of the font by the editors (include address, e-mail, ftp-site, etc.):							
	Anshuman F	Pandey					
6. References:							
<ul> <li>a. Are references (to other characte</li> </ul>			Yes				
<ul><li>b. Are published examples of use (s</li></ul>	such as samples from i	newspapers, magazines, or other	r sources)				
of proposed characters attached?		Yes					
7. Special encoding issues:							
Does the proposal address other as							
presentation, sorting, searching, inc	lexing, transliteration e	tc. (if yes please enclose informa	ation)? <u>Yes</u>				
8. Additional Information:							
Submitters are invited to provide any add	itional information abou	ut Properties of the proposed Cha	aracter(s) or Script				
that will assist in correct understanding of	and correct linguistic p	processing of the proposed chara	acter(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							
related information. See the Unicode sta							
see Unicode Character Database ( http://							
for information needed for consideration h	by the Unicode Technic	at Committee for inclusion in the	Unicode Standard				

 $<sup>^1\,\</sup>text{Form number: N4502-F (Original 1994-10-14; Revised 1995-01, 1995-04, 1996-04, 1996-08, 1999-03, 2001-05, 2001-09, 2003-11, 2005-01, 2005-09, 2005-10, 2007-03, 2008-05, 2009-11, 2011-03, 2012-01)}$ 

## C. Technical - Justification

Has this proposal for addition of character(s) been submitted before?							
If YES explain							
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?  Nicholas Sims-Williams <ns5@soas.ac.uk> Yutaka Yoshida <yutaka.yoshida@bun.kyoto-u.ac.jp></yutaka.yoshida@bun.kyoto-u.ac.jp></ns5@soas.ac.uk>							
If YES, available relevant documents:							
3. Information on the user community for the proposed characters (for example:							
size, demographics, information technology use, or publishing use) is included?							
Reference: See text of proposal							
	Common						
Reference: See text of proposal							
5. Are the proposed characters in current use by the user community?	Yes;						
If YES, where? Reference: Currently used by scholars of Sogdian and Central Asian s	tudies						
6. After giving due considerations to the principles in the P&P document must the proposed characters be	e entirely						
in the BMP?	N/A						
If YES, is a rationale provided?							
If YES, reference:							
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?	No						
If YES, is a rationale for its inclusion provided?							
If YES, reference:							
9. Can any of the proposed characters be encoded using a composed character sequence of either							
existing characters or other proposed characters?	No						
If YES, is a rationale for its inclusion provided?							
If YES, reference:							
10. Can any of the proposed character(s) be considered to be similar (in appearance or function)	Ma						
to, or could be confused with, an existing character?	No						
If YES, is a rationale for its inclusion provided?							
If YES, reference:							
11. Does the proposal include use of combining characters and/or use of composite sequences?	No						
If YES, is a rationale for such use provided?							
If YES, reference:	N//						
Is a list of composite sequences and their corresponding glyph images (graphic symbols) provided?	N/A						
If YES, reference:							
12. Does the proposal contain characters with any special properties such as control function or similar semantics?	No						
If YES, describe in detail (include attachment if necessary)	740						
II 1 L3, describe in detail (include attachment il necessary)							
13. Does the proposal contain any Ideographic compatibility characters?							
If YES, are the equivalent corresponding unified ideographic characters identified?							
If YES, reference:							