# Preliminary proposal to encode the Old Sogdian script in Unicode

Anshuman Pandey
Department of Linguistics
University of California, Berkeley
Berkeley, California, U.S.A.
anshuman.pandey@berkeley.edu

November 3, 2015

## 1   Introduction

This is a preliminary proposal for a Unicode encoding for the scripts used in the Sogdian 'Ancient Letters' and in Sogdian inscriptions of the Upper Indus. These scripts are tentatively designated as 'Old Sogdian' (see section 4.1). The intent is to enable users to input, display, and exchange content in the Old Sogdian script on digital devices.

This document describes an encoding model and character repertoire, and provides specimens of sources in which the script appears. In the proposed repertoire each graphical element of the script that possesses a distinctive semantic value is defined as a Unicode 'character'. The exception are letters with identical glyphs, which are merged into a single character in accordance with the character-glyph model of the Unicode standard. This character-encoding repertoire may differ from traditional scholarly analyses of the script. Such differences naturally arise from the requirements for representing the script electronically as opposed to writing it by hand. The representative glyphs shown here are illustrative and are not intended to be typographically aesthetic.

The author seeks feedback from scholars regarding the proposed encoding model, character repertoire, and representative glyphs. Issues requiring further discussion are enumerated in section 7. A formal proposal to encode the script is forthcoming, pending comments from experts.

## 2   Background

The Old Sogdian script was used primarily for writing Sogdian (ISO 639: sog), an Eastern Iranian language that was spoken in the region between the Amu Darya and the Syr Darya, which roughly corresponds to the modern provinces of Samarqand and Bokhara in Uzbekistan and the Sughd province of Tajikistan. The language and script were carried by Sogdian merchants from eastern Iran to western China.

The earliest manuscripts containing the Old Sogdian script are known as the 'Ancient Letters' (see figures 1–5). These paper documents were found in 1907 by Aurel Stein in Dunhuang, now in Gansu province, western China. Based upon evidence contained within the documents, it has been suggested that the 'Ancient Letters' were written in 312–313 CE (Sims-Williams 1985). A system of writing similar to that used in the 'Ancient

Letters' appears in some 600 rock carvings at Shatial in the Gilgit region of Pakistan. These 'Upper Indus inscriptions' have been dated to the 4th–7th centuries CE (Sims-Williams 1989, 2000; see figures 6 and 7 here). Based upon the structural and graphical similarities between the scripts of the 'Ancient Letters' and the 'Upper Indus inscriptions', they may be considered to be typologically identical.

The script bears resemblance to Aramaic, Inscriptional Parthian, and Inscriptional Pahlavi. A comparison of these scripts is provided in table 5.

The Old Sogdian script is the focus of active scholarship. The International Dunhuang Project at the British Library is currently digitizing texts in Old Sogdian, namely the 'Ancient Letters'. The ability to represent the contents of these documents in plain text for scholarly exchange, cataloguing, and the production of corpora requires an encoding for Old Sogdian in Unicode.

## 3   Roadmap for Sogdian scripts

More than a decade ago, in an overview on the encoding of additional Semitic scripts it was suggested that the Aramaic block would be expected to encompass 'Sogdian' (Everson 2001). The document did not define 'Sogdian'. It may be inferred that the reference was most likely to the 'formal' Sogdian script, which is the most well known. There are, however, several varieties of typologically related scripts that are known by designation 'Sogdian':

1. the script of the oldest Sogdian inscriptions at Kultobe
2. the script of the 'Ancient Letters' from Dunhuang
3. the script of the Upper Indus inscriptions at Shatial
4. the 'cursive' script of documents from Mount Mug
5. the formalized 'cursive' script used in Buddhist texts, known as the 'formal' or 'sūtra' script
6. the latest cursive form, known as the 'Sogdian Uyghur' or 'Turco-Sogdian' script

In relation to the above 'proper Sogdian' scripts, mention may be made here to a derivative of 'Sogdian Uyghur' known as 'Old Uyghur', which is written vertically. The Sogdian language was also written in Syriac, Manichaean, and Brahmi-based scripts. All of these scripts are distinctive writing systems in their own right and some have already been encoded in Unicode.

None of the scripts belonging to the 'proper Sogdian' family are encoded in Unicode. It is, therefore, beneficial to describe a possible roadmap for handling the encoding of these scripts. The available information suggests that these scripts may be divided across three encodings:

• 'Old Sogdian': to encompass scripts #2 and #3, and possibly #1
• 'Cursive Sogdian': for script #4
• 'Sogdian': a unified block for scripts #5 and #6

These divisions are based upon the orthographic structure, letterforms, and character repertoires of the scripts. Scripts #1–#3 are structurally non-conjoining, while #4–#6 are inherently conjoining. The letters of #1–#3 retain their basic shapes in words, while letters of #4–#6 have contextual shapes based upon their position in a word. The letterforms of these six scripts may be classified into two groups based upon similarity of graphical structures: letters of #1–#4 exhibit common features, but differ from letter shapes of #5–#6. Furthermore, scripts #5–#6 use diacritics to indicate different phonetic values of a particular letter. Although these scripts are typologically related, their differing structures require different technical implementations, and therefore separate encodings. A comparison of scripts #2 and #4–#6 is shown in figure 10.

## 4 Script Details

### 4.1 Script name

The name 'Old Sogdian' is assigned for the script in Unicode in order to distinguish it from related scripts, but also to emphasize their genetic relationship. There is no standard or conventional name for the script of the 'Ancient Letters' or the 'Upper Indus inscriptions'. The catalogue of the International Dunhuang Project at the British Library refers to the script in the 'Ancient Letters' generically as "Sogdian" and does not differentiate it from other related scripts, which are labelled using the same name. Skjærvø (1996) refers to the script as "Sogdian Aramaic".

### 4.2 Structure

Old Sogdian is an *abjad* that is written from right to left. Structurally, it is a non-joining *abjad*, similar to Hebrew. Letters retain their shapes within a word; however, some letters have distinctive word-final shapes (see section 5.1.1). The available sources show instances where adjacent letters are connected by joined or overlapping strokes, but such conjunctions result from the regular flow of writing rather than any intrinsic conjoining behavior of the letters, as is the case for Arabic, Mongolian, and the later 'cursive, 'sūtra' and 'Uyghur' Sogdian scripts. Similar to other *abjad* writing systems, vowels in Old Sogdian are represented using *aleph*, *yodh*, and *waw*; however, these are also used for short vowels, not only long forms.

### 4.3 Encoding model

The arrangement of core letters in the code block proposed for Old Sogdian is based upon a one-to-one correspondence with the Aramaic block. However, Old Sogdian analogues do not exist for every Aramaic letter. Only 20 Old Sogdian letters corresponding to the original 22 Aramaic are attested in the available sources. The repertoire lacks analogues for Aramaic *teth* and *qoph*. The number of distinctive letters is further reduced on account of the merging of glyphs for two or more letters into forms that are nearly identical or indistinguishable. Such is the case for *zayin* and *nun*; as well as for *ayin* and *resh*, and quite possibly *daleth*, which could be assimilated with these two. The result is a repertoire that contains 17 or 18 distinctive letters.

The encoding model for Old Sogdian differs from the Unicode models for other Iranian scripts derived from Aramaic. For the Inscriptional Parthian and Inscriptional Pahlavi encodings letters with nearly identical glyphic representations are merged into single letters. For instance, in Inscriptional Pahlavi the letters *waw*, *ayin*, and *resh* are represented using ✷ U+10B65 INSCRIPTIONAL PAHLAVI LETTER WAW-AYIN-RESH. A similar case occurs with Inscriptional Pahlavi *mem* and *qoph*, which are merged into the character ✷ U+10B6C INSCRIPTIONAL PAHLAVI LETTER MEM-QOPH.

The proposed arrangement for the Old Sogdian block departs from the pattern for Inscriptional Parthian and Inscriptional Pahlavi. These latter script blocks are not arranged in a one-to-one correspondence with Aramaic and do not have spaces reserved within them for potential disunification of merged letters. For Old Sogdian, spaces have been reserved for *teth* and *qoph* in the event that these characters are discovered in the future. The *nun* may be perceived to be identical to *zayin*, and therefore it is not encoded as separate character. But a space is reserved for it in the case that a distinctive form is identified later. Until that time, *zayin* is to be also used for *nun*. Similarly, there is no distinctive character for *ayin*, but a space is reserved for it. The letter *resh* is to be also used for *ayin*. This approach is preferred over merging these letters into single characters such as *OLD SOGDIAN LETTER ZAYIN-NUN and *OLD SOGDIAN LETTER AYIN-RESH.

The rationale for aligning the Old Sogdian code block with the full Aramaic repertoire is based upon a source from the late 7th century or early 8th century CE. An ostracon found at Panjikant, Tajikistan bears an

inscription that enumerates twenty-three letters of the 'cursive' Sogdian script; the last letter is a repetition of *lamedh* (Livshits 2015: 227). This source suggests that knowledge and memory of the twenty-two letter repertoire persisted among Sogdian scribes well after the period of the 'Ancient Letters' and the 'Upper Indus inscriptions'. Even if the repertoire etched upon the potsherd contains vestigial letters, such as *qoph*, a character-encoding standard for that script should provide a means for full representation. If a complete alphabet exists for 'cursive' Sogdian, then it is quite possible that something similar exists for Old Sogdian. For that reason it is practical to reserve spaces in the proposed code block for letters that appear to have disappeared in case that they are identified in future discoveries.

Word-final forms of several letters occur in the available sources. A complete list of distinctive word-final forms has yet to be produced.

### 4.4 Representative glyphs

Old Sogdian letters exhibit a high degree of glyphic uniformity across the available sources, as observed in the specimens of the 'Ancient Letters' and 'Upper Indus inscriptions' consulted by the proposal author. The representative glyphs for Old Sogdian characters are normalized forms of characters found in the 'Ancient Letters'. The glyphs for letters are derived from Skjærvø (1996), shown here in figure 9, and have been modified by the proposal author. Glyphs for numbers have been drawn by the proposal author.

## 5 Proposed characters

The proposed repertoire for Old Sogdian contains 30 characters: 18 basic letters, 2 alternate letters, 9 numbers, and 1 fraction. Word-final forms of letters are not included at present. Names for basic letters align with corresponding characters of the 'Imperial Aramaic' block.

### 5.1 Letters

The following 18 basic letters are included in the proposed repertoire:

| | Character name | Phonetic values |
|---|---|---|
| ⵏ | OLD SOGDIAN LETTER ALEPH | a, ā |
| �misc | OLD SOGDIAN LETTER BETH | β |
| ⵏ | OLD SOGDIAN LETTER GIMEL | γ |
| ⵏ | OLD SOGDIAN LETTER DALETH | d |
| ⵏ | OLD SOGDIAN LETTER HE | Ø, a |
| ⵏ | OLD SOGDIAN LETTER WAW | w, u, ū, o, ō |
| ⵏ | OLD SOGDIAN LETTER ZAYIN | z, ž ; n |
| ⵏ | OLD SOGDIAN LETTER HETH | x |
| ⵏ | OLD SOGDIAN LETTER YODH | y, i, ī, e, ē |

| ﻜ | OLD SOGDIAN LETTER KAPH | k, g |
| | | |
| ﻝ | OLD SOGDIAN LETTER LAMEDH | δ, l |
| ﻡ | OLD SOGDIAN LETTER MEM | m |
| ﻥ | OLD SOGDIAN LETTER SAMEKH | s |
| ﻭ | OLD SOGDIAN LETTER PE | p, b, f |
| ﺻ | OLD SOGDIAN LETTER SADHE | č, ǰ |
| ﻱ | OLD SOGDIAN LETTER RESH | r ; Ø, a |
| ﺵ | OLD SOGDIAN LETTER SHIN | š |
| ﺕ | OLD SOGDIAN LETTER TAW | t, d |

The following 2 alternate letters are included:

| | Character name | Phonetic value |
| --- | --- | --- |
| ﺡ | OLD SOGDIAN LETTER ALTERNATE HE | a |
| ﻉ | OLD SOGDIAN LETTER ALTERNATE AYIN | Ø |

**Notes on the repertoire**

1. The letter ﻯ DALETH is used only in Aramaic heterograms.

2. The letter ﻩ HE is used only at the end of words. It has a variant form ﺡ that occurs concurrently in various documents, which is proposed for encoding as the separate character ﺡ ALTERNATE HE in order to enable the representation of both.

3. The letters *zayin* and *nun* have nearly identical, if not indistinguishable shapes. For this reason, a *nun* is not proposed for encoding as a separate letter. Instead, the letter ZAYIN is to be used for both *zayin* and *nun*.

4. The letters *ayin* and *resh* also have indistinguishable forms. A separate character for *ayin* is not proposed. Instead, the letter RESH is to be used for both *ayin* and *resh*.

5. In the 'Ancient Letters', the letter ﻉ ALTERNATE LETTER AYIN occurs only in Aramaic heterograms, such as ﻯﻉ *OD* <ﻉ ALTERNATE LETTER AYIN, ﻯ DALETH>, a preposition meaning "to" that is used in terms of address in correspondence.

6. Some scholars proposed that *qoph* was retained and reassigned for the number 100 (Sims-Williams 1985); however, the common position is that *qoph* was never used in the Old Sogdian script.

### 5.1.1 Word-final shapes

In the 'Ancient Letters' certain letters are written using glyphically distinctive forms when they occur in word-final position. These final forms are distinguished from regular forms by the shape of the terminal stroke. A curved stroke may be modified into a vertical line or the curve may be rounded upwards; a horizontal stroke may be elongated.

For some letters the shape of the terminal stroke produces a discernable distinction between the regular and final forms. For instance, ﻼ SAMEKH may be written as ﻪ in final position; ﻟ *nun* (ZAYIN) may be written as ﺍ in final position (*zayin* retains the form ﻟ when final). Final forms of other letters are not as distinctive. For example, ﻼ ALEPH may occur as ﻚ; ﺥ TAW may occur as ﺦ; For such letters, the terminal stroke may be interpreted as a stylistic flourish rather than a specific 'final' form. Moreover, there is no regularity in usage. For instance, regular and 'final' forms of SAMEKH are observed in final position in adjacent words.

At present, it is unclear if the differences between regular and word-final forms of letters should be distinguished at the character level, as is the case for the Hebrew and Nabatean encodings in Unicode.

### 5.2 Numbers

Distinctive signs are attested for 'one', 'ten', 'one hundred', and the fraction 'one-half'. Other numbers are produced using repetitions and compounds of these signs, and using letters and words. The following characters are proposed for representing numbers:

| | Character name | Numeric value |
|---|---|---|
| ﻟ | OLD SOGDIAN NUMBER ONE | 1 |
| ﻠ | OLD SOGDIAN NUMBER TWO | 2 |
| ﻡ | OLD SOGDIAN NUMBER THREE | 3 |
| ﻢ | OLD SOGDIAN NUMBER FOUR | 4 |
| ﺝ | OLD SOGDIAN NUMBER TEN | 10 |
| ﺟ | OLD SOGDIAN NUMBER TWENTY | 20 |
| ﺠ | OLD SOGDIAN NUMBER THIRTY | 30 |
| ﺤ | OLD SOGDIAN NUMBER ONE HUNDRED | 100 |
| ﻮ | OLD SOGDIAN NUMBER ONE THOUSAND | 1000 |

The ordering of numbers follows the right-to-left directionality of the script. The expression of numbers is additive. Compound numbers of different decimal orders are written by placing the larger numbers first.
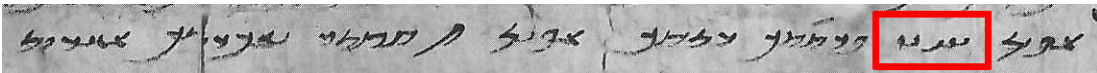
### 5.2.1 Primary units

The numbers 2–9 are produced using repetitions of the number ﻟ ONE. These sequences are generally joined together at the left edge of ONE, eg. '3' ﻡ. The numbers 4–9 are written using sequences of ONE arranged in

groups of three or four separated by a space, eg. '6' سس; '7' سس سس . The representation of the primary units is modeled upon the encoding for Inscriptional Parthian, which has separate characters for the numbers 1–4, eg. ل U+10B58 ɪɴꜱᴄʀɪᴘᴛɪᴏɴᴀʟ ᴘᴀʀᴛʜɪᴀɴ ɴᴜᴍʙᴇʀ ᴏɴᴇ ... ⦀⦀ U+10B5B ɪɴꜱᴄʀɪᴘᴛɪᴏɴᴀʟ ᴘᴀʀᴛʜɪᴀɴ ɴᴜᴍʙᴇʀ ꜰᴏᴜʀ. Similarly, the following characters are proposed for Old Sogdian: ر ᴏɴᴇ, ru ᴛᴡᴏ, سس ᴛʜʀᴇᴇ, سس ꜰᴏᴜʀ.

The number 4 written سس (from 'Ancient Letter 5', line 26).



The number 5 written سن (from 'Ancient Letter 5', line 10).



The number 8 written سس سس (from 'Ancient Letter 2', line 26).



### 5.2.2    Tens

The number ⌣ ᴛᴇɴ resembles a vertically compressed form of the letter ⌣ ʟᴀᴍᴇᴅʜ. Multiples of 10 are produced using horizontal sequences or vertical stacks of ⌣. The only attested stacked forms are 20 ⅔ and 30 ⅔. These are encoded as the atomic characters ⅔ ᴛᴡᴇɴᴛʏ and ⅔ ᴛʜɪʀᴛʏ.

The number 20 ⅔ (from 'Ancient Letter 5', line 21):



The number 30 has two different representations in the available sources. In one instance it occurs as ⅔ (from 'Ancient Letter 5', line 32):



In another source, 30 is represented as ⌣⅔, which is a compound of ⅔ ᴛᴡᴇɴᴛʏ and ⌣ ᴛᴇɴ. This form occurs in the number ⅔دن '32' (from 'Ancient Letter 2', line 62):



The proposed repertoire offers a means for representing both forms of 30.

### 5.2.3   Hundreds

The number 100 is written using ꣽ ONE HUNDRED.  The hundreds unit is represented using primary numbers followed by ꣽ ONE HUNDRED, eg. 500 is expressed as ꣽدد سد <سد THREE, دد TWO, ꣽ ONE HUNDRED>.

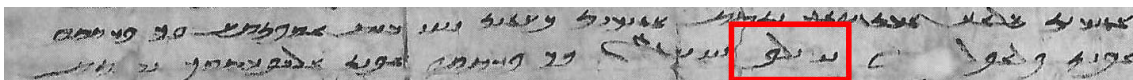The number 500 written ꣽدد سد (from 'Ancient Letter 5', line 9).



### 5.2.4   Thousands

The number 1000 is written using ولد ONE THOUSAND.  This numerical sign is a ligature consisting of the number د ONE joined to the sequence ود لـ <لـ LAMEDH, و PE>.  The word ود represents the Aramaic hetero-gram *LP* "thousand" (= ﻝ U+1085A IMPERIAL ARAMAIC NUMBER ONE THOUSAND).  Although ولد can be represented using the sequence <د ONE, لـ LAMEDH, و PE>, it is encoded as an atomic character on account of its conceptualization as a distinctive unit and the numerical value it possesses.

The thousands are represented using the appropriate repetitions of د ONE followed by ولد ONE THOUSAND, eg. ولد دد '2,000'.  The presence of د ONE in the glyph for ONE THOUSAND does not have a separate value.

The number 1,000 written ولد (from 'Ancient Letter 2', line 1):



The number 2,000 written ولد دد (from 'Ancient Letter 5', line 9):



### 5.2.5   Ten thousand

There is no distinctive sign for the number 10,000.  It is represented using the word ʸⱻϬ𐼿ϓ *βryw'r* (from 'Ancient Letter 2", line 1):
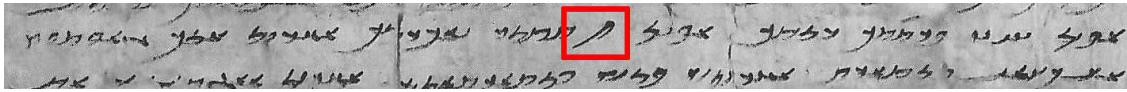


### 5.2.6   Fractions

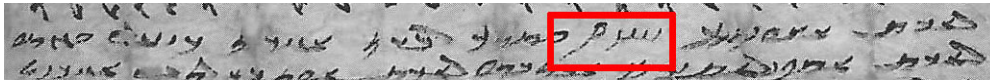The following fraction is attested:

| Character name | Numeric value |
|---|---|
| ↗   OLD SOGDIAN FRACTION ONE HALF | ½ |

The ⟋ FRACTION ONE HALF occurs in 'Ancient Letter 5'. This character was identified as an alternate form of 100 by early scholars of Sogdian, but Grenet, et al (1998) have interpreted it as a sign for the fraction ½.

The fraction ½ written ⟋ (from 'Ancient Letter 5', line 10):



The fraction 4½ written ⟋ שש (from 'Ancient Letter 5', line 24):



### 5.3  Punctuation

There are no script-specific marks of punctuation. The available records show the usage of spaces for separating words. There are no punctuation signs for indicating the end of a sentence or larger sections of text.

### 5.4  Linebreaking

Linebreaks appear to occur after the end of a word. Words are not broken across lines and there are no continuation marks, such a hyphen. A sequence of numbers should be treated similarly, without being broken across lines.

### 5.5  Collation

The sort order for Old Sogdian is as follows:

⟋ ALEPH  <  ⟋ BETH  <  ⟋ GIMEL  <  ⟋ DALETH  <  ⟋ HE  ≪  ⟋ ALTERNATE HE  <

⟋ WAW  <  ⟋ ZAYIN  <  ⟋ HETH  <  ⟋ YODH  <  ⟋ KAPH  <  ⟋ LAMEDH  <  ⟋ MEM  <

⟋ SAMEKH  <  ⟋ ALTERNATE AYIN  <  ⟋ PE  <  ⟋ SADHE  <  ⟋ RESH  <  ⟋ SHIN  <

⟋ TAW

## 6  Character Data

**Character Properties**   Properties in the format of `UnicodeData.txt`:

```
xx00;OLD SOGDIAN LETTER ALEPH;Lo;0;R;;;;;N;;;;;
xx01;OLD SOGDIAN LETTER BETH;Lo;0;R;;;;;N;;;;;
xx02;OLD SOGDIAN LETTER GIMEL;Lo;0;R;;;;;N;;;;;
xx03;OLD SOGDIAN LETTER DALETH;Lo;0;R;;;;;N;;;;;
xx04;OLD SOGDIAN LETTER HE;Lo;0;R;;;;;N;;;;;
xx05;OLD SOGDIAN LETTER WAW;Lo;0;R;;;;;N;;;;;
xx06;OLD SOGDIAN LETTER ZAYIN;Lo;0;R;;;;;N;;;;;
xx07;OLD SOGDIAN LETTER HETH;Lo;0;R;;;;;N;;;;;
xx08;<reserved>
```

9

```
xx09;OLD SOGDIAN LETTER YODH;Lo;0;R;;;;;N;;;;;
xx0A;OLD SOGDIAN LETTER KAPH;Lo;0;R;;;;;N;;;;;
xx0B;OLD SOGDIAN LETTER LAMEDH;Lo;0;R;;;;;N;;;;;
xx0C;OLD SOGDIAN LETTER MEM;Lo;0;R;;;;;N;;;;;
xx0D;<reserved>
xx0E;OLD SOGDIAN LETTER SAMEKH;Lo;0;R;;;;;N;;;;;
xx0F;<reserved>
xx10;OLD SOGDIAN LETTER PE;Lo;0;R;;;;;N;;;;;
xx11;OLD SOGDIAN LETTER SADHE;Lo;0;R;;;;;N;;;;;
xx12;<reserved>
xx13;OLD SOGDIAN LETTER RESH;Lo;0;R;;;;;N;;;;;
xx14;OLD SOGDIAN LETTER SHIN;Lo;0;R;;;;;N;;;;;
xx15;OLD SOGDIAN LETTER TAW;Lo;0;R;;;;;N;;;;;
xx16;OLD SOGDIAN LETTER ALTERNATE HE;Lo;0;R;;;;;N;;;;;
xx17;OLD SOGDIAN LETTER ALTERNATE AYIN;Lo;0;R;;;;;N;;;;;
xx18;OLD SOGDIAN NUMBER ONE;No;0;R;;;;1;N;;;;;
xx19;OLD SOGDIAN NUMBER TWO;No;0;R;;;;2;N;;;;;
xx1A;OLD SOGDIAN NUMBER THREE;No;0;R;;;;3;N;;;;;
xx1B;OLD SOGDIAN NUMBER FOUR;No;0;R;;;;4;N;;;;;
xx1C;OLD SOGDIAN NUMBER TEN;No;0;R;;;;10;N;;;;;
xx1D;OLD SOGDIAN NUMBER TWENTY;No;0;R;;;;20;N;;;;;
xx1E;OLD SOGDIAN NUMBER THIRTY;No;0;R;;;;30;N;;;;;
xx1F;OLD SOGDIAN NUMBER ONE HUNDRED;No;0;R;;;;100;N;;;;;
xx20;OLD SOGDIAN NUMBER ONE THOUSAND;No;0;R;;;;1000;N;;;;;
xx21;OLD SOGDIAN FRACTION ONE HALF;;No;0;R;;;;1/2;N;;;;;
```

**Linebreaking**    Linebreaking properties in the format of `LineBreak.txt`:

```
xx00..xx07;AL      # Lo     [8] OLD SOGDIAN LETTER ALEPH..OLD SOGDIAN LETTER HETH
xx09..xx0C;AL      # Lo     [4] OLD SOGDIAN LETTER YODH..OLD SOGDIAN LETTER MEM
xx0E;AL            # Lo         OLD SOGDIAN LETTER SAMEKH
xx10..xx11;AL      # Lo     [2] OLD SOGDIAN LETTER PE..OLD SOGDIAN LETTER SADHE
xx13..xx15;AL      # Lo     [3] OLD SOGDIAN LETTER RESH..OLD SOGDIAN LETTER TAW
xx16;AL            # Lo         OLD SOGDIAN LETTER ALTERNATE HE
xx17;AL            # Lo         OLD SOGDIAN LETTER ALTERNATE AYIN
xx18..xx20;AL      # No     [9] OLD SOGDIAN NUMBER ONE..OLD SOGDIAN NUMBER ONE THOUSAND
xx21;AL            # No         OLD SOGDIAN FRACTION ONE HALF
```

## 7   Questions

1. *Glyphic representations*    Are the glyphic representations for the proposed characters acceptable?

2. *Any abecedary?*    Is there any primary source containing the letter repertoire of Old Sogdian, similar to the Panjikant ostracon for 'cursive' Sogdian?

3. *Merged letters*    Is there any evidence to suggest that scribes perceived a distinction between letters such as *zayin* and *nun*, or *ayin* and *resh* when writing Old Sogdian? Did they knowingly use the same glyph for multiple letters?

4. *Missing characters?*    Are there any characters in the 'Ancient Letters' that are not identified here? Alternate forms of other letters?

5. *he*    Are there rules that specify the usage of ⌣ ᴀʟᴛᴇʀɴᴀᴛᴇ ʜᴇ in place of ↘ ʜᴇ? Or, are the two forms interchangable?

6. *teth*    Is an Old Sogdian analogue for Aramaic *teth* attested in heterograms?

10

7. *qoph*    Is an Old Sogdian analogue for Aramaic *qoph* attested in heterograms? Is there any relation between the ꝁ FRACTION ONE HALF and Aramaic ꝗ *qoph*?

8. *Final forms*    Is there is list of letters that have distinctive final forms? Secondly, should final forms be represented using regular letters, or should they be encoded separately, as has been done in Unicode for Hebrew and Nabatean?

## 8    References

Everson, Michael. 2001. "Roadmapping early Semitic scripts" (L2/01-024). `http://www.unicode.org/L2/L2001/01024-n2311.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, edited by Osmund Bopearachchi, Carol Altman Bromberg, and Frantz Grenet, pp. 91–104.

Livshits, V. A. 2015. *Sogdian epigraphy of Central Asia and Semirech'e*. Corpus Inscriptonum Iranicarum, pt. II (Inscriptions of the Seleucid and Parthian periods of Eastern Iran and Central Asia), v. III (Sogdian), no. II. Translated by Tom Stableford, edited by Nicholas Sims-Williams. London: Published on behalf of Corpus Inscriptionum Iranicarum by School of Oriental and African Studies.

———. 2015. "A Sogdian alphabet from Penjikent", in Livshits (2015), pp. 227–232.

Sims-Williams, Nicholas. 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; and Franz 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.

———. 2006. "Iran. VI. Iranian Languages and Scripts. (3) Writing Systems." *Encyclopædia Iranica*, vol. XIII, fasc. 4, pp. 366–370. `http://www.iranicaonline.org/articles/iran-vi3-writing-systems`

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

## 9   Acknowledgments

| | 10E0 | 10E1 | 10E2 |
|---|---|---|---|
| 0 | 10E00 | 10E10 | 10E20 |
| 1 | 10E01 | 10E11 | 10E21 |
| 2 | 10E02 | | |
| 3 | 10E03 | 10E13 | |
| 4 | 10E04 | 10E14 | |
| 5 | 10E05 | 10E15 | |
| 6 | 10E06 | 10E16 | |
| 7 | 10E07 | 10E17 | |
| 8 | | 10E18 | |
| 9 | 10E09 | 10E19 | |
| A | 10E0A | 10E1A | |
| B | 10E0B | 10E1B | |
| C | 10E0C | 10E1C | |
| D | | 10E1D | |
| E | 10E0E | 10E1E | |
| F | | 10E1F | |

*This block unifies the scripts used in the Ancient Letters and Upper Indus inscriptions. Spaces are reserved for Aramaic analogues, which are unattested at present.*

## Letters

| | | |
|---|---|---|
| 10E00 | ⟨glyph⟩ | OLD SOGDIAN LETTER ALEPH |
| 10E01 | ⟨glyph⟩ | OLD SOGDIAN LETTER BETH |
| 10E02 | ⟨glyph⟩ | OLD SOGDIAN LETTER GIMEL |
| 10E03 | ⟨glyph⟩ | OLD SOGDIAN LETTER DALETH |
| | | • used only in Aramaic heterograms |
| 10E04 | ⟨glyph⟩ | OLD SOGDIAN LETTER HE |
| | | • only used in word-final position |
| 10E05 | ⟨glyph⟩ | OLD SOGDIAN LETTER WAW |
| 10E06 | ⟨glyph⟩ | OLD SOGDIAN LETTER ZAYIN |
| | | = zayin, nun |
| 10E07 | ⟨glyph⟩ | OLD SOGDIAN LETTER HETH |
| 10E08 | ⟨reserved⟩ | <reserved> |
| 10E09 | ⟨glyph⟩ | OLD SOGDIAN LETTER YODH |
| 10E0A | ⟨glyph⟩ | OLD SOGDIAN LETTER KAPH |
| 10E0B | ⟨glyph⟩ | OLD SOGDIAN LETTER LAMEDH |
| 10E0C | ⟨glyph⟩ | OLD SOGDIAN LETTER MEM |
| 10E0D | ⟨reserved⟩ | <reserved> |
| 10E0E | ⟨glyph⟩ | OLD SOGDIAN LETTER SAMEKH |
| 10E0F | ⟨reserved⟩ | <reserved> |
| 10E10 | ⟨glyph⟩ | OLD SOGDIAN LETTER PE |
| 10E11 | ⟨glyph⟩ | OLD SOGDIAN LETTER SADHE |
| 10E12 | ⟨reserved⟩ | <reserved> |
| 10E13 | ⟨glyph⟩ | OLD SOGDIAN LETTER RESH |
| | | = daleth, ayin, resh |
| 10E14 | ⟨glyph⟩ | OLD SOGDIAN LETTER SHIN |
| 10E15 | ⟨glyph⟩ | OLD SOGDIAN LETTER TAW |

## Alternate letters

| | | |
|---|---|---|
| 10E16 | ⟨glyph⟩ | OLD SOGDIAN LETTER ALTERNATE HE |
| | | • used only in word-final position |
| 10E17 | ⟨glyph⟩ | OLD SOGDIAN LETTER ALTERNATE AYIN |
| | | • used only in Aramaic heterograms |

## Numbers

| | | |
|---|---|---|
| 10E18 | ⟨glyph⟩ | OLD SOGDIAN NUMBER ONE |
| 10E19 | ⟨glyph⟩ | OLD SOGDIAN NUMBER TWO |
| 10E1A | ⟨glyph⟩ | OLD SOGDIAN NUMBER THREE |
| 10E1B | ⟨glyph⟩ | OLD SOGDIAN NUMBER FOUR |
| 10E1C | ⟨glyph⟩ | OLD SOGDIAN NUMBER TEN |
| 10E1D | ⟨glyph⟩ | OLD SOGDIAN NUMBER TWENTY |
| 10E1E | ⟨glyph⟩ | OLD SOGDIAN NUMBER THIRTY |
| 10E1F | ⟨glyph⟩ | OLD SOGDIAN NUMBER ONE HUNDRED |
| 10E20 | ⟨glyph⟩ | OLD SOGDIAN NUMBER ONE THOUSAND |

## Fraction

| | | |
|---|---|---|
| 10E21 | ⟨glyph⟩ | OLD SOGDIAN FRACTION ONE HALF |

| | Old Sogdian | Inscriptional Pahlavi | Inscriptional Parthian | Imperial Aramaic |
|---|---|---|---|---|
| ALEPH | ⳤ | ᴧ | ᴧ | א |
| BETH | ⳡ | ᒉ | ⳓ | נ |
| GIMEL | Ⳣ | ᒐ | ⳝ | ⳑ |
| DALETH | ⳗ | ⳝ | ⳝ | ר |
| HE | Ⳓ | ⳍ | ⳍ | ꞁ |
| WAW | ⳛ | 2 | ⳛ | ⳗ |
| ZAYIN | ⳛ | ⳛ | ⳛ | ᴵ |
| HETH | ⳨ | ⳛ | ᴧ | ꞁꞁ |
| TETH | — | ⳛ | ⳛ | ⳛ |
| YODH | ⳛ | ⳛ | ᴵ | ᴧ |
| KAPH | ⳛ | ⳛ | ⳛ | ⳛ |
| LAMEDH | ⳛ | ⳛ | ⳛ | ᴵ |
| MEM | ⳛ | ⳛ | ⳛ | ⳛ |
| NUN | (ⳛ) | ⳛ | ᴧ | ⳛ |
| SAMEKH | ⳛ | ⳛ | ⳛ | ⳛ |
| AYIN | (ⳛ) | (2) | ⳛ | ⳛ |
| PE | ⳛ | ⳛ | ⳛ | ⳛ |
| SADHE | ⳛ | ⳛ | ᴧ | ⳛ |
| QOPH | — | (ⳛ) | ⳛ | ⳛ |
| RESH | ⳛ | (2) | ⳛ | ⳛ |
| SHIN | ⳛ | ⳛ | ⳛ | ⳛ |
| TAW | ⳛ | ⳛ | ⳛ | ⳛ |

Table 5: Comparison of Old Sogdian and other Iranian scripts derived from Aramaic. Parenthesis indicate that the respective letter is not encoded as a distinct character and that it is represented in the encoded repertoire using another character. In the encoding for Inscriptional Pahlavi, *ayin* and *resh* are represented using WAW, and *qoph* is represented using MEM. Dashes indicate the absence of a given letter in the respective script.

| | Old Sogdian | Inscriptional Pahlavi | Inscriptional Parthian | Imperial Aramaic |
|---|---|---|---|---|
| ONE | ل | ؟ | ا | ا |
| TWO | لل | ؟؟ | اا | ٧ |
| THREE | للل | ؟؟؟ | ااا | ٧٧ |
| FOUR | لللل | ؟؟؟؟ | اااا | — |
| TEN | ا | ٦ | ע | ٦ |
| TWENTY | ٤ | ٤ | ه | ٦ |
| THIRTY | ٤ | — | — | — |
| ONE HUNDRED | ٢ | ٥ | ٢ | ٦ |
| ONE THOUSAND | للو | ٩ | ٦ | ٧ |
| TEN THOUSAND | — | — | — | ٢ |
| ONE HALF | ٬ | — | — | — |

Table 6: Comparison of numerical characters proposed for Old Sogdian and those encoded in Unicode blocks for Aramaic and Iranian scripts derived from it. Note: There is no distinctive numerical sign for 10,000 in Old Sogdian, instead this number is represented using the word ܝܘܣܝܘ *βrewar*.

Figure 1:  Excerpt of '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).
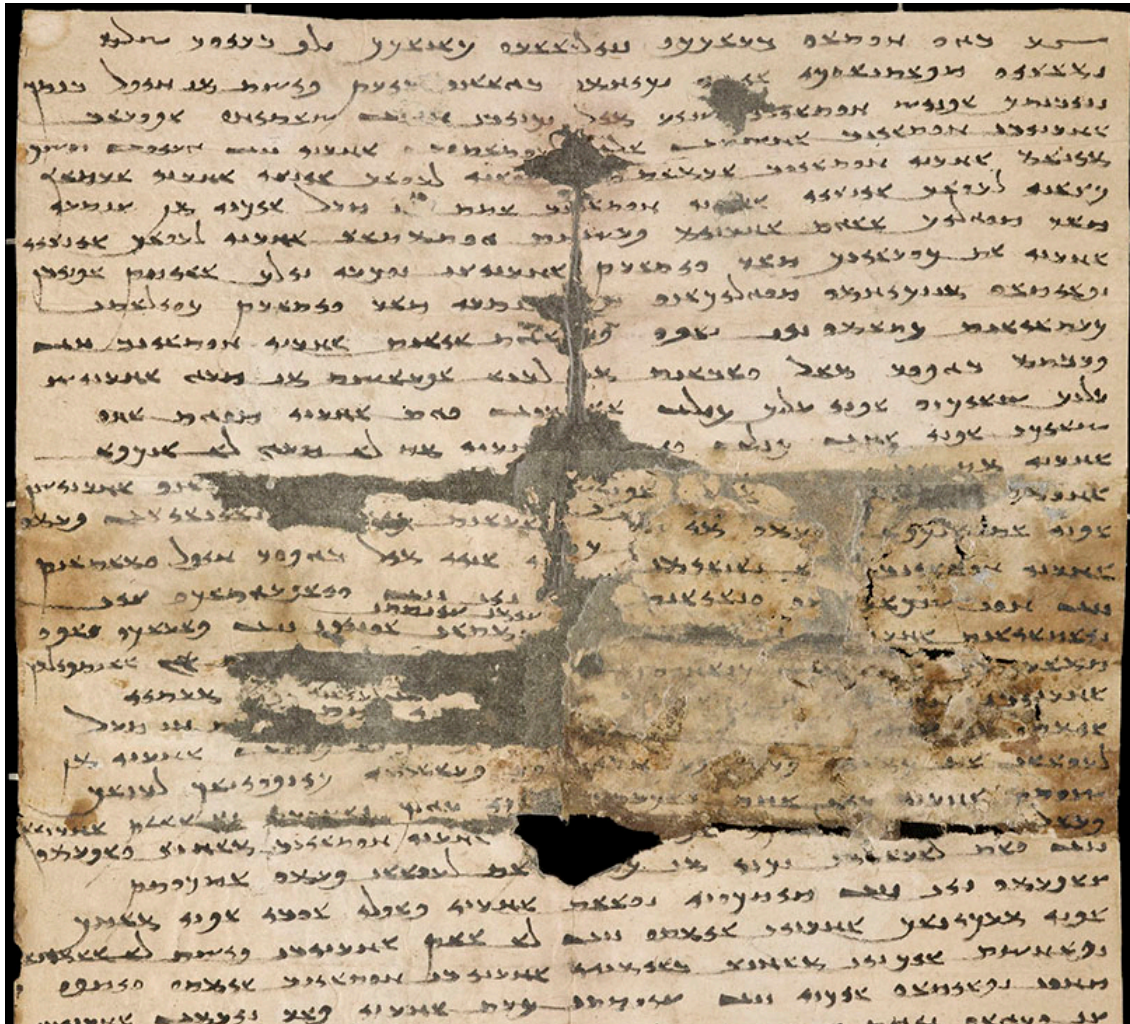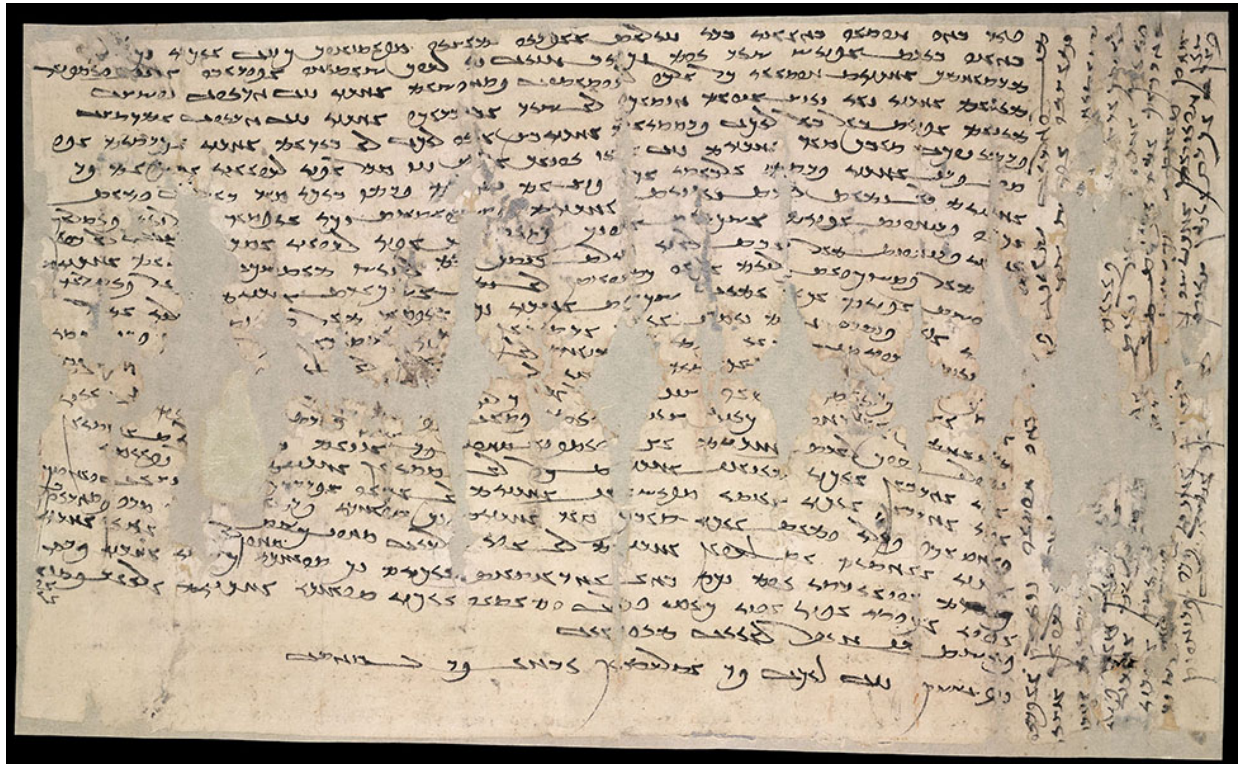
Figure 2: Excerpt of 'Ancient Letter 2' (British Library, International Dunhuang Project: Or. 8212/95 side a). "To the noble lord Varzakk (son of) Nanai-thvar (of the family) Kanakk. Sent [by] his servant Nanai-vandak." (translation by Sims-Williams in Waugh 2004).

Figure 3: Excerpt of 'Ancient Letter 3' (British Library, International Dunhuang Project: Or. 8212/98 recto 1). "From (his) daughter Shayn to the noble lord Nanai-dhat." (translation by Sims-Williams in Waugh 2004).
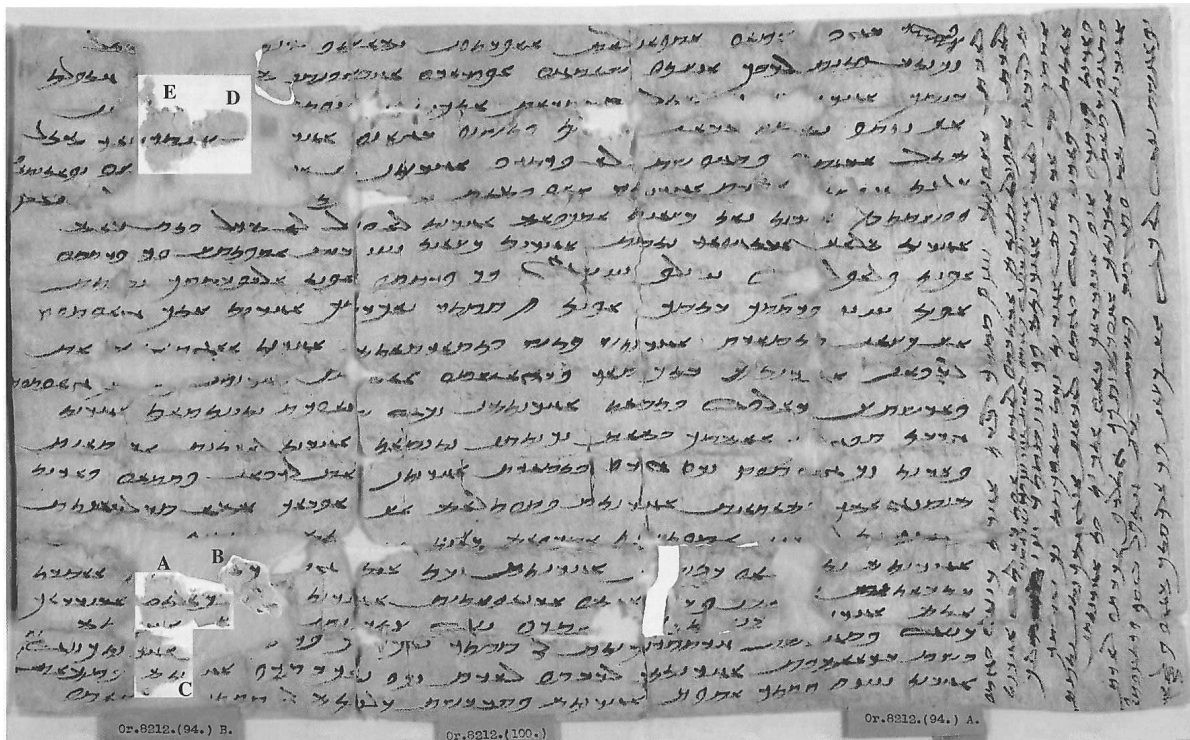
Figure 4: Excerpt of 'Ancient Letter 5' (from Grenet et al. 1998: 94). "To the noble lord, the chief merchant Aspandhãt. [Sent] by your servant [Frī-khwatãw]."
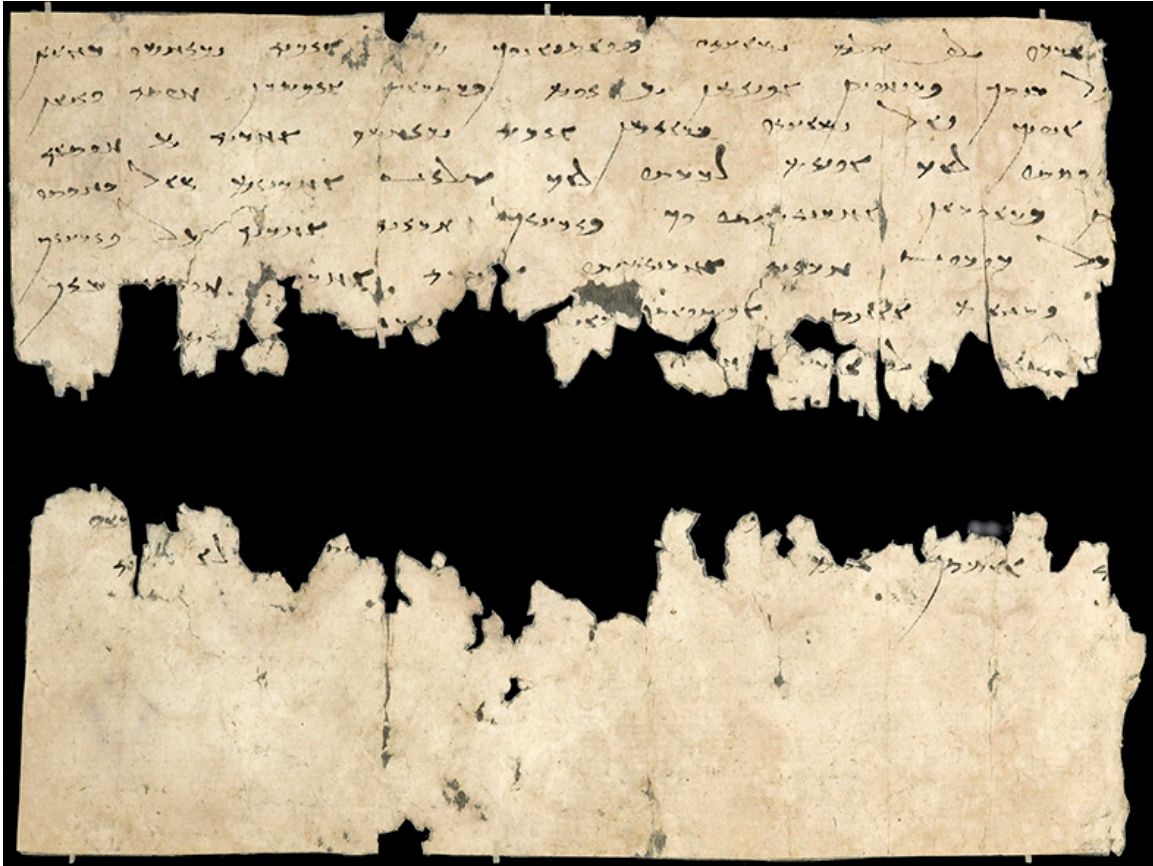
Figure 5: Excerpt of Ancient Letter 6 (British Library, International Dunhuang Project: Or. 8212/97).

Figure 6: Sogdian rock inscription from Shatial (from Sims-Williams 1989: plate 10b) The inscription reads ⟨ꢔꢔ⟩ *nny'kk ZK* (top line), ⟨ꢔꢔꢔ⟩ *sw'ßr* (middle), ⟨ꢔꢔ⟩ *BRY* (bottom). Latin transcription from *ibid*: 14. The inscription in the bottom right-hand corner is shown in detail in figure 7.



Figure 7: Sogdian rock inscription from Shatial (from Sims-Williams 1989: plate 10a). The central inscription reads ⟨ꢔꢔꢔ⟩ *p'p'kk* (top line), ⟨ꢔꢔꢔ⟩ *ZK kwš''n* (middle), ⟨ꢔꢔ⟩ *BRY* (bottom). Latin transcription from *ibid*: 14. The inscription in the top left-hand corner is shown in detail in figure 6.
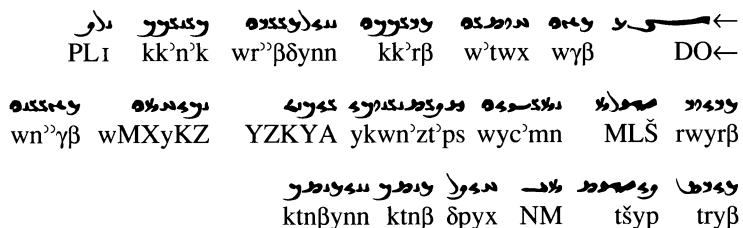
## Sogdian script

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, ' = **n**, '**n** = **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: γ/x/s/š, r/β/y. Some letters are distinguished only in final position (by some scribes), e.g., **n ~ z, x ~ γ**.

z is sometimes distinguished from **n** or z from ž by a diacritical point ⸲, and the foreign sound *b* was noted as ⸲ **ṗ**.

### SAMPLES OF SOGDIAN

ANCIENT LETTERS

| | | | | | |
|---|---|---|---|---|---|
| PLɪ | kk'n'k | wr'ʾβδynn | kk'rβ | w'twx | wγβ | DO← |

| | | | | | |
|---|---|---|---|---|---|
| wn'ʾγβ | wMXyKZ | YZKYA ykwn'zt'ps | wyc'mn | MLŠ | rwyrβ |

| | | | | |
|---|---|---|---|---|
| ktnβynn | ktnβ | δpyx | NM | tšyp | tryβ |

| | | | | | | |
|---|---|---|---|---|---|---|
| *1. Transliteration:* | OD | βγw | xwt'w | βr'kk | nnyδβ'ʾrw | k'n'kk |
| *2. Normalization:* | at | βaγu | xutāw | βarak | nanē-θβār | kanak |
| *3. Gloss:* | to | lord.ACC | master | Barak | Nana's-gift | Kanak |

| | | | | | | |
|---|---|---|---|---|---|---|
| *1.* ɪLP | βrywr | ŠLM | nm'cyw | sp'tz'nwky | AYKZY |
| *2.* (ēw-)zār | βrēwar | *āfrīwan | namācyu | spātzānūk | kaδ-uti |
| *3.* thousand | ten.thousand | greeting(?) | reverence.ACC | bended.knee | when-that.and |

| | | | | | | |
|---|---|---|---|---|---|---|
| *1.* ZKyXMw | βγ'ʾnw | βyrt | pyšt | MN | xypθ | βntk | nnyβntk |
| *2.* wēšanu | βaγān(u) | βyart | pišt | con | xēpθ | βantē | nanē-βantē |
| *3.* them.OBL | lords.OBL | 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.'
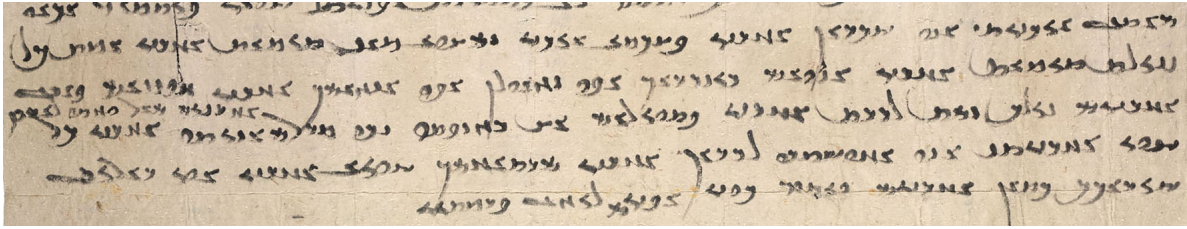
*—From the Old Sogdian "Ancient Letters" found in a mailbag in the Great Wall*
                                        *(AL II, Reichelt 1931: 12 and pl. 2).*

Figure 8: Description and specimen of the script in the 'Ancient Letters' (from Skjærvø 1996: 529).
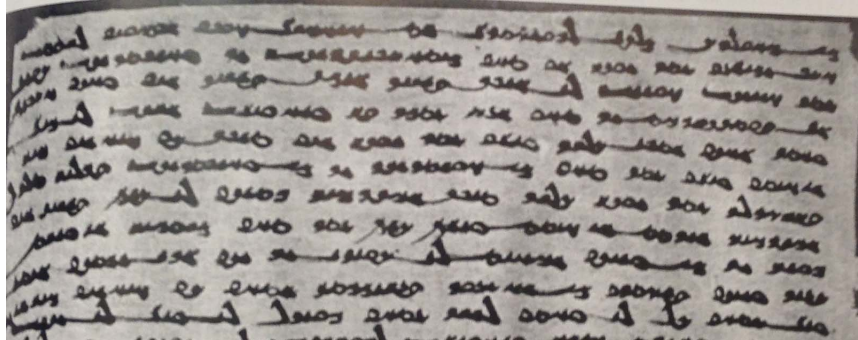
TABLE 48.2: *Main East Iranian Scripts Developed from Aramaic*

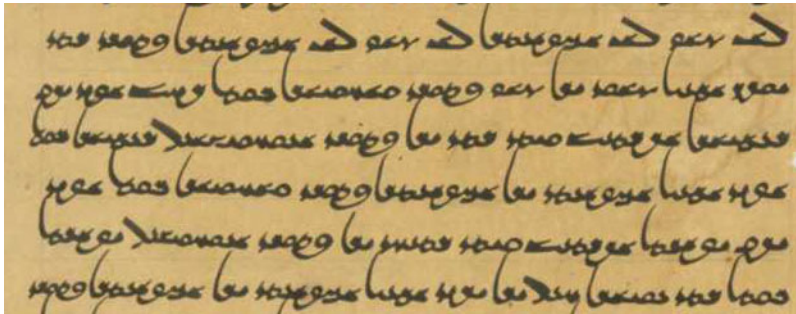| Aramaic | Sogdian Ancient Letters | Sogdian sutra script | Manichean Sogdian | Christian Sogdian | Principal Phonetic Values (Sogdian) |
|---|---|---|---|---|---|
| ʾ | | | | | a, ā |
| b | | | | | b, β |
| (β) | | | | | β |
| g | | | | | g, γ |
| (γ) | | | | | γ |
| d | | | | | d, δ |
| h (ḥ) | | | | | a, Ø |
| w | | | | | w, ŏ, ŭ |
| z | | | | | z |
| (j) | | | | | ž |
| (ž) | | | | | ž |
| ḥ (h) | | | | | γ, x, h |
| ṭ | | | | | t |
| y | | | | | y, ĕ, ĭ |
| k | | | | | k |
| (x) | | | | | x |
| l (δ) | | | | | δ |
| m | | | | | m |
| n | | | | | n |
| s | | | | | s |
| ʿ | | | | | Ø |
| p | | | | | p |
| (f) | | | | | f |
| ṣ (c) | | | | | č, ǰ |
| q | | | | | k |
| r | | | | | r |
| š | | | | | š |
| t | | | | | t, θ |

Figure 9: Table showing various scripts for writing Sogdian (from Skjærvø 1996: 519).
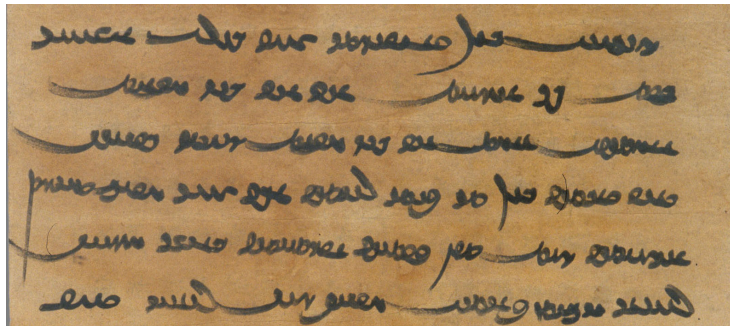
Sogdian script of the 'Ancient Letters' (excerpt from BL IDP: Or. 8212/92.1 recto 1)



'Cursive' Sogdian script of the Mount Mug documents (excerpt of B-18, from Livshits 2015: 99)



Sogdian 'sūtra' script (excerpt from Digitales Turfan-Archiv: So 14851, recto 1)



Sogdian 'Uyghur' script (excerpt from Digitales Turfan-Archiv: Ch/So 14744, recto 1)

Figure 10: Comparison of Old Sogdian (top) and later Sogdian scripts.