

Draft technical note:**Text representation and shaping specification of the Mongolian script****技术笔记草案: 蒙古文字的文本标记与变形规范**

To: Unicode Technical Committee
From: Liang Hai / 梁海 <lianghai@gmail.com>
Date: 8 October 2020

Document history 文档历史

This is a revision of the following:

- L2/19-368, *Draft technical note: Text representation and shaping specification of the Mongolian script* / 技术笔记草案: 蒙古文字的文本标记与变形规范
- L2/19-130 (WG2 5104, MWG/3-N8R), *Towards a well-formed Mongolian specification that allows interoperable implementations*

The following are the major changes since L2/19-368:

- Highlighted significant additions and changes with yellow background.
- *Unfortunately*, as the author is bilingual only in Chinese and English, Cyrillic Mongolian translations will only be available for title and section headings once document content is stable. A Hudum Mongolian translation will also be available for title.
- SSBM for decomposing vowel letters' marked medial forms.
- Removed default variants.
- Remove Phoenician names for written units.
- Rearranged written units O and U.
- [...]

And the major changes in L2/19-368 since L2/19-130:

- Rearranged sections.
 - Also added Chinese translations for title and section headings.
- Clarified design considerations.
- Revised character usage and shaping rules.
 - Unified U+183A KA into U+183B KHA.
 - Removed the special treatment of *η* as not involving a bowed written unit.
 - Clarified that ZWJ and ZWNJ should be kept from average users.
 - Preliminarily introduced data files.

- Rephrased the four encoding principles.
 - Briefly introduced the most confusable phonetic letter groups.
 - Rewrote the section about particles (previously “enclitics”).
- Revised terminology.
 - Certain concepts have been substituted with more commonly used ones: grapheme → written unit, enclitic → particle.
 - Certain translated native terms have been substituted with transcriptions: splash → chachlag, pad → devsger. [Add cross references or add glossary, probably in Introduction.]
 - Merged consonant coda and vowel offglide into the traditional concept devsger.
 - Revised transliteration scheme: $\ddot{a} \rightarrow e \rightarrow \acute{e}$, $x \rightarrow h \rightarrow \hbar$, $\mathbf{x} \rightarrow \mathbf{H} \rightarrow \mathbf{H}$.

Contents 目录

1. Introduction 引言
 2. Representation guidelines 表记指南
 3. Written units 书写单位
 4. Data files and the standard tooling 数据文件与标准工具
 5. Shaping and font production 变形及字体制作
 - 5.1 Encoded phonetic letters 编码的语音字母
 - 5.2 Orthographical shaping 正字法变形
 - 5.2 Particles 助词
 - 5.4 Uncaptured forms 漏网形
 6. Acknowledgements 致谢
 7. References 参考文献
-
- A. Principles 原则
 - A.1 Phonetic letters are encoded as characters 语音字母编码为字符
 - A.2 Contextual shaping is orthographical 自动变形基于正字法
 - A.3 Particles are grammatically encoded 助词根据语法编码
 - A.4 Uncaptured forms are manually requested 漏网形受手动请求
 - B. Architecture 架构
 - B.1 Character set 字符集
 - B.2 Shaping process 变形流程
 - C. The Twelve Syllabaries 十二字头
 - D. Orthography background 正字法背景
 - E. Generally recognized particles 普遍承认的助词

1. Introduction 引言

Two decades ago, Unicode/10646 added support for the Mongolian script with a seemingly sensible character set. Little was known at the time about how exactly these characters were meant to be rendered. Various vendors have since struggled to make their own senses of the characters, while users have been suffering from both the poor interoperability between vendor implementations and the consequent lack of native support on major platforms.

Experts have been uncovering issues of the Mongolian encoding and have proposed various patches, especially over the last few years. However, until very recently (MGC/01-01, 2018), few complete specifications have ever been available for the community to discuss and evolve in order to eventually reach an agreement.

This specification 本规范

This *draft* specification demonstrates and proposes how to specify both comprehensive guidelines for text representation and coherent rules for text shaping. It mainly deals with the first three layers in the overall technical architecture of rendering texts in the Mongolian script:

- The Unicode Standard and ISO/IEC 10646, a synchronized pair of standards, specify identities of encoded characters and standardized variation sequences.
- The Unicode Standard and its various supplementary standards further provide the characters with additional behavioral specifications, including character properties (general category, cursive joining type, etc.) and algorithms (normalization, collation, line breaking, text segmentation, bidirectional, vertical text layout, etc.).
- The required complex shaping is implemented in fonts and text shaping engines, according to OpenType Layout (OTL) or other complex shaping models, such as AAT (Apple Advanced Typography) and Graphite.
- Hudum text is treated inline as horizontal and left-to-right, while text layout engines are responsible for setting lines vertically and arranging multiple lines with the preferred left-to-right order.

Although the architectural defects cannot be resolved without migrating to a radically different encoding model, a meticulous specification can help eliminate unwanted differences between vendor implementations.

Design considerations of this specification include:

- Systematic (with the Twelve Syllabaries as the foundation) and unambiguous guidelines for text representation
- Moderate normalization of orthographical disagreements
- Minimal unnecessary confusability
- Explicitly modularized and ordered rules for text rendering, with well-defined fallback mechanisms
- Minimal dictionary-based shaping

How to read

[...]

Initial scope 初步范围

As a beginning, the scope has been restricted to the typical style of Hudum (ᠬᠣᠳᠤᠮ *budum*; Хүдәм *khudam*; “old” in Oirat Mongolian) writing system that is contemporarily used by the Mongolian language.

Eventually a single specification should cover all the major writing systems unified under the Unicode Mongolian encoding (i.e., Hudum, Todo, and Manchu-Sibe), as well as their Sanskrit-Tibetan extensions (Ali Gali letters and writing systems) and historical forms of the writing systems (early Hudum ones, early Todo, Old Manchu, etc.).

1 Strategies

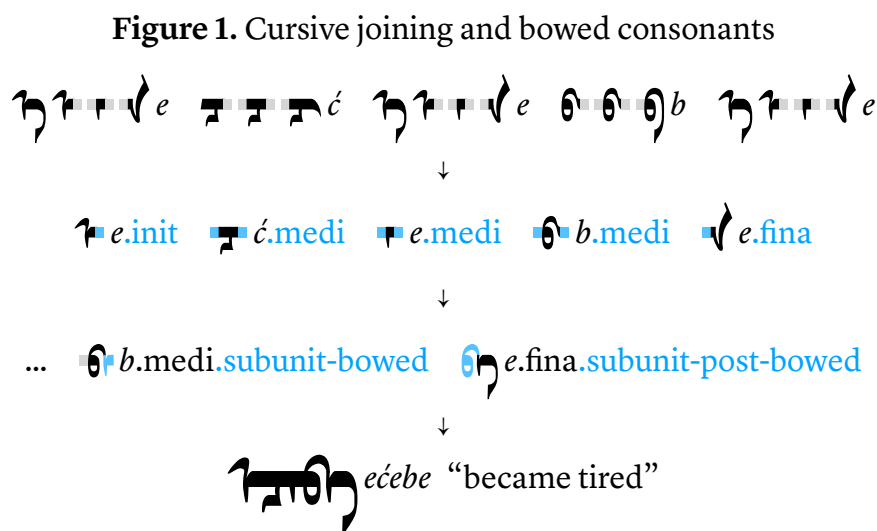
1. Maintain the phonetic encoding principle of the existing model.
 - That is to say: encode the underlying phonetic letters as characters, then contextually shape and manually control the encoded characters to reproduce the written forms.
 - If this principle were not maintained, we would be better off migrating to a graphic/graphetic encoding model designed from scratch.
 - The syllabic categories masculine/feminine/neuter vowel and consonant are critical to contextual shaping.
2. Minimize confusability by restricting the possible ways how a certain written form can be graphically encoded.
 - In particular, limit overlaps between the variant sets of phonetic letters.
3. Prioritize the modern orthographies in day-to-day usage.
 - Such a problematic encoding model needs to be standardized to be as restrictive as possible, and cannot be that extensible anymore for marginal and scholarly use cases.
 - Compromise historical orthographies when phonetically representing them harms the text exchange of modern orthographies.

A. Principles 原则

To a limited extent, Hudum is encoded in a way similar to how other cursive scripts and writing systems are encoded in Unicode. That is to say, instead of directly encoding written units that are already cursively joined to each other, certain underlying abstract units are considered to have the potential of *cursive joining* and are encoded as characters, then the actually displayed form of a character is contextually determined.

There is also an obligatory contextual variation process (which involves so-called *bowed consonants*) that is beyond cursive joining, but it still within the usual scope of text encoding (cf., obligatory lam–alef ligatures of Arabic).

For more information about these ordinary contextual shaping mechanisms, see Section 3.2, *Cursive joining*, and Section 3.3, *Written units*. Figure 1 shows a simple word that only involves these two mechanisms, where abstract characters with cursive joining potential and indefinite shapes contextually interact with each other, and get resolved into certain positional forms, then exhibit some additional interaction between a bowed consonant and its following letter.



What actually make the Hudum encoding unusual are the additional principles that extend far beyond ordinary text encoding:

1. **Phonetic letters are encoded as characters** 语音字母编码为字符
2. **Contextual shaping is orthographical** 自动变形基于正字法
3. **Particles are grammatically encoded** 助词根据语法编码
4. **Uncaptured forms are manually requested** 漏网形受手动请求

Specific rules are then derived from these principles for *how text should be encoded* (see Section 3, *Text representation*) and *how encoded text should be shaped* (Section 4, *Text rendering*).

2.1 Phonetic letters are encoded as characters

语音字母编码为字符

Hudum does not have a well-received system of typical letters (i.e., user-perceived primary units of writing) that is common to most writing systems. Instead, users are accustomed to identifying letters on a much more phonetic level, where letters are not reliably related to writing and are thus considered to be *phonetic letters* in this specification. Although the exact alphabet (i.e., the set and order of letters) varies considerably, Table 1 shows a typical version.

Table 1. Hudum phonetic letters

<i>a</i>	<i>e</i>	<i>é</i>	<i>i</i>	<i>o</i>	<i>u</i>	<i>ö</i>	<i>ü</i>									
a	ə	ə	i, ɪ	ɔ	ʊ	o	u									
<i>n</i>	<i>ŋ</i>	<i>b</i>	<i>p</i>	<i>h</i>	<i>g</i>	<i>m</i>	<i>l</i>	<i>s</i>	<i>ś</i>	<i>t</i>	<i>d</i>	<i>ć</i>	<i>j</i>	<i>y</i>	<i>r</i>	<i>w</i>
n	ŋ	p	p ^h	x, ɣ	k, q	m	l	s	ʃ	t ^h	t	tʃ ^h	tʃ	j	r	w
<i>f</i>	<i>k</i>	<i>c</i>	<i>z</i>	<i>h</i>	<i>ř</i>	<i>ł</i>	<i>ž</i>	<i>č</i>								
f	k ^h	ts ^h	ts	x	ɭ	ɬ	tɕ	tɕ ^h								

↖ Transliteration

← Representative written form

↙ Phoneme

Notes:

1. In this specification, phonetic letters are referred to with their single-letter transliterations that are always in *italics*.
2. The listed written forms here are not the full set, but are merely what commonly used as written representatives of phonetic letters.
3. Typical phonetic transcriptions of corresponding phonemes in the standard Chakhar Mongolian phonology are also provided for reference. Hudum phonetic letters do not however cover all phonemes in the contemporary spoken language.
4. The seven native vowel letters are often referred to as *vowel one* to *vowel seven*, in order to distinguish in writing the two visually identical pairs, *o/u* and *ö/ü*.
5. Light gray highlighted ones are loanword letters, thus often excluded from the alphabet. *p* and *w* were introduced early, and today are often considered to be native letters. *ŋ* is often excluded, because it is the only consonant letter that cannot occur in syllable onset positions. *ŋ* and *ł* are often excluded together when they are considered to be letter sequences *ng* and *lh*.
6. Letters *é* and *h* are often considered to be disambiguating written forms of *e* and *b*, respectively, and their duplicated phonemes are highlighted in light gray.
7. Letters *ř*, *ž*, and *č* are often excluded when they are considered to be special written forms of *r*, *j*, and *ć* for transcribing Chinese syllables each in only one syllable: (ᠷ, ᠵ, and ᠴ). Actually *ř* has extended usage today for other Chinese r-initial syllables and loanword [ɿ~ʒ] sounds from other languages.

Multi-to-one confusability. The system of phonetic letters is largely based on historical phonemes of the Mongolian written language reflected in the conservative orthography, instead of how the under-differentiated Hudum writing system actually

works with its limited set of written units. Many phonetic letters therefore do not have distinct written forms and can be confusable in writing, but are still identified as distinct letters because they are meant to be the abstract representatives of distinct phonemes.

Because modern pronunciations of words are actually distant from what spellings of Hudum literally suggest, identifying phonetic letters is a highly subjective process, and makes text representation fundamentally unreliable. Among others, the following groups of phonetic letters are especially prone to inconsistent identification:

- $a \sim e$
- $\acute{e} \sim w$
- $i \sim j \sim y$
- $o \sim u, \ddot{o} \sim \ddot{u}$, and $o \sim u \sim \ddot{o} \sim \ddot{u} \sim w$
- $h \sim g$ as in $(s|d)h$ vs. $(s|d)g$, due to phonetic alternation: $[\chi] \rightarrow [q] / [s, t] _$
- $s \sim \acute{s}$ as in si vs. \acute{si} , due to phonetic alternation: $[s] \rightarrow [ʃ] / _ [i]$
- $t \sim d$

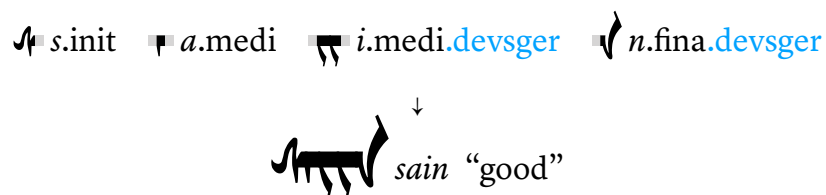
For the exact confusable forms of these phonetic letter groups, see Table 7, *Encoded phonetic letters and their written forms*.

One-to-multi unpredictability. Furthermore, the yellow highlighted phonetic letters in Table 1 can be written with multiple different written units and/or written unit sequences. As which written form exactly is used to represent a phonetic letter is determined with a combination of complex predictive rules and arbitrary variations, many phonetic letters have largely unpredictable correspondences with their written forms.

2.2 Contextual shaping is orthographical 自动变形基于正字法

Typically, orthographical features of a writing system are directly reflected in text representation (e.g., the English writing system requires the final consonant letter of certain words to be doubled when a suffix is joined). However in the Hudum encoding, as the encoded phonetic letters do not directly represent written forms, a number of orthographical rules are utilized to predict most written forms in order to minimize required manual controls in encoding. See Section 4.2, *Orthographical shaping*.

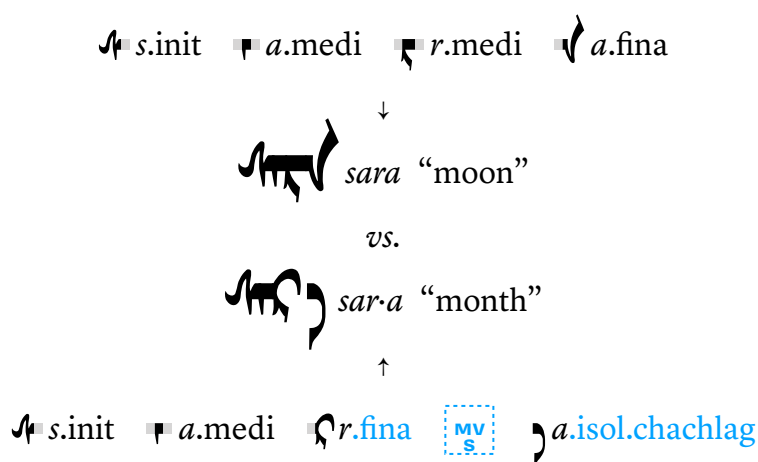
Figure 2. Orthographical shaping



Subjective and incomplete. The one-to-multi unpredictable nature of phonetic letters, combined with the orthographical flexibility of Hudum, means the usage of various written forms is ultimately unpredictable from each word’s phonetic letter sequence alone. Orthographical rules are also naturally subjective and incomplete because they are scholars’ summary of their observation on the writing system. The result is, manual overriding on the predictive orthographical shaping rules is inevitable.

Chachlag. The structure *chachlag* (𐌆𐌆𐌆𐌆𐌆𐌆𐌆𐌆; *čáčulg-a*; цацлага *tsatslaga*; “sprinkle”), a non-joining leftward tail form of *a* and *e* that only appears at the end of certain words, is a common lexical variation. Certain letters take special forms when followed by it. For more information about how chachlag is requested with a dedicated format control, see “Chachlag” in Section 4.2.

Figure 3. Chachlag

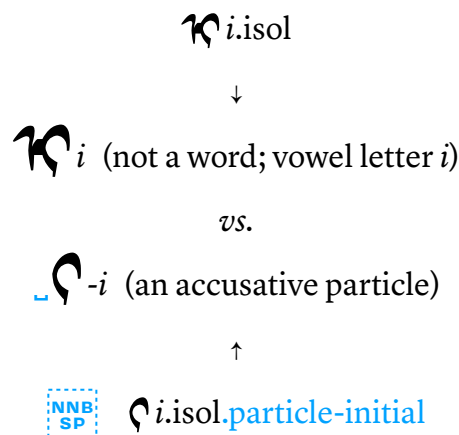


2.3 Particles are grammatically encoded

助词根据语法编码

There is a group of special, frequently used *particles* that are written separately from their modified word (as postposition words), but are conventionally considered to be grammatically suffixes (and are thus considered to be part of the modified word). Many of them exhibit special orthographical features that make them distinct from ordinary words. But as not all of them are graphically distinct, the scope of such words relies on highly subjective grammatical categorization.

Figure 4. Particles



All of these particles are then all marked at the encoding level by a special whitespace that connects a particle to its preceding word (the modified word or another particle). This whitespace signals the particles' special shaping requirement (which cannot be orthographically predicted from underlying phonetic letters alone), and aids special text processing treatments (for their close relationship with the modified word). See Section 4.3, *Particles*.

2.4 Uncaptured forms are manually requested

漏网形受手动请求

When an expected written form is not captured by predictive shaping rules, one of the three last-resort format controls, Free Variation Selectors (FVSes), is used to request the desired written form. Unlike MVS and NNBS which mark complex lexical or grammatical features then have the expected variations derived, an FVS only affects the base character it is applied to, having no extra effects. See Section 4.4, *Uncaptured forms*.

Theoretically, as all written forms on a given cursive position are assigned with an FVS, the predictive rules would be merely syntactical sugar built upon the FVS mechanism for minimizing the number of FVSes used. However, in order to avoid the pollution of unnecessary format controls, FVSes are explicitly rendered invalid wherever the predictive shaping is adequate.

Figure 5. Manual overriding

ᶑᶞ e.init ᶑᶞ d.fina.devsger

↓

ᶑᶞ ed (not a word; cf., ᶑᶞ on “year”)

vs.

ᶑᶞ ed “thing”

↑

ᶑᶞ e.init ᶑᶞ d.fina.uncaptured FVS
1

2 Representation guidelines 表记指南

Due to the encoding model's inherent complexity, Mongolian text representation is highly entangled with shaping rules.

This section specifies how exactly Hudum text should be encoded, and how properly encoded text should be shaped. The section starts with a thorough introduction of the fundamental building blocks:

- Written units as a concrete low-level abstraction for describing written forms
- Phonetic letters and their written forms

Then in order to fill the gap between encoded phonetic letters and actual written units:

- Variation patterns considered for reducing phonetic letters to written forms

The exact text representation of a word is thus determined by:

- Identifying the underlying sequence of phonetic letters by analyzing possible written forms and variation patterns, as well as ultimately, consulting a dictionary.
- Apply FVSes wherever phonetic letters are not reduced to the desired written forms by predictive rules.

Then tightly specified shaping rules are responsible for reproducing the written form of the word.

Informally, users usually first decide a sequence of phonetic letters, then interact with a comformant font until required format controls are placed.

The following are the formal steps to encode Mongolian texts:

1. Graphically segment the analog text into the written units for later comparison.
 - Section 3, *Written units*
2. Decide the basic characters by the assumed underlying phonetic letters.
 - Section A.1, *Encoded phonetic letters*
3. Apply MVS wherever a *chachlag* is in the analog text.
 - Section A.2, *Contextual shaping is orthographical*
4. Apply NNBSF wherever special written units present in the analog text in a grammatical particle.
 - Section A.3, *Particles are grammatically encoded*
5. Apply FVS whenever the shaped written units still do not match the analog text.
 - Section A.4, *Uncaptured forms are manually requested*






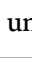



















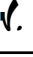





























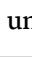







3 Written units 书写单位

Instead of the confusingly encoded characters for phonetic letters, the written units of Hudum are first introduced for building a concrete analysis model.

Before examining the encoded phonetic letters, encoding-independent *written units* are defined in Table 4. All written forms in this specification are analyzed with and formally referred to as sequences of written units on certain cursive positions, for the sake of clarity and accuracy.

Written units are assigned single-letter transliterations that are always in SMALL-CAPS, and are further marked by <...> when referred to inline. Phonetic letters *η* and *ł*, as well as certain variants of other phonetic letters, are written as sequences of written units, thus do not appear in the table.

Table 4. Written units

Written unit		Positional forms: .init, .medi, .fina, .isol		Sub-written unit variants		Represented phonetic letters	Note
A		   		• 		a/e 'a/ e/(n) 'a/e/n a/e/n	unified: 
Á	Aa					a/e	a
I		   		• 		i/j i/j/y i/y i/y	
İ	Ix					(i)	early modern orthographies
O		  		 		u/ü o/u/ ö/ü o/u/ ö/ü	
Ü	Ue						ö/ü
U						u/ü	o/u/ö/ ü/w
U̇	Ux					(u/ü)	early modern orthographies
N		  		• ?		n n n	see <A>
B		  	 			b b b	unified: 
P		  	 			p p p	
H		  		• •		b h/g h/g	
Ğ	Gh	  		• •		g g g	
G		 	 	•		h/g h/g	g
Ğ	Gx	 	 	•		(h) (h)	early modern orthographies
M		  				m m m	unified:  ; see <A>
L		  				l l l	see <A>
S		  		• •		s s s	

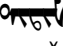
Š	Sh		• •	š š š	
T				t/d t t	
D			• ?	d t/d d	
Đ	Dd			d d	
Ć	Ch			ć ć ć	
J				j j	
Y			•	y y	
R			• ?	r r r	
W			• •	w w/é w/é	
F				f f f	
K				k k k	unified:
C				c c c	
Z				z z z	
H	Hh			h h h	see <A>
Ř	Rh			ř ř ř	see <A>
Ž	Zr			ž	
Č	Cr			č	

Positional forms. Positional forms of written units are denoted with an additional suffix. True isolated forms are absent for most written units, and many written units are not attested on all of the rest three positions either. An explicit cursive joining break is transliterated as “|”, and an explicit joining as “-”.

Sub-written unit variation involving bowed written units. Bowed-looking written unit groups <G> and , namely, <G>, <K>, , <P>, and <F> (on their initial and medial positions) cause certain written units (on their medial and final positions) to join perpendicularly without a stem when following. Attested perpendicularly joinable medial and final forms are marked in the Table 4 with either a bullet or the perpendicularly joined form (if significant structural change is observed).

Note the commonly mentioned “leftward tail” (or “feminine tail”) of <A> is here considered merely the perpendicularly joined form of <A.fina>. Written unit <U> lacks a perpendicularly joined final form, as the forms it would orthographically take (<O.fina> or, marginally, <Ü.fina>) have been disunified.

Stylistic variants. Certain commonly used stylistic variants of written units exhibit significant structural differences:

- All written units that involve a *titem* ( *titim*; ТИТЭМ *titem*; “crown”) on certain positional forms (i.e., <A>, <N>, <M>, <L>, <Ž>, and <H>) have a historical preferred, non-titem variant.
- Among the two variants of written unit <K>, the Ali Gali *kha* letterform (the default in the table) and the Manchu velar *k’a* letterform (which also appears to be similar to Ali Gali *ga* and Todo velar *ga*) are largely interchangeably used in China; while the Ali Gali *kha* letterform is strongly preferred in Mongolia.
- Written units <M> and both have a final variant that was historically preferred. Note that the Todo writing system still prefers these two variants and consequently has disunified characters (U+184B MONGOLIAN LETTER TODO BA and U+184F MONGOLIAN LETTER TODO MA).

4 Data files and the standard tooling

数据文件与标准工具

The Mongolian shaping data files formally captures Mongolian characters' additional identity information and shaping behavior in a machine-readable format. Section 5, *Shaping and font production*, demonstrates how to utilize the data files in font production.

I don't think we really need to squeeze Mongolian's shaping data into UCD. An independent set of data files can allow much cleaner data structure to be expressed in, say, JSON.

Table 7. Phonetic letter classes

<i>Class</i>	<i>Members</i>
masculine vowel	<i>a o u</i>
feminine vowel	<i>e é ö ü</i>
neuter vowel	<i>i</i>
derived: vowel	class <i>masculine vowel</i> + class <i>feminine vowel</i> + class <i>neuter vowel</i>
consonant	<i>' n ɣ b p h g m l s ś t d ć j y r w f k c z ħ ř ł ž č</i>

Phonetic letter classes can be recorded in a data file as part of the Unicode Character Database (UCD). For example in the following format:

```
# ...
```

The crucial information provided in Table 8, *Conditional forms of phonetic letters* (Section 4.2, *Orthographical shaping*) can also be recorded in a data file, with a different format:

```
# ...
```

5 Shaping and font production 变形及字体制作









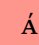
























This section introduces the shaping rules by providing a tutorial about how to produce a font.

5.1 Encoded phonetic letters 编码的语音字母

Phonetic letters are encoded as abstract characters that may be used to represent any written forms of a given phonetic letter. Attested written forms are defined as positional forms of written unit sequences, as shown in Table 5.

Yellow highlighted phonetic letters exhibit additional written forms (also highlighted in yellow) besides a single set of positional forms. For such complex characters, their default positional forms are specified to be the most reasonable forms (which are stray forms for consonant letters that are affected by syllable structure) when cursive positions are caused by medial forms of an abstract, ordinary letter. Red highlighted written forms are the ones only accessible with manual overriding.

Table 5. Encoded phonetic letters and their written forms

Phonetic letter & character	Written forms: positional forms, including contextual variants & manual variants	Note
' 1807 SSBM	 A	
a 1820 A	 	<ul style="list-style-type: none"> • A.isol/init—<i>particles</i> • Á.isol—<i>chachlag</i> • Á.fina—<i>Ali Gali</i>
	   	
	 	
e 1821 E	   	• Á.isol— <i>chachlag</i>
		
é 1827 EE	 	Cross-written unit simple letter
	 	
i 1822 I	 	<ul style="list-style-type: none"> • I.isol/init—<i>particles</i> • II.medi—<i>devsger</i> • I.isol—<i>early modern orthographies</i>
	   	
		
		
o 1823 O	 	<ul style="list-style-type: none"> • O.fina—<i>initial-body, post-bowed, disambiguating</i>
	 	
		
u 1824 U	 	

			ᵀ	ᵀ	ᵀ		O 2	O	O 1	<ul style="list-style-type: none"> • U.isol/init—<i>particles</i> • O.fina—<i>initial-body, post-bowed</i>
		ᵇ			ᵇ	U 2			U 2	<ul style="list-style-type: none"> • ᵇ.isol—<i>early modern orthographies</i>
		ᵇ̇				ᵇ̇ 3!				
ö	1825 OE	ᵀ	ᵀ			AÜ	AOI			
			ᵀ	ᵀ				OI 1	Ü 1	<ul style="list-style-type: none"> • UI.medi, Ü.fina—<i>initial-body</i> • O.fina—<i>post-bowed</i>
			ᵀ	ᵀ				O 2	O 2	
				ᵇ					U 3	
ü	1826 UE	ᵀ	ᵀ			AÜ 1	AOI 1			
			ᵀ	ᵀ				OI 1	Ü 1	<ul style="list-style-type: none"> • U.isol/init—<i>particles</i> • UI.medi, Ü.fina—<i>initial-body, disambiguating</i>
			ᵀ	ᵀ	ᵀ			O 2	O 2	O 2
		ᵇ̇				AU 2!				<ul style="list-style-type: none"> • AU.isol—<i>lexical</i> • O.fina—<i>post-bowed</i>
		ᵇ			ᵇ	U 3			U 3	<ul style="list-style-type: none"> • ᵇ.isol—<i>early modern orthographies</i>
		ᵇ̇				ᵇ̇ 4!				
n	1828 NA		ᵇ̇	ᵇ̇	ᵇ̇		N	N 1	N 1	<ul style="list-style-type: none"> • N.fina—<i>pre-chachlag</i> • A.medi—<i>devsger</i> • A.init—<i>early modern orthographies</i>
			ᵇ̇	ᵇ̇	ᵇ̇		A 1	A 2	A 2	
ŋ	1829 ANG		ᵇ̇	ᵇ̇				AG	AG	.init unattested
b	182A BA		ᵇ̇	ᵇ̇	ᵇ̇		B	B	B	
p	182B PA		ᵇ̇	ᵇ̇	ᵇ̇		P	P	P	
h	182C QA		ᵇ̇	ᵇ̇	ᵇ̇		H 1	H 1	H 1	<ul style="list-style-type: none"> • G.init/medi—<i>feminine</i> • H.fina—<i>pre-chachlag</i>
			ᵇ̇	ᵇ̇			G 2	G 2		
			ᵇ̇	ᵇ̇			Ğ 3	Ğ 3		
g	182D GA		ᵇ̇	ᵇ̇	ᵇ̇		Ğ 1	Ğ 1	Ğ 1	<ul style="list-style-type: none"> • G.init/medi/fina—<i>feminine</i> • H.medi—<i>devsger</i> • Ğ.fina—<i>pre-chachlag</i>
			ᵇ̇	ᵇ̇	ᵇ̇		G 2	G 2	G 2	
			ᵇ̇	ᵇ̇				H 3	H 3	
m	182E MA		ᵇ̇	ᵇ̇	ᵇ̇		M	M	M	
l	182F LA		ᵇ̇	ᵇ̇	ᵇ̇		L	L	L	
s	1830 SA		ᵇ̇	ᵇ̇	ᵇ̇		S	S	S	
ś	1831 SHA		ᵇ̇	ᵇ̇	ᵇ̇		Ś	Ś	Ś	

<i>t</i>	1832 TA				T	T 1	T	• T.medi— <i>disambiguating</i>
						D 2		
<i>d</i>	1833 DA				T 1			• D.init/fina— <i>disambiguating</i> • D.medi— <i>devsger</i>
						D 2	D 1 D 1	
							D 2 D 2	
<i>ć</i>	1834 CHA				Ć	Ć	Ć	
<i>j</i>	1835 JA					J	J	Lack of an explicit fallback .isol. • I.isol— <i>pre-chachlag</i>
					I 1	I		
<i>y</i>	1836 YA					Y 1	Y 1	Lack of an explicit fallback .fina. • I.init/medi— <i>particles</i> • I.fina— <i>pre-chachlag</i>
						I 2	I 2 I 1	
<i>r</i>	1837 RA				R	R	R	
<i>w</i>	1838 WA				W	W	W 1	• U.fina— <i>pre-chachlag</i>
							U 2	
<i>f</i>	1839 FA				F	F	F	
<i>k</i>	183B KHA				K	K	K	
<i>c</i>	183C TSA				C	C	C	
<i>z</i>	183D ZA				Z	Z	Z	
<i>h</i>	183E HAA				H	H	H	
<i>ř</i>	183F ZRA				Ř	Ř	Ř	
<i>l</i>	1840 LHA				LH	LH		.fina unattested
<i>ž</i>	1841 ZHI				Ž			.medi/fina unattested
<i>č</i>	1842 CHI				Č			.medi/fina unattested

Notes:

1. The written unit sequence II.medi (in written forms of *öi/üi* or *devsger i*, etc.) has a stylistic variant that looks similar to AI.medi , particularly in handwritten styles. The written unit I.init of *y* also has a stylistic variant that looks similar to A.init , particularly in handwritten particle *yin* (which consequently looks similar to an ordinary *in*). Text representation should not be affected by such stylistic variations.

These characters are more abstract than what are commonly encoded for other cursive joining scripts, as not only do they involve positional forms for cursive joining, but some of them also involve multiple possible written forms on a single cursive position and thus require examining additional variation patterns for determining a character. By comparison, cross-written unit variations are also involved in various Arabic-script

writing systems, but they are always constrained within the cursive joining mechanism in the Unicode Arabic encoding.

Cursive joining fallback: isol > init > medi > fina in general. But because there is not a default positional form, each fallback form is explicitly specified.

Variant identities. Variant register.

The contemporary orthographies need to be prioritized. Avoid increasing confusability when covering historical orthographies. Prescribe graphical encoding strategies for historical spellings, as long as the consonant–vowel categorization isn’t crossed.

Because characters are encoded on a phonetic basis, many written forms are shared across multiple characters (some character pairs even have fully overlapping sets of written forms). Therefore, identifying a word’s underlying sequence of phonetic letters from its written form is highly subjective, and requires knowledge such as orthography, grammar, and ultimately a dictionary, which are far beyond the basic script behavior. Hudum text representation is thus even more complicated than that of Indic scripts, which also involve multiple variation mechanisms but are generally contextually predictable and self explanatory. Faced with the overly complicated variation mechanisms of Hudum, users often turn to piece together a written word graphically with phonetically arbitrary characters.

Ordinary letters. For the sake of this specification, an ordinary letter is defined as a consonant letter that does not involve a bowed written unit.

5.2 Orthographical shaping 正字法变形

Table 8. Conditional forms of phonetic letters

Phonetic letter & character	Condition	Conditional forms				Note
<i>a</i> 1820 A	Chachlag	ا		أ		
	Particle	ا	آ	أ	ا	
	Fallback	ا	آ	AA	AA	
<i>e</i> 1821 E	Chachlag	ه		أ		
	Particle	ه		أ		
	Fallback	ه		أ		
<i>i</i> 1822 I	Particle	ي	ي	I	I	
	Devsgar		ي		II	
	Fallback	ي	ي	AI	AI	I
	Marked		و			O

<i>o</i> 1823 O	Post_Bowed								
	Fallback								
<i>u</i> 1824 U	Particle								
	Marked								
	Post_Bowed								
	Fallback								
<i>ö</i> 1825 OE	Marked								
	Post_Bowed								
<i>ü</i> 1826 UE	Particle								
	Marked								
	Post_Bowed								
<i>n</i> 1828 NA	Onset								
	Devsgger								
	Chachlag_Onset								
<i>b</i> 182C QA	Masculine_Onset								
	Feminine								
	Chachlag_Onset								
<i>g</i> 182D GA	Masculine_Onset								
	Masculine_Devsgger								
	Feminine								
	Chachlag_Onset								
<i>d</i> 1833 DA	Onset								
	Devsgger								
	Particle								
<i>j</i> 1835 JA	Chachlag_Onset								
<i>y</i> 1836 YA	Chachlag_Onset								
	Dictionary_Particle								
<i>w</i> 1838 WA	Chachlag_Onset								

Table 6. Hudum shaping: phonetic and written unit-level

<i>Shaping step</i>	<i>Subject letters</i>	<i>Rules</i>	<i>Resulted condition</i>
1. Chachlag	<i>a/e</i>	if follows MVS:	Chachlag
2. Syllabic	<i>o/u/ö/ü</i>	if follows an initial <i>consonant</i> :	Marked
	<i>n/j/y/w</i>	if precedes MVS that precedes <i>a/e.isol</i> :	Chachlag_Onset
	<i>h/g</i>	if precedes MVS that precedes <i>a.isol</i> :	Chachlag_Onset
	<i>n/d</i>	if precedes a <i>vowel</i> :	Onset
		else if follows a <i>vowel</i> :	Devsgger
	<i>h/g</i>	if precedes a <i>masculine vowel</i> :	Masculine_Onset
		else if precedes a <i>feminine vowel</i> or <i>neuter vowel</i> :	Feminine
	<i>g</i>	else if follows a <i>masculine vowel</i> :	Masculine_Devsgger
		else if follows a <i>feminine vowel</i> :	Feminine
		else if remotely follows a <i>masculine vowel</i> without a blocking <i>feminine vowel</i> :	Masculine_Devsgger
		else:	Feminine
3. Particle	<i>a/i/u/ü/d</i>	if follows NNBSP:	Particle
	<i>u/ü</i>	if follows an initial <i>consonant</i> that follows NNBSP:	Particle
	<i>y</i>	if is in a word <i>yin/yi/iyar/iyer/iyen/iyen</i> that follows NNBSP:	Dictionary_Particle
4. Devsgger <i>i</i>	<i>i</i>	if follows a <i>vowel</i> written form that does not end with a written unit <i>i</i> :	Devsgger
5. Post-bowed	<i>o/u/ö/ü</i>	if is in the written form of U and follows a bowed written unit G/K/B/P/F:	Post_Bowed

When a specific positional written form of a letter is decided by a font to be in one of the conditions above, it is converted to the specified conditional form in Table 8.

Unspecified positions in the table do not affect the subject written form.

5.3 Particles 助词

MVS-NNBSP merger.

The special whitespace between a particle and its preceding word (either the modified word or a preceding particle) is sometimes considered an internal gap in a grammatical word, and thus is sometimes preferred to be non-line-breaking, non-word-breaking, and

narrower than an ordinary space. Typical particles exhibit at least one of the following variations:

- Absence of onset placeholder and initial-body variation (because the beginning of a particle it is not the beginning of a stem).
- An initial *d* takes its disambiguating lamedh form.
- The under-differentiated, historical form of *y*.
- Letters *h* and *g* show gender harmony with the preceding word. The gender features are only observable on these two letters because the first feature has neutralized gender-distinguishing features of all vowel letters.

NNBSP is used both for representing this whitespace and for requesting the special variations shown in particles. As the set of particles and usage of NNBSP are decided grammatical, a particle may or may not exhibit special variations but is still encoded with a preceding NNBSP. See Appendix C for a reference list and comparison.

The step 7 does not involve contextual effects, as an FVS only affects the base character it is applied to. FVSes are only used to request desired written forms that are not captured by all the previous predictive shaping rules in steps 1–6. [*To be elaborated in a later revision.*]

5.4 Uncaptured forms 漏网形

Unicode Technical Report #54, *Unicode Mongolian 12.1 Snapshot*.

The step 7 does not involve contextual effects, as an FVS only affects the base character it is applied to. FVSes are only used to request desired written forms that are not captured by all the previous predictive shaping rules in steps 1–6. [*To be elaborated in a later revision.*]

Fallback forms of invalid FVSes: unnecessary vs. undefined.

6. Acknowledgements 致谢

The author would like to thank Shen Yilei / 沈逸磊 for his always enlightening feedback. Wang Yihua / 王奕桦 contributed to the analysis of letter *k*'s variants.

Liang Jinbao / 梁金宝 and Narantsogt Baatarkhuu / Баатархүүгийн Наранцогт provided the title's Hudum and Cyrillic Mongolian translations, respectively.

The major font used for Hudum examples is Menk Vran Tig (version 1.02), one of the fonts freely released by Menksoft (<http://font.menksoft.com>). Menk Garqag Tig (1.03) and Menk Scnin Tig (1.02) are also used.

7. References 参考文献

- Bao Yuzhu/ 宝玉柱 and Menghebaoyin/ 孟和宝音. 2011. 现代蒙古语正蓝旗土语音系研究. 北京: 民族出版社.
- Mongolian Research Institute, School of Mongolian Studies, Inner Mongolia University / 内蒙古大学蒙古学研究院蒙古语文研究所. 1999. 蒙汉词典. 呼和浩特: 内蒙古大学出版社.
- Nicholas Poppe. 1954. *Grammar of Written Mongolian*. Wiesbaden: Harrassowitz Verlag.
- Qinggeertai/ 清格尔泰. 1991. 蒙古语语法. 呼和浩特: 内蒙古人民出版社.
- Quejingzhabu/ 确精扎布. 2000. 蒙古文编码. 呼和浩特: 内蒙古大学出版社.
- Rita Kullmann and D. Tserenpil. 1996. *Mongolian Grammar*. Hong Kong: Jensco Ltd.

Standards 标准

- GB/T 25914-2010, 信息技术 传统蒙古文名义字符、变形显现字符和控制字符使用规则 / *Information technology—Traditional Mongolian nominal characters, presentation characters and use rules of controlling characters*. 北京: 中华人民共和国国家质量监督检验检疫总局 and 中国国家标准化管理委员会. Accessed from <http://www.gb688.cn/bzgk/gb/newGbInfo?hcno=62808E0BCB8246A287CFD9CF795ECF94>.
- MGC/O1-01 (2018, version 1.0.2), 信息技术 传统蒙古文名义字符到变形显现字符的转换补充规则 / *Information technology—The Transferring Rules of Traditional Mongolian Nominal Form to Variant Form*. 呼和浩特: 内蒙古自治区民族事务委员会 and 内蒙古大学. <http://nmgmzw.gov.cn/nmmwh/gsgg/201808/5938899e00fc43aebd189acaa5c6f9e4.shtml>.
- MNS 4932: 2000, *Монголжин бичгийн кодыг хэрэглэх дүрэм* / *Use of Mongolian Character Encoding*. Улаанбаатар: Стандартчилал, хэмжилзүйн үндэсний төв. Accessed on 7 September 2018 from http://estandard.gov.mn/index.php?module=standart&cmd=standart_desc&sid=7813.
- The Unicode Standard*, Version 13.0 (2020). Mountain View: The Unicode Consortium. <https://unicode.org/versions/Unicode13.0.0/>.

UTC Document Registry UTC 文档记录处

- L2/17-036, *Encode Mongolian Suffix Connector (U+180F) To Replace Narrow Non-Breaking Space (U+202F)*. Greg Eck, Andrew West, Badral Sanlig, Siqinbilige, and Ou Rileke.
- L2/17-332, *Positional Mismatches in Mongolian Encoding*. 沈逸磊 / Shen Yilei.
- L2/18-293, *Solution for NNBS Issues*. Badral Sanlig and Munkh-Uchral Enkhtur.
- L2/19-247 (WG2 N5054), *Recommendations from WG 2 meeting 68*.

Other documents 其他文档

- Unicode Technical Report #54, *Unicode Mongolian 12.1 Snapshot*. Ken Whistler. <https://www.unicode.org/reports/tr54/>
- Myatav Erdenechimeg, Richard Moore, and Yumbayar Namsrai. 1999. UNU/IIST Report No. 170, *Traditional Mongolian Script in the ISO/IEC 10646 and Unicode Standards*. Macau: UNU/IIST. Accessed from <http://babelstone.co.uk/Mongolian/>

Report170.pdf, <http://babelstone.co.uk/Mongolian/Report170A.pdf>, and
<http://babelstone.co.uk/Mongolian/Report170B.pdf>.

B. Architecture 架构

This section introduces in summary how Hudum text is encoded and shaped. The following aspects are covered:

- Required characters
- Expected shaping mechanisms
- Cursive script shaping mechanism
- Hudum-specific shaping mechanism

For exact, Hudum-specific specification of text. representation and shaping, see Section 4, *Rules*.

3.1 Character set 字符集

A typical Hudum implementation requires the characters shown in Table 2, including Mongolian-specific ones and characters that are shared with other scripts.

Table 2. Required characters

<i>Script</i>	<i>Type of characters</i>	<i>Character or range of characters</i>	<i>Note</i>
General	Space (1)	0020 SPACE	
	Format controls (3)	200C ZERO WIDTH NON-JOINER	
		200D ZERO WIDTH JOINER	
		202F NARROW NO-BREAK SPACE	
	Misc	·()«»<>0123456789 <i>etc.</i>	
Mongolian	Punctuation (6)	1800 MONGOLIAN BIRGA	BIRGA and FOUR DOTS have limited modern usage
		..1805 MONGOLIAN FOUR DOTS	
	Format controls (5)	180A MONGOLIAN NIRUGU	
		..180E MONGOLIAN VOWEL SEPARATOR	
	Digits (10)	1810 MONGOLIAN DIGIT ZERO	Limited modern usage
		..1819 MONGOLIAN DIGIT NINE	
	Phonetic letters (35)	1807 MONGOLIAN SIBE SYLLABLE BOUNDARY MARKER	U+1807 SSBM represents the syllable-initial consonant placeholder for medial forms; U+183A KA is excluded as its represented letterform is now unified into U+183B KHA
		1820 MONGOLIAN LETTER A	
		..1839 MONGOLIAN LETTER FA	
		183B MONGOLIAN LETTER KHA	
		..1842 MONGOLIAN LETTER CHI	

Yellow highlighted characters are involved in the complex shaping of Hudum. All the Mongolian-specific characters listed in the table above are encoded in the main Mongolian block (U+1800..U+18AF). The other block, Mongolian Supplement (U+11660..U+1167F), currently only have 13 characters for variants of the *birga* sign (U+1800 MONGOLIAN BIRGA), which are not typically used in day-to-day text.

For information about phonetic letters, see Section 3.3, *Written units* and Section 3.4, *Encoded phonetic letters*. The following subsections introduce other types of characters.

*3.1.0 Notable changes

- U+183A KA
- SSBM

Potential additional improvements:

- Unify h/g, t/d, then use existing characters to encode them graphically.
- Deprecate o and ö.

3.1.1 Format controls 格式控制符

[Clarification: including characters that have implication on shaping and forming of Hudum, not just Unicode Format_Control]

Zero Width Non-Joiner (ZWNJ), Zero Width Joiner (ZWJ), and nirugu. U+200C ZERO WIDTH NON-JOINER and U+200D ZERO WIDTH JOINER are Unicode’s standard cursive joining controls. Note that ZWJ also breaks interaction (such as ligation) between consecutive two letters as it is treated as an invisible letter. U+180A MONGOLIAN NIRUGU is a Mongolian-specific modifier letter that behaves exactly like ZWJ but is visible as a piece of stem stroke. See Section 3.2, *Cursive joining*.

In particular, for the Hudum-specific shaping steps, ZWNJ acts like an ordinary space (U+0020), while ZWJ and nirugu act like medial forms of an ordinary letter (defined in Section 2.1).

ZWNJ and ZWJ, however, should be kept from average users’ access on common keyboard layouts, as day-to-day text does not need these tricky characters.

The visible character nirugu should be used for causing joining in day-to-day text. A common use case is terminating a patronymic abbreviation, which the initial syllable body (i.e., an optional onset plus the first vowel) or merely the initial consonant letter of one’s father’s name.

Vowel Separator (MVS) and Narrow No-Break Space (NNBSP). MVS is a Mongolian-specific format control for requesting the chachlag variation. It is transcribed as “.” (a middle dot). See “Chachlag” in Section 3.5, *Additional variation patterns*. NNBSP is a Mongolian-specific whitespace and format control for marking and shaping particles,

and note that it is also used as a general whitespace by other scripts. See “Particles” in Section 3.5.

In terms of cursive joining behavior, MVS and NNBS are both non-joining inline characters, like an ordinary space.

Free Variation Selectors (FVSes). Mongolian-specific format controls. As combining marks, they are applied to certain characters for requesting the forms not captured by predictive shaping rules. See Section 4.4, *Uncaptured forms*.

3.1.2 Other characters 其他字符

Mongolian-specific punctuation characters do not yet have well-defined spacing behavior. For example, it is inconsistent in implementations if U+1802 MONGOLIAN COMMA, U+1803 MONGOLIAN FULL STOP, and other punctuation characters (as they tend to have significant and balanced spacing on both sides) are both preceded and followed by whitespace characters (e.g., U+0020 SPACE) or just have the preceding spacing as part of their glyphs.

The choices of non-Mongolian-specific punctuation characters have been heavily influenced by what characters are used in Chinese text, and are often not the ideal choices. Certain CJK punctuation characters are indeed beneficial because of their upright, non-rotated appearance in vertical text, but it is a problem to actually clarify how spacing and positioning of these fullwidth characters should be adapted for Mongolian usage.

3.2 Shaping process 变形流程

The shaping process of Hudum is based on the well-implemented technology foundation for general scripts and cursive scripts, while an additional phase of Hudum-specific shaping steps is inserted into the ordinary shaping process required by cursive scripts. The minimal shaping process consists of a number of steps as shown in the table below.

Table 3. Overview of shaping process

<i>Shaping phase</i>		<i>Shaping step</i>
IA. General		• Basic character-to-glyph mapping
IIA. Cursive script		• Initiation of cursive positions
III. Hudum-specific <i>Reduction of phonetic letters to written units</i>	<i>Phonetic:</i>	1. Chachlag 2. Syllabic 3. Particle
	<i>Graphemic:</i>	4. Devsger <i>i</i> 5. Post-bowed

Uncaptured: 6. FVS-selected	
IIb. Cursive script (continued) <i>Sub-written unit variations</i>	<ul style="list-style-type: none"> • Variation involving bowed written units • Cleanup of format controls • Optional treatments
IB. General (continued) <i>Typography</i>	<ul style="list-style-type: none"> • Vertical forms of punctuation marks • Optional treatments

For details about the phase III, see Section 3.3, *Hudum-specific shaping phase*.

3.2.1 General shaping phases 通用变形阶段

These are the basic mechanisms in fonts, applicable for all scripts.

Basic character-to-glyph mapping (phase IA) is typically controlled by the TrueType/OpenType table “cmap”. The Unicode representative glyphs may be used here as the default glyph mappings for phonetic letters, however these representative glyphs are not actually kept in the final rendering in typical implementations.

Vertical forms of punctuation marks (phase IB) are critical to proper typesetting of Hudum text, but are not part of the complex shaping between letters and format controls.

3.2.2 Cursive script shaping phases 连写文种变形阶段

On top of the general shaping mechanisms, complex scripts require additional shaping phases to be inserted after the basic character-to-glyph mapping and before typographical treatments. In particular, cursive scripts all undergo the cursive joining mechanism.

Cursive joining. Written forms exhibit the cursive joining mechanism. Both sides of a written form can be either joined to a neighboring written form or not, exhibiting up to four different statuses. Or, abstractly speaking, each written form is in one of the four cursive positions:

- *isolated*, abbreviated as “isol”: not joined above, not joined below
- *initial*, “init”: not joined above, joined below
- *medial*, “medi”: joined above, joined below
- *final*, “fina”: joined above, not joined below

Cursive positions are irrelevant to word boundaries, although they are usually consistent with word-wise positions in Hudum because cursive joining breaks inside a word are limited in the writing system.

For the exact algorithm used here for initiating cursive positions, see “Arabic Cursive Joining” in Section 9.2, *The Unicode Standard, Version 12.0—Core Specification*.

Implementation. The originally mapped glyphs from the last phase are converted to default positional forms in phase IIA. Although these default positional forms are not necessarily kept till the end of shaping process, the cursive positions are immutable once initiated.

Unattested default positional forms, especially default isolated forms of consonants, are recommended to be explicitly marked invalid in rendering. It is recommended to implement variation of bowed written units with contextual glyph variants, although many fonts use ligatures.

Graphemic variation after bowed written units. Before the sub-written unit variation introduced in Section 3.3, *Written units*, bowed written units first cause a following U.fina vowel form to change to O.fina.

3.2.3 *Hudum-specific shaping phase* 胡都木特有的变形阶段

The phase III consists a series of steps for Hudum-specific shaping requirements, and inside each step there may be more than one set of non-overlapping rules, each for a different group of letters.

Forms not captured by the predictive conditions introduced below are requested with FVSes. See Section 4.4, *Uncaptured forms*.

Phonetic and written unit-level shaping steps. A number of commonly recognized variation patterns are considered for the Hudum text presentation. Some patterns correspond to predictive shaping rules formalized as the 11 conditions below, while the rest are not executed in fonts. In the phonetic and written unit-level shaping steps (1-5), certain letters are analyzed as subjects and a set of contextual rules determine if the subject letters are in one of the 11 conditions.

Chachlag
Marked
Onset
Devsger
Chachlag_Onset
Masculine_Onset
Masculine_Devsger
Feminine
Particle
Dictionary_Particle
Post_Bowed

In_Isolation
Fallback

Contextual rules for each step and condition are specified in Section 4.2.1, *Orthographical shaping rules* (Table 6). The execution order of shaping steps is critical, as they have dependencies of previous ones, and a later decided condition overwrites an earlier one's output. For a reference list of particles that are commonly recognized as

NNBSP-applicable (step 3), see Appendix C. For predefined classes of phonetic letters that are used in contextual rules, see Table 7.

An alternative model for contextual rules. [Multi-pass analyses then single-pass substitution: advanced OTL rules (attested performance issue).]

[Compare with the current multi-pass substitutions: tightly coupled to rule order, maintenance of glyph classes of letters, ordinary OTL rules, arbitrary default positional forms and other intermediate forms, some conditions (e.g., onset) do not lead to forms different from the arbitrary default positional forms, arbitrary composition of conditions.]

C. The Twelve Syllabaries 十二字头

C.1 The First Syllabary 第一字头

	-a	-e	-i	-o/u	-ö/ü
Written forms of syllables: word-wise beginning (isolated or initial) vs. continuing (medial or final)					
Ø-	𐰇 𐰆 —	𐰇 𐰆 —	𐰇 𐰆 —	𐰇 𐰆 —	𐰇 𐰆 —
n-	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆
b-	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆
p-	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆
b-	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆
g-	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆
m-	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆
l-	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆
s-	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆
ś-	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆
t-	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆
d-	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆
ć-	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆
j-	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆
y-	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆
r-	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆
w-	𐰇 𐰆 𐰇 𐰆	𐰇 𐰆 𐰇 𐰆			

Notes:

1. The Twelve Syllabaries only exhibit an incomplete, core set of orthographical rules.
2. Only the 16 native onsets, seven native vowels, and 11 native devsgers are enumerated in the basic (C)V(C) syllable structure, making up the Twelve Syllabaries. This means not only loanword syllables, even native syllables as common as *sain* are out of scope.
3. Only the enumerated syllables' *marked* beginning forms and *unmarked* continuing forms are included, however they also have *unmarked* beginning forms (in particles) and *marked* continuing forms (in non-initial stems of compound words). Also, various marginal written forms are not included.

- Cells with light gray background contain only repeated written forms, which introduce great ambiguity into the orthography.
- [...]

C.2 Patterns exhibited in the First Syllabary

第一字头展现的模式

[...]

C.3 The other eleven syllabaries 其余十一字头

Syllabary	Devsgar	Written forms: medial & final		Examples shown with na/ne-	Note
2nd	-i				Medial form appends only a single I.medi after an existing I.medi .
3rd	-u/ü				
4th	-n				
5th	-ŋ				
6th	-b				
7th	-g				Masculine form for nucleus -a/i/o/u-.
					Feminine form for nucleus -e/i/ö/ü-.
8th	-m				
9th	-l				
10th	-s				
11th	-d				Alternative forms D.medi/fina are used for disambiguating.
12th	-r				

D. Orthography background 正字法背景

Words, particles, stems, etc. A whitespace-separated word (a morphological word) may be either an ordinary word (which can be modified by one or more particles as a host word) or a particle. Particles are prosodically part of their host words, forming a single prosodic word. Ordinary words contain one or more word stems, and may receive one or more suffixes.

Syllabic variations





















A written *syllable* in Hudum has the structure of C?V+C? (i.e., an optional *onset* consonant letter, one or more vowel letters forming the *nucleus*, and an optional *coda* consonant letter). The leading C?V part (a syllable excluding the trailing V?C?; i.e., the First Syllabary structure) has a special place in Hudum orthography patterns and is termed a *body*.




A consonant letter between two vowel letters belongs to the latter syllable, while a consonant letter excluded by the syllable structure is considered a *stray*. The stray forms (if distinct) are used as default positional forms.

Boundaries of written syllable always occur at morphological word boundaries (i.e., whitespaces) and between adjacent word stems (but not required between a word stem and a suffix). Note especially that syllable boundaries do not occur at cursive joining breaks inside a word, which are typically observed only before a chachlag.

Devsgers-specific forms. Traditionally, one of the eleven appendable phonetic letters from *the Second to Twelfth Syllabaries* (i.e., a postvocalic phonetic letter in a syllable) is recognized as a *devsger* (ᠳᠡᠪᠰᠢᠭᠡᠷ *debisger*; ДЭВСГЭᠷ *devsger*; “pad”). Consonant *n*, *g*, and *d* as well as vowel *i* have devsger-specific forms.

Table 9. The eleven devsgers

Type & phonetic letters	Written forms		Note	
Soft devsgers	<i>n</i>	 	A	A
	<i>m</i>	 	M	M
	<i>l</i>	 	L	L
	<i>η</i>	 	AG	AG
Hard devsgers	<i>b</i>	 	B	B
	<i>g</i>	 	H	H
		 	G	G
	<i>r</i>	 	R	R
	<i>s</i>	 	S	S
	<i>d</i>	 	ᠳ	ᠳ

Vocalic devsgers	<i>i</i>			II	I	Considered <i>y</i> or <i>yi.medi</i> / <i>y.fina</i> in alternative theories.
	<i>u/ü</i>			O	U	Considered <i>w</i> in alternative theories.

A medial devsger *i* takes its devsger-specific form II unless the preceding vowel's written form already ends with a written unit I. The two vocalic devsgers are not analyzed as *y* and *w*. *Predictive*

Phonetic Letters *h* and *y* are not eligible for a devsger position. *Predictive*

Variations in stem domains

A special *stem domain* is a word stem together with its following zero or more suffixes and particles. Certain variation patterns appears to be affected by boundaries of such scopes, however these scopes and boundaries are obscure to native users and are thus not fully utilized in encoding.

Onset placeholder. The syllable onset is required at the beginning of a stem domain, and a historical consonant letter aleph is used as the onset placeholder when an initial consonant is absent. This onset placeholder is generally considered by native users to be a part of the nucleus vowel letter's written form:

- It is a part of the default isolated and initial forms.
- It is considered to be a part of uncaptured forms on medial positions, which typically occur at beginning of non-first stems in a compound word. Such a form is unattested for the loanword vowel letter *e*.
- It is unattested on final positions although theoretically possible if the last stem of a compound word is a single vowel letter.

Initial-body forms of rounded vowels. Vowel letters *o*, *u*, *ö*, and *ü* are disambiguated when they are in the initial (i.e., not preceded by any other letter) body (C?V) of a stem domain:

- The default isolated and initial forms of *o*, *u*, *ö*, and *ü* are already initial-body forms.
- For the first stem in a word, initial-body medial and final forms (initial-body medial forms are only applicable to *ö* and *ü*) are predictable when one of these vowels follows an initial consonant. *Predictive*
- For the non-first stem in a compound word, medial and final initial-body forms are uncaptured forms.

Note that loanwords exhibit two major types of violation: Chinese loanwords tend to write a final u sound as U.fina even when it is in the initial body, and as Ü.fina after a bowed written unit; o-type sounds are generally normalized to other vowels. Other loanwords, where o sounds and u sounds are contrasted, tend to always write o as U.medi/O.fina and u as UI.medi/Ü.fina.

Gender-specific forms of h and g

Consonant letters h and g both have contrasted masculine and feminine classes of written forms, which have a distribution related to vowel genders (i.e., vowel harmony classes). Vowels a , o , and u are masculine; e , $é$, $ö$, and $ü$ are feminine; i is neuter.

In principle, a stem domain is internally gender harmonious, and thus should have a determinate gender value. However, acquiring a stem domain's actual harmony status is non-trivial as a stem can easily have mixed genders (e.g., loanwords), thus the gender-specific variation patterns of h and g are broken down into smaller patterns:

- A stray g takes the feminine form, which is handled by the default positional forms.
- An onset h or g agrees with its following vowel's gender (masculine or feminine); the neuter vowel i is treated like a feminine vowel for the onset. *Predictive*
- A devsgger g agrees with its preceding vowel's gender (masculine or feminine). A devsgger g that follows i takes the feminine form by default, and takes the masculine form if it remotely follows a masculine vowel and there is no feminine vowels in between. Note that this specification only considers single-directional, forward gender propagation for ig *Predictive*
- Other situations are considered uncaptured forms.

Chachlag

This non-joining written unit of a and e is usually observed after n , masculine h , masculine g , m , l , y , r , w , and only occasionally after j . The format control MVS is used to break cursive joining between the preceding consonant and a/e , as well as to request the special forms required for certain consonants:

- Letters n , h , and g take what appear to be their special, onset final forms. *Predictive with MVS*
- Letters y , w , and j take their under-differentiated forms (i.e., yodh, waw, and yodh, respectively). *Predictive with MVS*
- Writing in the form of chachlag is orthographically mandatory for a to appear after a word-medial h/g .
- Note that h and y normally do not have final forms as they are not used as devsggers.

E. Generally recognized particles

普遍承认的助词

Particle lists provided in the following documents and standards have been normalized and compared in Table 10:

- *The Users' Convention*, abbreviated as “UC” in the table, published as UNU/IIST Report No. 170 (Myatav Erdenechimeg et al. 1999) and MNS 4932: 2000.
- GB/T 25914-2010, “GB/T”.
- MGC/01-01 (2018, version 1.0.2), “EAC”.
 - Note that only the words that require special shaping are listed in the standard.
- L2/17-036 (Greg Eck et al. 2017), “17-036”.
- L2/18-293 (Badral Sanlig and Munkh-Uchral Enkhtur 2018), “18-293”.
 - Its listed particles are marked with white bullets in the table, because their proposed text representations are not explicitly given, and the document has inconsistencies in its Latin transliteration column.

Yellow and red highlighted are character sequences that would rely on NNBSF for achieving their special written forms. Red ones, in particular, involve lexical variations that are not predictable even if recognized as particles.

Table 10. Particles commonly recognized as NNBSF-applicable

Written form	Text representation following an NNBSF	UC	GB/T	EAC	17-036	18-293	Note
ᠶᠢᠨ	IIN	•	•	•	•	○	Genitive
ᠤᠨ	OA	•	•	•	•	○	
ᠤ	U	•	•	•	•	○	
ᠳᠤ	DU	•	•	•	•	○	Dative
ᠲᠤ	TU	•	•	•	•	○	
ᠳᠣᠷ	DOR	•	•	•	•	○	
ᠲᠦᠷ	TOR	•	•	•	•	○	
ᠠ	Á	•	•	•	•	○	Exclamatory
ᠠ	A	•	•	•	•	○	
ᠶᠢ	II	•	•	•	•	○	Accusative

ᠶ	I	<i>i</i>		•	•	•	•	○	Accusative
ᠠᠴᠠ	AČA	<i>ača</i>	<i>eče</i>	•	•	<i>ača</i>	•	○	Ablative
ᠶᠠᠷ	IAR	<i>iyar</i>	<i>iyer</i>	•	•	•	•	○	Instrumental
ᠪᠠᠷ	BAR	<i>bar</i>	<i>ber</i>	•	•		•	○	
ᠲᠠᠢ	TAI	<i>tai</i>	<i>tei</i>	•	•		•	○	Comitative
ᠯᠣᠭᠠ	LOG Á	<i>lug-a</i>		•	•		•	○	
ᠯᠦᠭᠡ	LOGA		<i>lüge</i>	•	•	•	•	○	Reflexive
ᠪᠠᠭ	BAA	<i>ban</i>	<i>ben</i>	•	•		•	○	
ᠶᠠᠨ	IAA	<i>iyān</i>	<i>iyen</i>		•	•	•	○	Reflexive genitive
ᠨᠠᠢᠬᠠᠭ	NAIIHAA	<i>naihan</i>							
ᠨᠠᠢᠭᠠᠭ	NAIIGAA		<i>neihen</i>						
ᠤᠪᠠᠭ	OBAA	<i>uban</i>	<i>üben</i>					○	
ᠳᠠᠭᠠᠭ	DAĞAA	<i>dagan</i>		rendering only	•	•	•	○	Reflexive dative
ᠳᠠᠭᠠᠭ	DAGAA		<i>degen</i>	encoding only	•	•	•	○	
ᠲᠠᠭᠠᠭ	TAĞAA	<i>tagan</i>			•		•	○	
ᠲᠠᠭᠠᠭ	TAGAA		<i>tegen</i>	•	•		•	○	
ᠳᠤᠷᠢᠶᠠᠭ	DORIYAA	<i>duriyan</i>	<i>düriyen</i>					○	Reflexive accusative
ᠶᠣᠭᠠᠭ	YOĞAA	<i>yugan</i>			•		•	○	
ᠶᠣᠭᠠᠭ	YOGAA		<i>yügen</i>		•	•	•	○	Reflexive ablative
ᠠᠴᠠᠭᠠᠭ	AČAĞAA	<i>ačagan</i>			•	•	•	○	
ᠠᠴᠠᠭᠠᠭ	AČAGAA		<i>ečegen</i>		•		•	○	Reflexive comitative
ᠲᠠᠢᠭᠠᠭ	TAIIĞAA	<i>taigan</i>					<i>tayigan</i>	○	
ᠲᠠᠢᠭᠠᠭ	TAIIGAA		<i>teigen</i>				<i>teyigen</i>	○	Directive
ᠠᠣᠷᠣᠭᠤ	AOROĞU	<i>urugu</i>					•/not	○	
ᠮᠢᠨᠢ	MINI	<i>mini</i>						○	Possessive
ᠴᠢᠨᠢ	ČINI	<i>čini</i>						○	
ᠮᠠᠨᠢ	MANI	<i>mani</i>	<i>meni</i>					○	
ᠲᠠᠨᠢ	TANI	<i>tani</i>	<i>teni</i>					○	
ᠠᠨᠤ	AANU	<i>anu</i>						○	

ᐱᐢᐅᐅ	AINU	<i>inü</i>			○	
ᐢᐅ	NI	<i>ni</i>			○	
ᐢᐅᐢ	HINI	<i>bini</i>				
ᐅᐱᐢ	DONI	<i>duni</i>	<i>düni</i>	•		Possessive dative
ᐅᐱᐢᐅ	TONI	<i>tuni</i>	<i>tüni</i>	<i>tüni</i>		
ᐢᐅ	HI	<i>hi</i>			•	Empty noun
ᐢᐅᐢ	HIA	<i>bin</i>			•	
ᐅᐱᐢ	DAHI	<i>dahi</i>	<i>debi</i>	•	•	Empty noun dative?
ᐅᐱᐢᐅ	TAHI	<i>tabi</i>	<i>tebi</i>			
ᐢᐅᐱᐅᐅᐅ	NOĠOĐ	<i>nugud</i>			• ○	Plural
ᐢᐅᐱᐅᐅᐅᐅ	NOGOĐ		<i>nügüd</i>	•	• ○	
ᐅᐅ	OĐ	<i>ud</i>	<i>üd</i>	• • •	• ○	
ᐢᐅᐢ	NAR	<i>nar</i>	<i>ner</i>	<i>nar</i> •	• ○	
ᐅᐱᐢᐅᐢ	DOĠAR	<i>dugar</i>		•	•	Ordinal
ᐅᐱᐢᐅᐢᐅ	DOGAR		<i>düger</i>	•	•	
ᐢᐅᐢ	SIH	<i>sig</i>				Like
ᐢᐅᐢᐅ	SIG	<i>sig</i>				
ᐢᐅᐢᐅ	ĆIH	<i>ćig</i>				Even
ᐢᐅᐢᐅᐅ	ĆIG	<i>ćig</i>				
ᐢᐅᐢ	ĆU	<i>ću</i>	<i>ćü</i>	•		
ᐢᐅᐢ	LA	<i>la</i>	<i>le</i>			Just
ᐅᐱᐢ	DAH	<i>dag</i>		•	•/not	?
ᐅᐱᐢᐅ	DAG		<i>deg</i>	•	•/not	
ᐅᐢᐅ	TA	<i>ta</i>	<i>te</i>			?
ᐢᐅᐢᐅ	SAA	<i>san</i>	<i>sen</i>			?
ᐅᐅ	OU	<i>uu</i>	<i>üü</i>	not	not	Interrogative
ᐅᐢ	DA	<i>da</i>	<i>de</i>	•		Modal particle?
ᐢᐅᐢᐅᐢᐅ	AOIGAI	<i>ügei</i>		•	•/not ○	Negative
ᐢᐅᐢᐅᐢᐅᐅ	NOGAA		<i>nügen</i>	•	• ○	

𐄎𐄌	YOM		<i>yüm</i>	•		
𐄎𐄌𐄎	YOMSAA		<i>yümsen</i>	•		
𐄎	GO		<i>bü</i>	•		
𐄎𐄎	BOU		<i>büü</i>	not		

Copyright © 2020 Liang Hai and Unicode, Inc. All Rights Reserved. The Unicode Consortium and Liang Hai make no expressed or implied warranty of any kind, and assume no liability for errors or omissions. No liability is assumed for incidental and consequential damages in connection with or arising out of the use of the information or programs contained or accompanying this technical note. The Unicode [Terms of Use](#) apply.

Unicode and the Unicode logo are trademarks of Unicode, Inc., and are registered in some jurisdictions.