

An improved graphetic model for the Mongolian encoding

改进的蒙古文编码字形模型

Liang Hai (梁海)
lianghai@gmail.com

3 April 2018
2018年4月3日

1 Introduction

1.1 Related earlier documents

- [L2/17-328](#): *Script Ad Hoc Group Recommendations on Mongolian Text Model*, 9 August 2017.
- [L2/17-335](#): *A graphetic approach for the Mongolian encoding model (draft)*, 31 August 2017.
- [L2/17-334](#): *On migration issues of the graphetic model*, 18 September 2017.
- [L2/17-347](#): *Mongolian Ad Hoc Report (Hohhot, Inner Mongolia)*, 29 September 2017.

1.2 Acknowledgements

The author received feedback from Shen Yilei (沈逸磊) and the Script Ad Hoc group. The prototype Mongolian font is based on Bolorsoft's open-source font MongolianScript.

1.3 Scope

This document discusses¹ a crucial subset of the Mongolian script's usage, the so-called *Hudum* (ᠬᠣᠳᠤᠮᠤ or худам) writing system of the modern Mongolian language. Various other writing systems and historical stages (Hudum Ali Gali, historical Hudum, Todo, Todo Ali Gali, Manchu, Manchu Ali Gali, historical Manchu, Sibe, etc) are excluded.

¹For a precise comparison with the current encoding in the Unicode Standard, note this document covers the equivalence of letters U+1820..U+1842 and format control characters U+180A..U+180E, in terms of representation capability.

Table 2: (continued)

	-a	-e	-i	-o	-u	-ö	-ü
t-	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ
d-	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ
ć-	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ
j-	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ
γ-	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ
r-	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ
w-	ᵀ ᵀ ᵀ ᵀ	ᵀ ᵀ ᵀ ᵀ					

From the First Syllabary it's already observable that encoding the phonetic letters as characters (the *phonetic* approach) leads to has serious inherent defects:

1. A large number of combinations are confusable in written forms, and such ambiguity means low reliability. For example, text inputting and editing are highly prone to visually invisible errors. Notable confusable groups include *a/e, o/u/ö/ü* (in particular, *o* is always written the same as *u*, and *ö* the same as *ü*), *x/g, t/d*, etc.
2. Even though complex contextual rules can decide most written forms, additional information is required when an expected form is grammatical or lexical thus is not predictable with only a context of phonetic letters. This uncertainty is underrepresented in the Twelve Syllabaries (where the uncertainty is only observable in *a/e* and *d*), but is common in actual words. Thus a phonetic encoding requires inserting arbitrary format control characters to select writing forms, leading to low usability for users.

1.5 Graphetic approach

In light of the unsuitability of the phonetic approach, a group of experts have been exploring a *graphetic*³ approach that encodes cursive joining graphemes instead of phonetic letters, still utilizing the Unicode–OpenType cursive joining model, but avoiding the aforementioned ambiguity and uncertainty inherent in the phonetic approach. After analyzing written forms of phonetic letters according to their alternation between different cursive joining positions, revealed cursive joining graphemes are further decomposed and unified to avoid highly confusable characters. Such a relatively simple encoding model provides solid support for higher-level text processes.

Taking the graphetic approach, a specific encoding model has been designed for discussing and testing, and is introduced in this document.

³The rather special term *graphetic* is used to represent a concept opposed to *phonetic*. It means “relating to graphemes or graphs”. Just like *phonetic* is often preferred in text encoding discussions as it's a broader concept than *phonemic* (strictly relating to phonemes), the broader wording *graphetic* is here chosen over *graphemic* (strictly relating to graphemes).

2 Character repertoire

The model requires a repertoire of 30 characters. See *Table 3* for a provisional list of representative glyphs and character names. Note representative glyphs are not relevant to actual shaping. Slashes (/) instead of the word “OR” are used in provisional names for better readability. Transliterations of phonetically definite characters are simply borrowed from their corresponding phonetic letters (note these transliterations are all in lowercase, and colored in blue to distinguish from phonetic transliterations), while phonetically ambiguous characters have distinct transliterations in uppercase letters or other symbols.

Table 3: Character repertoire

No.	Repr.	Name	Translit.
1.	•	MONGOLIAN CHARACTER NIRUGU	-
2.	ᠠ	MONGOLIAN CHARACTER ALEPH/A/E/NA	•
3.	ᠡ	MONGOLIAN CHARACTER A/E	A
4.	ᠢ	MONGOLIAN CHARACTER I/JA/YA	I
5.	ᠣ	MONGOLIAN CHARACTER O/U/OE/UE	U
6.	ᠤ	MONGOLIAN CHARACTER O/U/OE/UE/WA	Ú
7.	ᠥ	MONGOLIAN CHARACTER OE/UE	Û
8.	ᠨ	MONGOLIAN CHARACTER NA	<i>n</i>
9.	ᠪ	MONGOLIAN CHARACTER BA	<i>b</i>
10.	ᠫ	MONGOLIAN CHARACTER PA	<i>p</i>
11.	ᠬ	MONGOLIAN CHARACTER XA/GA	X
12.	ᠭ	MONGOLIAN CHARACTER XE/GE	G
13.	ᠭᠦ	MONGOLIAN CHARACTER GA	<i>g</i>
14.	ᠮ	MONGOLIAN CHARACTER MA	<i>m</i>
15.	ᠯ	MONGOLIAN CHARACTER LA	<i>l</i>
16.	ᠰ	MONGOLIAN CHARACTER SA	<i>s</i>
17.	ᠱ	MONGOLIAN CHARACTER SHA	<i>ś</i>
18.	ᠲ	MONGOLIAN CHARACTER TA/DA	T
19.	ᠳ	MONGOLIAN CHARACTER DA/TA	D
20.	ᠴ	MONGOLIAN CHARACTER CHA	<i>č</i>
21.	ᠵ	MONGOLIAN CHARACTER JA	<i>j</i>
22.	ᠶ	MONGOLIAN CHARACTER YA	<i>y</i>
23.	ᠷ	MONGOLIAN CHARACTER RA	<i>r</i>
24.	ᠰ	MONGOLIAN CHARACTER WA/EE	W
25.	ᠰ	MONGOLIAN CHARACTER FA	<i>f</i>

Table 3: (continued)

No.	Repr.	Name	Translit.
26.	᠎	MONGOLIAN CHARACTER KA	<i>k</i>
27.	᠎ᠠ	MONGOLIAN CHARACTER CA	<i>c</i>
28.	᠎ᠢ	MONGOLIAN CHARACTER ZA	<i>z</i>
29.	᠎ᠡ	MONGOLIAN CHARACTER HA/ZHI	<i>H</i>
30.	᠎ᠢᠨ	MONGOLIAN CHARACTER RHA	<i>ř</i>

3 Text representation

The model requires two stages of contextual shaping to select the orthographically correct form for every character: first the *cursive joining* stage (Section 3.1), then the *round consonant* stage (Section 3.2).

Note certain written forms of phonetic letters are encoded as character sequences:

1. The *aleph*-ed (prepended with an extra cap or tooth) forms of vowel letters (*a*: ᠠ, ᠡ, ᠢ; *i*: ᠢᠨ, ᠡᠨ, ᠢᠨ; etc) have their aleph encoded as a character: ᠠ ALEPH/A/E/NA.
2. The *long tooth*-ed forms of ö/ü (ᠠᠨ, ᠡᠨ, ᠢᠨ) have their long tooth encoded as a character: ᠠ I/JA/YA.
3. *η* (ᠠᠨ, ᠡᠨ) is encoded as a sequence: < ᠠ ALEPH/A/E/NA, ᠨ XE/GE >.
4. The syllable coda form of *d* (ᠠᠨ, ᠡᠨ) is encoded as a sequence: < ᠠ O/U/OE/UE, ᠠ ALEPH/A/E/NA >.
5. The word-initial form of *h* (ᠠᠨ) is encoded as a sequence: < ᠠ ALEPH/A/E/NA, ᠨ HA/ZHI >.
6. *ł* (ᠠᠨ, ᠡᠨ) is encoded as a sequence: < ᠠ LA, ᠨ HA/ZHI >.
7. *č* (ᠠᠨ) is encoded as a sequence: < ᠠ O/U/OE/UE, ᠠ O/U/OE/UE >.

3.1 Cursive joining stage

According to joining types of adjacent characters, the appropriate positional form is selected for every character. Note cursive joining positional forms are irrelevant to the grammatical definition of “word”.

See Table 4, Table 5, and Table 6 for positional forms of characters of the three joining types.⁴ Joined ends of positional forms are emphasized with an extra stroke in blue (•). Positions with parentheses are unattested. A couple of recommended fallbacks are in the parentheses, but further discussion is necessary for the unattested isolate positions, which are the most problematic and are all left empty for now.

⁴The standard terminology and notation in the Unicode Standard for joining types and positional forms are avoided here, because they are confusing for vertically written scripts. For the background, see “Cursive Joining” in Section 14.4, *Phags-pa, The Unicode Standard, Version 10.0*.

Table 4: Dual-joining characters

Character	Translit.	Isolate	Initial	Medial	Final
MONGOLIAN CHARACTER NIRUGU	-	.	-	-	-
MONGOLIAN CHARACTER ALEPH/A/E/NA	.	ᠠ	ᠠ	ᠠ	ᠠ
MONGOLIAN CHARACTER I/JA/YA	<i>I</i>	ᠢ	ᠢ	ᠢ	ᠢ
MONGOLIAN CHARACTER O/U/OE/UE	<i>U</i>	()	ᠣ	ᠣ	ᠣ
MONGOLIAN CHARACTER NA	<i>n</i>	()	ᠨ	ᠨ	ᠨ
MONGOLIAN CHARACTER BA	<i>b</i>	()	ᠪ	ᠪ	ᠪ
MONGOLIAN CHARACTER PA	<i>p</i>	()	ᠫ	ᠫ	ᠫ
MONGOLIAN CHARACTER XA/GA	<i>X</i>	()	ᠬ	ᠬ	ᠬ
MONGOLIAN CHARACTER XE/GE	<i>G</i>	()	ᠭ	ᠭ	ᠭ
MONGOLIAN CHARACTER GA	<i>g</i>	()	ᠭ	ᠭ	ᠭ
MONGOLIAN CHARACTER MA	<i>m</i>	()	ᠮ	ᠮ	ᠮ
MONGOLIAN CHARACTER LA	<i>l</i>	()	ᠯ	ᠯ	ᠯ
MONGOLIAN CHARACTER SA	<i>s</i>	()	ᠰ	ᠰ	ᠰ
MONGOLIAN CHARACTER SHA	<i>ś</i>	()	ᠰ	ᠰ	ᠰ
MONGOLIAN CHARACTER TA/DA	<i>T</i>	()	ᠲ	ᠲ	ᠲ
MONGOLIAN CHARACTER DA/TA	<i>D</i>	()	ᠲ	ᠲ	ᠲ
MONGOLIAN CHARACTER CHA	<i>č</i>	()	ᠴ	ᠴ	ᠴ
MONGOLIAN CHARACTER JA	<i>j</i>	()	(ᠵ)	ᠵ	ᠵ
MONGOLIAN CHARACTER YA	<i>y</i>	()	ᠶ	ᠶ	(ᠶ)
MONGOLIAN CHARACTER RA	<i>r</i>	()	ᠷ	ᠷ	ᠷ
MONGOLIAN CHARACTER WA/EE	<i>W</i>	()	ᠸ	ᠸ	ᠸ
MONGOLIAN CHARACTER FA	<i>f</i>	()	ᠸ	ᠸ	ᠸ
MONGOLIAN CHARACTER KA	<i>k</i>	()	ᠶ	ᠶ	ᠶ
MONGOLIAN CHARACTER CA	<i>c</i>	()	ᠶ	ᠶ	ᠶ
MONGOLIAN CHARACTER ZA	<i>z</i>	()	ᠶ	ᠶ	ᠶ
MONGOLIAN CHARACTER HA/ZHI	<i>H</i>	()	ᠬ	ᠬ	ᠬ
MONGOLIAN CHARACTER RHA	<i>ř</i>	()	ᠬ	ᠬ	(ᠬ)

Table 5: Top-joining characters

Character	Translit.	Isolate	Final
MONGOLIAN CHARACTER O/U/OE/UE/WA	ᠣ	ᠤ	ᠶ
MONGOLIAN CHARACTER OE/UE	ᠥ	()	ᠦ

Table 6: Non-joining characters

Character	Translit.	Isolate
MONGOLIAN CHARACTER A/E	A	ᠠ

See Table 7 for examples of text representation that only depend on this stage of contextual shaping. Both phonetic (including notable alternatives) and graphetic transliterations are provided.

Table 7: Examples

Phonetic	Output	Graphetic	Character sequence
<i>a/en</i>	ᠠᠨ	..	< ᠠᠨ ALEPH/A/E/NA, ᠠᠨ ALEPH/A/E/NA >
<i>e/a</i>	ᠠ	.	< ᠠᠨ ALEPH/A/E/NA >
<i>a/e</i>	ᠠ	A	< ᠠᠨ A/E >
<i>i/ei</i>	ᠢᠨ	·I	< ᠠᠨ ALEPH/A/E/NA, ᠢᠨ I/JA/YA >
<i>i</i>	ᠢ	I	< ᠢᠨ I/JA/YA >
<i>o/u</i>	ᠣᠨ	·U	< ᠠᠨ ALEPH/A/E/NA, ᠣᠨ O/U/OE/UE >
<i>ö/ü</i>	ᠣᠨ	·Ü	< ᠠᠨ ALEPH/A/E/NA, ᠣᠨ OE/UE >
<i>ü/eü</i>	ᠣᠨ	·Ú	< ᠠᠨ ALEPH/A/E/NA, ᠣᠨ O/U/OE/UE/WA >
<i>u/ü</i>	ᠣ	Ú	< ᠣᠨ O/U/OE/UE/WA >
<i>ordo/urtu</i>	ᠣᠷᠳᠣᠤᠷᠲᠤ	·URDÚ	< ᠠᠨ ALEPH/A/E/NA, ᠣᠨ O/U/OE/UE, ᠷᠠ RA, ᠳᠠ DA/TA, ᠣᠨ O/U/OE/UE/WA >
<i>ada/ende</i>	ᠠᠳᠠᠨᠡ	..D·	< ᠠᠨ ALEPH/A/E/NA, ᠠᠨ ALEPH/A/E/NA, ᠳᠠ DA/TA, ᠠᠨ ALEPH/A/E/NA >
<i>xana</i>	ᠬᠠᠨᠠ	X·n·	< ᠬᠠ XA/GA, ᠠᠨ ALEPH/A/E/NA, ᠨᠠ NA, ᠠᠨ ALEPH/A/E/NA >
<i>xana</i>	ᠬᠠᠨᠠ	X·nA	< ᠬᠠ XA/GA, ᠠᠨ ALEPH/A/E/NA, ᠨᠠ NA, ᠠᠨ A/E >
<i>čerig</i>	ᠴᠡᠷᠢᠭ	ć·RIG	< ᠴᠢ CHA, ᠠᠨ ALEPH/A/E/NA, ᠷᠠ RA, ᠢᠨ I/JA/YA, ᠷᠢ XE/GE >
<i>jarlig</i>	ᠵᠠᠷᠯᠢᠭ	I·rlIX	< ᠢᠨ I/JA/YA, ᠠᠨ ALEPH/A/E/NA, ᠷᠠ RA, ᠯᠠ LA, ᠢᠨ I/JA/YA, ᠬᠠ XA/GA >
<i>altan odo</i>	ᠠᠯᠲᠠᠨ ᠣᠳᠣ	..lD·ᠤ·UDÚ	< ᠠᠨ ALEPH/A/E/NA, ᠠᠨ ALEPH/A/E/NA, ᠯᠠ LA, ᠳᠠ DA/TA, ᠠᠨ ALEPH/A/E/NA, ᠠᠨ ALEPH/A/E/NA, ᠤ SPACE, ᠠᠨ ALEPH/A/E/NA, ᠣᠨ O/U/OE/UE, ᠳᠠ DA/TA, ᠣᠨ O/U/OE/UE/WA >
<i>altanodo</i>	ᠠᠯᠲᠠᠨᠣᠳᠣ	..lD...UDÚ	< ᠠᠨ ALEPH/A/E/NA, ᠠᠨ ALEPH/A/E/NA, ᠯᠠ LA, ᠳᠠ DA/TA, ᠠᠨ ALEPH/A/E/NA, ᠠᠨ ALEPH/A/E/NA, ᠠᠨ ALEPH/A/E/NA, ᠣᠨ O/U/OE/UE, ᠳᠠ DA/TA, ᠣᠨ O/U/OE/UE/WA >
<i>bilig batu</i>	ᠪᠢᠯᠢᠭ ᠪᠠᠲᠤ	bilIGᠤ·b·DÚ	< ᠪᠠ BA, ᠢᠨ I/JA/YA, ᠯᠠ LA, ᠢᠨ I/JA/YA, ᠷᠢ XE/GE, ᠤ SPACE, ᠪᠠ BA, ᠠᠨ ALEPH/A/E/NA, ᠳᠠ DA/TA, ᠣᠨ O/U/OE/UE/WA >

Table 7: (continued)

Phonetic	Output	Graphetic	Character sequence
<i>biligbatu</i>		<i>bilIGb·DÚ</i>	< ᠪ BA, ᠢ I/JA/YA, ᠯ LA, ᠢ I/JA/YA, ᠭ XE/GE, ᠪ BA, ᠠ ALEPH/A/E/NA, ᠲ DA/TA, ᠤ O/U/OE/UE/WA >
<i>ed</i>		<i>·D</i>	< ᠠ ALEPH/A/E/NA, ᠲ DA/TA >
<i>gxir/xxir</i>		<i>GGIR</i>	< ᠭ XE/GE, ᠭ XE/GE, ᠢ I/JA/YA, ᠷ RA >
<i>radio/radiu</i>		<i>r·DIU</i>	< ᠷ RA, ᠠ ALEPH/A/E/NA, ᠲ DA/TA, ᠢ I/JA/YA, ᠤ O/U/OE/UE >
<i>ᠲᠢᠷᠢ/ᠲᠡᠭᠢ</i>		<i>T·GᠷI</i>	< ᠲ TA/DA, ᠠ ALEPH/A/E/NA, ᠭ XE/GE, ᠷ RA, ᠢ I/JA/YA >
<i>naima</i>		<i>n·Im·</i>	< ᠨ NA, ᠠ ALEPH/A/E/NA, ᠢ I/JA/YA, ᠮ MA, ᠠ ALEPH/A/E/NA >
<i>sayixan</i>		<i>s·yIX·</i>	< ᠰ SA, ᠠ ALEPH/A/E/NA, ᠶ YA, ᠢ I/JA/YA, ᠬ XA/GA, ᠠ ALEPH/A/E/NA, ᠠ ALEPH/A/E/NA >
<i>sain/sayin</i>		<i>s·II·</i>	< ᠰ SA, ᠠ ALEPH/A/E/NA, ᠢ I/JA/YA, ᠢ I/JA/YA, ᠠ ALEPH/A/E/NA >
<i>uu/üü</i>		<i>UÚ</i>	< ᠤ O/U/OE/UE, ᠤ O/U/OE/UE/WA >

For controlling cursive joining, although the general characters U+200C ZERO WIDTH NON-JOINER and U+200D ZERO WIDTH JOINER can be used, in day-to-day usage it is recommended to produce joined forms in isolation using · MONGOLIAN CHARACTER NIRUGU, which as a visible character is easier for users to manipulate. See Table 8 for examples.

Table 8: Examples

Phonetic	Output	Graphetic	Character sequence
<i>ga</i>		<i>g·-</i>	< ᠭ GA, ᠠ ALEPH/A/E/NA, · NIRUGU >
<i>ga</i>		<i>-g·-</i>	< · NIRUGU, ᠭ GA, ᠠ ALEPH/A/E/NA, · NIRUGU >
<i>ga</i>		<i>-gA</i>	< · NIRUGU, ᠭ GA, ᠠ A/E >
<i>ba</i>		<i>b·-</i>	< ᠪ BA, ᠠ ALEPH/A/E/NA, · NIRUGU >
<i>na</i>		<i>-nA</i>	< · NIRUGU, ᠨ NA, ᠠ A/E >

3.2 Round consonant stage

Certain combinations of positional forms that involve the five so-called round consonants (ᠪ BA, ᠫ PA, ᠭ XE/GE, ᠮ FA, ᠬ KA) require special treatment. The rules for selecting the required variants forms are formally defined below. See Table 9 for explanation of the notation.

- Rule 1 Cr + Va → Cr + Va'
- Rule 2 Cr + Vu → Cr' + Vu

Table 9: Notation

Notation	Positional forms	Note
Cr	𐤁 𐤂 𐤃 𐤄 𐤅 𐤆 𐤇 𐤈 𐤉 𐤊 𐤋 𐤌	Initial and medial forms of round consonants: BA, PA, XE/GE, FA, KA.
Cr'	𐤁𐤀 𐤂𐤀 𐤃𐤀 𐤄𐤀 𐤅𐤀 𐤆𐤀 𐤇𐤀 𐤈𐤀 𐤉𐤀 𐤊𐤀 𐤋𐤀 𐤌𐤀	Corresponding variant forms of Cr.
Va	𐤀 𐤁	Final forms of certain vowels: ALEPH/A/E/NA, I/JA/YA.
Va'	𐤁𐤀 𐤂𐤀	Corresponding variant forms of Va.
Vu	𐤀𐤁 𐤀𐤂 𐤀𐤃	Medial and final forms of certain vowels: O/U/OE/UE, OE/UE.

See Table 10 for examples of affected text representation cases. A column of intermediate glyphs (the output positional forms from the last stage) is provided for comparison.

Table 10: Examples

Phonetic	Interm.	→	Output	Graphetic	Rule	Character sequence
<i>ba/be</i>	* 𐤁	→	𐤁	<i>b</i>	1	< 𐤁 ^{BA} , 𐤀 ^{ALEPH/A/E/NA} >
<i>bi</i>	* 𐤁	→	𐤁	<i>bI</i>	1	< 𐤁 ^{BA} , 𐤀 ^{I/JA/YA} >
<i>bo/bu/bö/bü</i>	* 𐤁	→	𐤁	<i>bU</i>	2	< 𐤁 ^{BA} , 𐤀 ^{O/U/OE/UE} >
<i>bö/bü</i>	* 𐤁	→	𐤁	<i>bÜ</i>	2	< 𐤁 ^{BA} , 𐤀 ^{OE/UE} >

Only the shape changes relevant to encoding are specified here. In actual fonts, depending on the specific style, glyphs (including the glyphs discussed above) might require more typographical adjustments that do not affect encoding. In particular, Vu glyphs usually need a change of connection point in order to properly join to Cr' glyphs, although the resulted glyph changes can be subtle.

Other optional variant forms, such as 𐤁𐤀 → 𐤁^{bl}, are also commonly seen in fonts.

4 Adjustments to the character repertoire

The character repertoire presented in Table 3 is a result of maximizing the “low ambiguity” advantage of the graphetic approach, therefore most confusable graphemes are unified or decomposed. Considering the benefits of recording more information of phonetic letters,⁵ it might be desirable to moderately disunify or recompose some graphemes while tolerating more ambiguity. Some candidates are shown in Table 11.

Table 11: Candidates for disunification

From	Disunify/recompose	Note
𐤀 ^{ALEPH/A/E/NA}	ALEPH	Only initial and medial forms are attested.
𐤀 ^{I/JA/YA}	LONG TOOTH OF OE/UE	Change joining type of 𐤀 ^{OE/UE} to dual-joining to have a medial form that has the disunified as a part.

⁵The key benefit is easier syllabification, which is helpful for collation, line-breaking, typographical preferences, etc.

Table 11: (continued)

From	Disunify/recompose	Note
ᠪ WA/EE	EE	Only medial and final forms are attested.
ᠠ ALEPH/A/E/NA	CODA NA	Only medial and final forms are attested.
< ᠪ O/U/OE/UE, ᠠ ALEPH/A/E/NA >	CODA DA	Only medial and final forms are attested.
ᠢ I/JA/YA	INITIAL JA	The disunified becomes initial ᠢ JA.
< ᠠ ALEPH/A/E/NA, ᠬ HA/ZHI >	INITIAL HA	The recomposed becomes initial ᠬ HA/ZHI, pushing the existing intial form to be disunified.
ᠬ HA/ZHI	ZHI	Only initial form is attested.
< ᠪ O/U/OE/UE, ᠬ O/U/OE/UE >	CHI	Only initial form is attested.

Note although the medial forms ᠬ XA/GA and ᠬ GA are theoretically confusable with < ᠠ ALEPH/A/E/NA, ᠠ ALEPH/A/E/NA > and < ᠠ NA, ᠠ NA >, respectively, it is not beneficial to decompose them, because they are well distinguished by users (especially since the related intial and final forms are distinct) and decomposing only the medial forms would introduce problematic medial-less dual-joining characters.

5 Text processes

Text encodings play a fundamental role in the life cycle of digitalized text: users → *input* (keyboards, OCR, speech-to-text, etc) → Unicode encoding → *output* (display, collation, text-to-speech, etc) → users. The encoding model's implications on various text processes are discussed here.

5.1 Output: collation

When idealistic phonetic collation is expected, for example in dictionaries, phonetic information should be carried as metadata.

Otherwise, in day-to-day use cases where text must be collated automatically without metadata, a folded collation (certain homographic written forms are folded into a single collation element) is a reliable solution. For an introduction, see Shen Yilei's document [MWG/2-N5 G2P Sorting: An Automated Natural Sorting Method for Graphetically Encoded Mongolian](#).

5.2 Input: keyboards and input methods

The preference of inputting phonetically should be fulfilled in the keyboard and input methods for average users. See *Section 7, Input* in [L2/17-334 On migration issues of the graphetic model](#) for a full discussion.

In order to prevent confusing fallbacks on undefined cursive joining positions (for example, if an undefined isolate BA were shown with its intial form ᠪ, it would be confused with ᠪ, an isolate O/U/OE/UE/WA), unattested positions must either have distinct and reasonable fallback forms or be explicitly marked as invalid.

Table 4 shows some recommended fallback forms in parentheses, but unattested isolate positions are especially problematic and need discussion. To explicitly mark invalid positions, visual aids similar to Indic scripts’ dotted circles (◌̣) can be considered, for example, arrows pointing to the direction where a character is required. See Table 12 for how arrows can improve user experience.

Table 12: Input experience

Keystroke	Output	Graphetic	Note
N	ጥ→	<i>n</i>	ጥ NA doesn’t have a valid initial form.
A	ጥ	<i>n·</i>	
I	ጥገ	<i>n·I</i>	
M	ጥገጌ	<i>n·Im</i>	
A	ጥገጌ	<i>n·Im·</i>	
Shift-N	ጥገጌ	<i>n·Im··</i>	
[Space]	ጥገጌ	<i>n·Im··</i> ◻	
Shift-U	ጥገጌ ፀ	<i>n·Im··</i> ◻ ሆ	
[Space]	ጥገጌ ፀ	<i>n·Im··</i> ◻ ሆ ◻	
N	ጥገጌ ፀ ጥ→	<i>n·Im··</i> ◻ ሆ ◻ <i>n</i>	ጥ NA doesn’t have a valid initial form.
I	ጥገጌ ፀ ጥገ	<i>n·Im··</i> ◻ ሆ ◻ <i>nI</i>	
Shift-G	ጥገጌ ፀ ጥገጌ	<i>n·Im··</i> ◻ ሆ ◻ <i>nIG</i>	
E	ጥገጌ ፀ ጥገጌ	<i>n·Im··</i> ◻ ሆ ◻ <i>nIG·</i>	

To avoid artificial visual aids, thus providing a more natural user experience, the output can be always appended with a · NIRUGU until an “end of word” signal (either a special keystroke or a character that is not bottom-joining) is detected. See Table 13.

Table 13: Alternative input experience

Keystroke	Output	Graphetic	Note
N	ጥ	<i>n-</i>	
A	ጥ	<i>n·-</i>	
I	ጥገ	<i>n·I-</i>	
M	ጥገጌ	<i>n·Im-</i>	
A	ጥገጌ	<i>n·Im·-</i>	
Shift-N	ጥገጌ	<i>n·Im··-</i>	
[End of Word]	ጥገጌ	<i>n·Im··</i> ◻	· NIRUGU replaced by a ◻ SPACE.
Shift-U	ጥገጌ ፀ	<i>n·Im··</i> ◻ ሆ ◻	[End of Word] implied by characters that are not bottom-joining.

