

Монгол бичгийн кодлох сайжруулсан загвар

An improved graphetic model for the Mongolian encoding

Лианг Хай 梁海, цахим хаяг: lianghai@gmail.com

2018 оны 3 сарын 24 өдөр

24March2018

1. Танилцуулага

Introduction

1.1 Өмнөх холбоотой баримтууд

Related earlier documents

- L2/17-328: 2017 оны 8-р сарын 9 өдрийн Монгол бичгийг дүрсээр кодлох загвар.
L2/17-328: ScriptAdHocGroupRecommendationsonMongolianTextModel,9August2017.
- L2/17-335: 2017 оны 8-р сарын 31-ны өдрийн Монгол бичгийг дүрсээр кодлох загвар (Ноорог буюу санал).
L2/17-335: AgrapheticapproachfortheMongolianencodingmodel(draft),31August2017.
- L2/17-334: Дүрсээр кодлоход үүссэн асуудалууд.
L2/17-334: Onmigrationissuesofthegraphicmodel,18September2017.
- L2/17-347: 2017 оны 9-р сарын 29 өдрийн Монгол бичгийн ажлын хэсгийн уулзалтын тайлан.
L2/17-347: MongolianAdHocReport(Hohhot,InnerMongolia),29September2017.

1.2 Хүлээн зөвшөөрсөн

Acknowledgements

Энэхүү баримт бичгийг боловсруулах явцад зохиогч Монгол бичгийн ажлын хэсгийн санал хүсэлтийг хүлээн авч боловсруулсан. Энэ баримт бичигт ашигласан монгол фонт нь Bolorsoft-ийн нээлттэй эх сурвалжын скрипт дээр үндэслэгдсэн.

The author received feedback from the Script Ad Hoc group while drafting this document. The prototype Mongolian font used in this document is based on Bolorsoft's open-source font Mongolian Script.

1.3 Хамрах хүрээ

Scope

Энэхүү загвар нь Монголчуудын бичгийн хэрэглээний хамгийн чухал хэсгийг багтаасан бөгөөд орчин үеийн монгол хэлний Худам (* уу буюу худам) бичгийн системийг ч хамаруулсан юм. Худам Али Гали, түүхийн Худам, Тодо, Тодо Али Гали, Манж, Манж Али Гали, Манжийн түүх, Шибэ, гэх мэт бусад бичгийн системүүд

Хэдийгээр Арванхоёр үеийн цагаан толгой нь бүх дүрслэх дүрмийг харуулах боломжгүй ч, дээрхи хүснэгтнээс авиан үсэгнээс бичгийн системийн хамгийн бага утгыг агуулсан харьцуулалтыг авиан үсэг дангаараа буулгах боломжгүй юм, тиймээс үүнийг тайлбарлах нэмэлт мэдээлэл шаарддаг.

Урьдаас тодорхойлсон ерөнхий дүрсүүдийг үүсгэх боломжгүй тохиолдолд заавал гаргаж авахыг хүсэж буй дүрсээ удирдах товчлуурын тусламжтайгаар үүсгэдэг тул, авиан үсгийг кодчилон оруулж дүрсийг гаргаж авхын тулд удирдах товчлуур шаардлагатай болдог.

Although the Twelve Syllabaries naturally don't exhibit all the shaping rules, it's already partly observable from the chart above how the mapping from phonetic letters to graphemes is not predictable from phonetic letters alone, but requires additional information. Encoding the phonetic letters as characters thus requires inserting arbitrary format control characters when predefined general shaping rules can't produce desired forms.

1.5 График арга

Graphetic approach

Тогтворгүй байгаа байдлын үүднээс, хэсэг бүлэг шинжээчид хураангуй, товчилсон авиан үсгийн оронд дүрсийн аргыг (авиа зүйн эсрэг утга) нэгтгэх аргыг хайж байна, ингэснээр Юникодын- Open Type Cursive Joining Model загварыг ашиглаж байгаа боловч кодлохдоо авиан үсгээс зөрүүтэй урьдчилан тааварлашгүй зураглал хийхээс зайлсхийж байна. График дизайныг авахын тулд, заримдаа түсгэй кодчиллын загварыг хэлэлцэж, туршихад зориулагдсан загварыг ашиглаж байгаа ба энэхүү баримтанд дэлгэрэнгүй танилцуулсан байгаа.

In light of the unsustainable status quo, a group of experts have been exploring a graphetic (as opposed to phonetic) approach that encodes cursive joining graphemes instead of abstract phonetic letters, thus still utilizing the Unicode – Open Type Cursive Joining Model but avoiding implementing the unpredictable mapping from phonetic letters to graphemes in encoding. Taking the graphetic approach, a specific encoding model has been designed for discussing and testing, and is introduced in this document in details.

2. Үсгийн нэр төрөл

Character repertoire

Энэ загвар нь 33 тэмдэгтийн нэр төрөл шаарддаг. Глифийн төлөөлөлүүд болон үсгийн нэр төрлийн жагсаалтыг доорх хүснэгтээс харна уу. Глифийн төлөөлөлүүд нь яг тохирсон хэлбэр дүрсүүд биш. Оруулсан үсгийн оронд Slashes (/) О Ховорууг уншихад хялбар болгоорой. Энд ашиглагдсан, О тэмдэг нь цөөн хэрэглэгддэг ба илүү сайн уншихад зориулагдсан, (/) тэмдэг нь үгийн оронд ашиглагдаж байна гэсэн үг. The model requires a repertoire of 33 characters. See the chart below for a provisional list of representative glyphs and character names. Note representative glyphs are not relevant to actual shaping. Slashes(/) instead of the word O Rare used herefor better readability.

Төлөөл.

Нэр, төрөл

- MONGOLIAN CHARACTER NIRUGU

MONGOLIAN CHARACTER ALEPH / A / E / NA ᠎

MONGOLIAN CHARACTER A / E ᠎

MONGOLIAN CHARACTER I / JA / YA

᠎ MONGOLIAN CHARACTER O / U / OE / UE

MONGOLIAN CHARACTER O / U / OE / UE / WA

MONGOLIAN CHARACTER OE / UE

᠎ MONGOLIAN CHARACTER NA

᠎ MONGOLIAN CHARACTER BA

᠎ MONGOLIAN CHARACTER PA

᠎ MONGOLIAN CHARACTER HA / GA FORM ONE

᠎ MONGOLIAN CHARACTER HA / GA FORM TWO

᠎ MONGOLIAN CHARACTER GA

᠎ MONGOLIAN CHARACTER MA

᠎ MONGOLIAN CHARACTER LA

᠎ MONGOLIAN CHARACTER SA

᠎ MONGOLIAN CHARACTER SHA

᠎ MONGOLIAN CHARACTER TA / DA FORM ONE 3

᠎ MONGOLIAN CHARACTER TA / DA FORM TWO

MONGOLIAN CHARACTER DA

᠎ MONGOLIAN CHARACTER CHA

MONGOLIAN CHARACTER JA

ᠮ MONGOLIAN CHARACTER YA

ᠮ MONGOLIAN CHARACTER RA

ᠮ MONGOLIAN CHARACTER WA / LOANWORD E

ᠮ MONGOLIAN CHARACTER LOANWORD FA

ᠮ MONGOLIAN CHARACTER LOANWORD KA

ᠮ MONGOLIAN CHARACTER LOANWORD TSA

ᠮ MONGOLIAN CHARACTER LOANWORD DZA

ᠮ MONGOLIAN CHARACTER LOANWORD HA

ᠮ MONGOLIAN CHARACTER LOANWORD HA / LOANWORD JA

ᠮ MONGOLIAN CHARACTER LOANWORD RA

ᠮ MONGOLIAN CHARACTER LOANWORD CHA

3. Текстийн дүрслэл

Text representation

Цогц текстийг кодчилох загвар нь 2 үе шаттайгаар нөхцөл байдлаас хамаарсан, зөв бичих дүрмийн дагуу тохирсон зөв үсэг, хувилбарыг сонгож хийдэг. Үүнд:

As a complex text encoding, the model requires two stages of contextual shaping to select the orthographically correct form for a character:

1. Товчилсон нийлбэрийн хүрээнд, үсэг нэгтгэх зөв дүрмийн дагуу үсэг тус бүр зөв тохирсонөөрийн онцлогтой байна.

At the cursive joining stage, characters are mapped to appropriate positional variants, according to joining types of surrounding characters.

2. Дараа нь дугуй дүрстэй гийгүүлэгчийнхүрээнд, гийгүүлэгч үсгийн хэв маяг мөн тус үсгийн бусад үсэгтэй нийлэхэд гарсан нарийн тоддорхой хэлбэр, хувилбарууд руу шилждэг.

Then, at the round consonant stage, in the specific context of a round consonant grapheme and its following grapheme, positional variants are further transformed to specific variants.

3.1 Товчилсон нийлбэр

Cursive joining stage

Товчилсон нийлбэрийн хэлбэрүүд нь “үг”-ийн дүрмийн тодорхойлолттой хамааралгүй байдаг. Цэнхэр өнгийн нуруу нь “” байрлалын хувилбартай нийлсэн хэсгийг тодотгон харуулж байна. Ихэнхи үсэгнүүд давхардсан хэлбэрээр нийлдэг ба харин хоёр үсгийн хувьд өөр байдаг. Үүнд: “MONGOLIAN CHARACTER A / E” үсгийн хувьд нийлдэггүй ба “MONGOLIAN CHARACTER O / U / OE / UE / WA” үсгүүдийн хувьд хэд хэдэн

давхардсан хэлбэрээр нийлсэн байдаг. (цэнхэр өнгийн “x” энэхүү тэмдэглэгээр дээрх хоёр үсгийн өөрсдийнх нь онцлог байдлыг тэмдэглэсэн). Хоосон нүднүүд нь тодорхойгүй хэлбэр, эсвэл үзүүлэнгүүдийг үүсгэхийн тулд зааж өгөөгүй байна.

Note cursive joining positional variants are irrelevant to the grammatical definition of “word”. Joined ends of positional variants are emphasized with a blue nirugu “-”. Most characters are dual -joining, except two characters: “MONGOLIAN CHARACTER A / E” is non-joining and “MONGOLIAN CHARACTER O / U / OE / UE / WA” is top-joining (both have their in valid positions marked with a blue“x”). Empty cells denote unattested positions, which should also be specified to produce certain forms or visual aids.

Үсгийн нэр	Эхэнд,	Д дунд,	Адагт ,	Салангид
MONGOLIAN CHARACTER NIRUGU	-----			
MONGOLIAN CHARACTER ALEPH / A / E / NA	а--а--а ^᠗			
MONGOLIAN CHARACTER A / E	x x x ^᠗			
MONGOLIAN CHARACTER I / JA / YA	i--i--i ^᠗			
MONGOLIAN CHARACTER O / U / OE / UE	᠋--u--u			
MONGOLIAN CHARACTER O / U / OE / UE / WA	x x x ^᠗			
MONGOLIAN CHARACTER OE / UE	-᠋--			
MONGOLIAN CHARACTER NA	᠗-- --			
MONGOLIAN CHARACTER BA	᠋-- --			
MONGOLIAN CHARACTER PA	᠋-- --			
MONGOLIAN CHARACTER HA / GA FORM ONE	᠋-- --			
MONGOLIAN CHARACTER HA / GA FORM TWO	᠋-- --			
MONGOLIAN CHARACTER GA	᠋-- --			
MONGOLIAN CHARACTER MA	᠋-- --			
MONGOLIAN CHARACTER LA	᠋-- --			
MONGOLIAN CHARACTER SA	᠋-- --			
MONGOLIAN CHARACTER SHA	᠋-- --			
MONGOLIAN CHARACTER TA / DA FORM ONE	᠋-- --			
MONGOLIAN CHARACTER TA / DA FORM TWO	᠋-- --			

MONGOLIAN CHARACTER DA - ᠳ -

MONGOLIAN CHARACTER CHA - ᠴ -

MONGOLIAN CHARACTER JA - ᠵ -

MONGOLIAN CHARACTER YA - ᠶ -

- MONGOLIAN CHARACTER RA - ᠷ -

MONGOLIAN CHARACTER WA / LOANWORD E - ᠠ -

MONGOLIAN CHARACTER LOANWORD FA - ᠮ -

MONGOLIAN CHARACTER LOANWORD KA - ᠨ -

MONGOLIAN CHARACTER LOANWORD TSA - ᠯ -

MONGOLIAN CHARACTER LOANWORD DZA - ᠰ -

MONGOLIAN CHARACTER LOANWORD HA - ᠬ -

MONGOLIAN CHARACTER LOANWORD HA / LOANWORD JA - ᠬᠵ -

MONGOLIAN CHARACTER LOANWORD RA - ᠷ -

- MONGOLIAN CHARACTER LOANWORD CHA - ᠴ -

Жишээ нь :

< a ALEPH / A / E / NA, a ALEPH / A / E / NA > → aa

< a ALEPH / A / E / NA > → ᠠ < ᠠ A / E > → ᠠ

< a ALEPH / A / E / NA, i I / JA / YA > → ai < i I / JA / YA > → ᠠᠢ

< a ALEPH / A / E / NA, ᠨ O / U / OE / UE > → au

< a ALEPH / A / E / NA, OE / UE > → a < O / U / OE / UE / WA > → ᠠᠤ

< a ALEPH / A / E / NA, a ALEPH / A / E / NA, ᠯ LA, ᠷ TA / DA FORM TWO,

a ALEPH / A / E / NA, a ALEPH / A / E / NA, a ALEPH / A / E / NA, ᠨ O / U / OE / UE, ᠷ TA / DA FORM TWO, O / U / OE / UE / WA > → aaaaau < ᠬ HA / GA FORM ONE,

a ALEPH / A / E / NA, ᠮ NA, a ALEPH / A / E / NA > → ᠮaa < ᠬ HA / GA FORM ONE,

a ALEPH / A / E / NA, ᠮ NA, ᠠ A / E > → ᠮᠠᠠ < ᠮ NA, a ALEPH / A / E / NA, i I / JA / YA, ᠮ MA,

a ALEPH / A / E / NA > → ʼaia < ʼ SA,

a ALEPH / A / E / NA, i I / JA / YA, i I / JA / YA,

a ALEPH / A / E / NA > → ʼaia < ʼ SA, a ALEPH / A / E / NA, ʼ YA, i I / JA / YA, ʼ HA / GA FORM ONE, a ALEPH / A / E / NA, a ALEPH / A / E / NA > → ʼaia < ʼ CHA,

a ALEPH / A / E / NA, ʼ RA, i I / JA / YA, ʼ HA / GA FORM TWO > → ʼai < i I / JA / YA,

a ALEPH / A / E / NA, ʼ RA, ʼ LA, i I / JA / YA, ʼ HA / GA FORM ONE > → ʼai < ʼ BA, i I / JA / YA, ʼ LA, i I / JA / YA, ʼ HA / GA FORM TWO, ʼ BA,

a ALEPH / A / E / NA, ʼ TA / DA FORM TWO, O / U / OE / UE / WA > → i

< a ALEPH / A / E / NA, O / U / OE / UE / WA > → a

< a ALEPH / A / E / NA, ʼ TA / DA FORM TWO > → a

< ʼ HA / GA FORM TWO, ʼ HA / GA FORM TWO, i I / JA / YA, ʼ RA > → ʼ < ʼ RA,

a ALEPH / A / E / NA, ʼ TA / DA FORM TWO, i I / JA / YA, ʼ O / U / OE / UE > → ʼai < ʼ TA / DA FORM ONE,

a ALEPH / A / E / NA, ʼ HA / GA FORM TWO, ʼ RA, i I / JA / YA > → ʼai < ʼ SA, a ALEPH / A / E / NA, i I / JA / YA, i I / JA / YA, a ALEPH / A / E / NA > → ʼaia < ʼ BA,

a ALEPH / A / E / NA, i I / JA / YA, i I / JA / YA, ʼ NA, ʼ A / E > → iiʼ

< ʼ O / U / OE / UE, O / U / OE / UE / WA > → ʼ

Удирдах товчлуурын хувьд, U + 200C ZERO WIDTH NON-JOINER болон U + 200D ZERO WIDTH JOINER гэсэн ерөнхий тэмдэгтүүдийг хэрэглэж болох боловч өдөр тутмын текстийг ашиглан MONGOLIAN CHARACTER NIRUGU-тай нийлсэн маягтуудыг үүсгэхийг зөвлөж байна. Учир нь NIRUGU нь хэрэглэгчидэд харагдахаас гадна ашиглахад хялбар байдаг. Жишээ нь:

For joining control, although the general characters U+200C ZERO WIDTH NON-JOINER and U+200D ZERO WIDTH JOINER can be used, in day-to-day text it is recommended to produce joined forms in isolation with MONGOLIAN CHARACTER NIRUGU, because NIRUGU is visible and easy for users to manipulate. Examples:

< ʼ GA, a ALEPH / A / E / NA, ʼ NIRUGU > → ʼaʼ

< ʼ NIRUGU, ʼ GA, a ALEPH / A / E / NA, ʼ NIRUGU > → ʼaʼ

< ʼ NIRUGU, ʼ GA, ʼ A / E > → ʼ

3.2 Дугуй дүрстэй гийгүүлэгч

Round consonant stage

4.1 Илүү олон график бүрэлдэхүүн

More graphetic

Ойлгомжгүй, хоёрдмол утгатай болгохын тулд авиан үсгүүдтэй илүү зохимжтой золиослохын тулд, илүү олон графикуудыг бүрэлдэхүүн хэсгүүдэд хуваах хэрэгтэй.

In order to further reduce visual ambiguity while sacrificing more relationship to phonetic letters, decompose more graphemes into components:

Одоогийн хувилбар	→	Хэсэгт хуваах	Тэмдэглэл
OE / UE	→	< ^œ O / U / OE / UE, i l / JA / YA >	Зөвхөн үгийн дунд орох хувилбарыг хуваах.
DA	→	< ^œ O / U / OE / UE, a ALEPH / A / E / NA >	
ᶯ гадаад үгэнд хэрэглэгдэх ХА	→	< a ALEPH / A / E / NA,	
ᶮ гадаад үгэнд хэрэглэгдэх ХА /ЖА			
ᶯ гадаад үгэнд хэрэглэгдэх ЧА	→	< ^œ O / U / OE / UE, ^œ O / U / OE / UE >	

4.2 Илүү олон авиа

Morephonetic

Авиа үсгийг авч хадгалахын тулд, тодорхойгүй, ойлгомжгүй байгаа болон тодорхой фонтой холбоогүй бүтцийг өөрчилж, нэгтгэх хэрэгтэй.

In order to retain more phonetic letters while tolerating more ambiguity and some cross-graphemecharacters, revoke the merger of certain phonetically unrelated structures:

Одоо	→	Салгах	Тэмдэглэл
a ALEPH / A / E / NA	→	ALEPH and NA FORM TWO	
i l / JA / YA	→	JA FORM TWO	Нийлүүлсэний дараагаар ЖА үсгийг салгах.
ᶯ WA / LOANWORD E	→	гадаад үгэнд хэрэглэгдэх E	
ᶮ гадаад үгэнд хэрэглэгдэх ХА /ЖА	→	гадаад үгэнд хэрэглэгдэх ЖА	Дараа нь үлдсэн үгийн дунд, адагт орох хэлбэрийг нэг тгэн "гадаад үгэнд хэрэглэгддэг ХА" үсгийг үүсгэнэ.