

G2P Ангилах арга: Автоматаар ангилан ялгах арга буюу Монгол бичгийн дүрсийг ангилан кодлох арга

G2P Sorting: An Automated Natural Sorting Method for Graphetically Encoded Mongolian

Зохиогч: Шэн Илэй

Author : SHEN Yilei

Огноо: 2018 оны 03 сарын 24

Date: 2018-03-24

G2P ангилах аргачилал буюу автоматаар ангилан ялгах арга буюу Монгол бичгийн дүрсийг ангилан кодлох аргачилалыг энэхүү баримтанд тусган харуулсан болно.

G2P sorting, as an automated natural sorting method for graphetically encoded Mongolian is presented in this document.

1. Монгол бичгийг ангилан ялгах конвенц

Mongolian sorting convention

1.1 Ангилчлах аргыг задлан харуулах нь.

Unfolded sorting

Хэдийгээр Монголчууд уугуул иргэддээ Монгол хэл зааж сургахдаа ерөнхийдөө үелэх хэлбэрийн аргачилалаар зааж сургадаг боловч, орчин үеийн Монгол хэлний толь бичиг нь цагаан толгойн дарааллаар ангилж оруулсанбайдаг. Гол төлөв зохиосон цагаан толгойн дараалалаар ангилсан байдал нь: Although Mongolian is generally taught to native Mongolians in a syllabic or semisyllabic way, modern Mongolian dictionaries are sorted alphabetically. The commonly assumed alphabet (unfolded alphabet) used in sorting is

a · e · i · o · u · ö · ü · ё · n · ᠨ · b · p · x · g · m · l · s · š · t · d · č · j · y · r · w · f · k · c · z · h · ř · † · ž · ċ

маш цөөн тооны үндсэн хувилбаруудтай. Үүнд:

with only a few minor variations:

- ё (ᠡ) нь биеэ даасан үсэг эсвэл (ᠨ) хувилбарын нэг гэж авч үздэг.
ё (ᠡ) is treated either as an independent letter or as a variant of e (ᠡ);
- ᠨ (ᠨ) нь мөн л биеэ даасан үсэг эсвэл (<ᠨ, ᠨ/>) хувилбарын нэг болгон авдаг.
ᠨ () is treated either as an independent letter or as a letter sequence ng (<ᠨ, ᠨ/>).

Энэхүү Монгол цагаан толгойн үсгийн тэмдэглэгээ нь тод бус байдлаар тодорхойлогддог. Мэдээжээр, энэ тохиолдолд цагаан толгойн үсгийн хэд хэдэн бүлэг үсэг байдаг бөгөөд энэ нь тодорхой тохиолдлуудад ижилхэнбичсэн (a / e / n, ё / w, o / u / ö / ü, x / g, t / d гэх мэт) хувилбарууд нь өөр хоорондоо адил байдаг. Мөн үсгийн өвөрмөц шинж чанарууд нь ерөнхийдөө ярианы хэллэгийн дуудлаган дээр суурьлаж харьцуулсан байдаг. Гэсэн хэдий ч, хувь хүмүүс тухайн үсгийг тодорхойлж чадахгүй бэрхшээлтэй тулгардаг байна.

This Mongolian alphabet is characterized by its opaque identification of letters. As is known, there are several groups of letters in this alphabet which share the same written forms in certain cases (a/e/n, ё/w, o/u/ö/ü, x/g, t/d, etc.), and the distinctions of letters are established generally on the contrasts in colloquial pronunciations. However, individual cases of identifying letters can be problematic:

- Түүхийн явц дах авиа зүйн өөрчлөлттэй холбоотойгоор, зөвхөн ярианы хэл болон дуудлага дээр тулгууралсан үсгийг тодорхойлоход боломжгүй байдаг. Авиа зүй буюу фонологийн мэдлэгээс гадна, зөв бичих дүрэм, үгийн гарал зүй, бичигийн үсэг болон хэлний үг зүйн бүтэц ч мөн үүнд хамааралтай байдаг. Тиймээс хүн бүрийн гаргаж байгаа тодорхойлолтууд нь үргэлж өөр хоорондоо ялгаатай байх ба санал зөрөлддөг учраас хүн хүний гаргасан тодорхойлолтууд өөр өөр байдаг. Жишээ нь: "ᠨᠢ" гэдэг

Монгол үгийг бичихдээ *čino_a*, *činu_a*, эсвэл *činw_a* гэж бичигдэж байгаа хувилбар, шалгуурууд нь өөр хоорондоо зөрчилдөж хүн бүр өөр өөр тодорхойлж бичиж болдог.

Due to drastic phonological changes in history, identifying letters relying merely on colloquial pronunciation is not always possible. Besides phonological knowledge, one's orthographical, etymological, and morphological knowledge of the language/script also take part. Therefore it is quite often that one's identification may differ widely from another's as these criteria do not always converge. For example, the Mongolian word *ᠠᠮᠣᠯᠤᠯ* "wolf" may be identified as *čino_a*, *činu_a*, or *činw_a* according to different criteria that contradict each other.

- Зөвхөн фонологийн буюу авиа зүйн мэдлэгийн хүрээнд, янз бүрийн нутгийн хэл, аялагаар ярьдаг хүмүүс өөрсдийнхөө орон нутгийн дуудлагуудын дагуу үсгээ тогтоох хандлагатай байдаг. Хэдийгээр стандарт дуудлагад бидний ярилцлагыг хязгаарласан ч гэсэн, Хятад болон Монгол улсын стандарт дуудлагууд нь олон талаараа өөр хоорондоо маш их ялгаатай болохыг тэмдэглэх нь зүйтэй. Гэсэн хэдий ч, өнөөгийн Хятад улс дахь Монголын стандартчилагдсан дуудлага нь Хятадад байгаа Монгол хүмүүсийн дунд өргөн хэрэглэгдэхгүйгээс гадна хамгийн оновчтой авиа зүй хэрэглэгддэггүй байна.

As for phonological knowledge alone, people of various dialectal backgrounds tend to identify letters according to their local pronunciations. Even if confining our discussion to standard pronunciations, it should be noted that the standard pronunciations of China and Mongolia differ in many aspects. Worse still, the standard Mongolian pronunciation of China is an idealized phonology which is native to nowhere and is poorly popularized in China today.

Үүний үр дагавар нь, ангилчлалыг задлан харах үйл явц нь мөн эргэлзээтэй байна гэсэн үг.

As a consequence, the usability of unfolded collation is questionable.

- Нэг үг нь яг тодорхой үсгийн дараалалтай байх ба толь бичигт үсгийн дараалалын дагуу бичигдсэн болон зохиогч нь тухайн үгийг өөр өөр тэмдэглэгээгээр тэмдэглэсэн тохиолдолд, тус үгийн үсгийн жинхэнэ дараалалыг олоход бүр ч их хэцүү байдаг.

One word may be identified as various letter sequences and thus sorted variously among dictionaries, and one will not be able to find the entry when they has a different identification of the underlying letters from that of the dictionary compiler.

- Үгийг судлах явцад тухайн үгийн зөв дуудлагыг мэдэхгүй бол, толь бичигт оруулахын тулд зөв бичих дүрмийг нэг бүрчлэн эхнээс нь үзэж судлах шаардлагатай болно. Гэсэн хэдий ч, бичигдсэн нэг үгийг 4-өөс доошгүй хэлбэрээр унших нь элбэг тохиолддог. Жнь: *talai*, *telei*, *dalai*, and *delei*.

When encountering a word they does not know how to pronounce correctly, one will have to go over every orthographically possible reading in order to find it in a dictionary. However, it is quite common for a single written form to have no less than four possible readings like *talai*, *telei*, *dalai*, and *delei*.

1.2 Ангилчлах аргыг нэгтгэх нь

Folded sorting

Гэсэн хэдий боловч, Монгол бичгийг ангилчлан нэгтгэхэд тулгарч байдаг асуудлууд ба дээр дурдсан орчин үе нь өмнөх үеийнхээс зайлсхийхэд хүрч байна. Хэний зарчим нь хамгийн энгийн зөв болох, нэг үеийн бүрэлдэхүүн, (эхэнд, дунд, адагт г.м.), нэг хэлбэр нь зөвхөн нэг үсгэн дотор багтаах, өөрөөр хэлбэл ангилах үйл явц нь үелэх бүтэц болон бичих аргачилал нь тодорхой болсон тохиолдолд бүрэн тодорхойлогдож байна гэсэн үг. Энэхүү аргыг "нэгтгэх арга" гэж нэрлэнэ. Учир нь ангилалын явцад үсэг болон үгүүд нь адил бичигддэг хэдий ч өөр утгыг агуулгыг ялгаж ангилахад хэрэглэгддэг. Доорхи хэсэгт:

Nevertheless, there has been another approach to Mongolian sorting dating back to pre-modern times which can evade these above-mentioned shortcomings, whose principle is quite intuitive: for each component of the syllable (onset, nucleus, or coda), one shape is subsumed under only one letter, which means the sorting behavior for a specific written form is fully determined

as long as the syllable structure is known. This approach is termed here as “folded” because a syllabary/alphabet that folds homographic written forms is adopted in sorting. In this folded approach,

- (ᠲᠡᠯᠡᠢ) гэдэг үгийг telei and dalai гэж ангилж бичсэн байна. Учир нь энэ 2 адилхан бичигдсэн хэдий ч хоёрдмол утгыг хамгийн ихээр үгүйсгэсэн байна. telei and dalai (ᠲᠡᠯᠡᠢ) are sorted like *talai because they are written alike, which has maximally eliminated ambiguous reading;
- (ᠪᠡᠶ) гэдэг үгийг *bay_a, биш харин *bai_a (cf. bai (ᠪᠠᠢ)) ангилж бичсэн байна. Үечлэх бүтэц нь хэл шинжлэлийн хамгийн чухал ач холбогдолтой ангилалд багтдаг. bey_e (ᠪᠡᠶ) is sorted like *bay_a, not *bai_a (cf. bai (ᠪᠠᠢ)), where the syllable structure is retained to make the sorting linguistically significant.

Ангилчлах аргыг нэгтгэхэд үүсэх санал зөрчилдөөний улсаас хоёр өөр үзэл бодол гарах магадлалтай. Үүнд: Arguments against folded sorting may come in two aspects:

- Энэ нь олон жилийн туршид хэрэглэгдэж ирсэн ялгах ангилалын нийтлэг практиктай зөрчилддөг. It goes against the common practice of unfolded sorting which has been established for several decades.
- Эдгээр ангилан нэгтгэх арга нь хэрэглэгчдэд цээжлэхэд хэцүү байдаг тул эдгээр нь ярианы хэл болон дуудлага дээр тулгуурлаагүй байдаг. The folded letters are difficult for users to memorize because they are not based on colloquial pronunciations.

Миний батлан харуулах гэж байгаа ангилчлах арга мөн хоёр талыг тусган харуулж байна. Үүнд My defense of folded sorting also comes in two aspects:

- Өнөөдөр өргөнөөр хэрэглэгддэг ангилан ялгах аргачилал нь уламжлалт аргыг сайн тогтоож өгөөгүй. The unfolded sorting widely adopted today itself is not a well-established tradition.
 - Үечлэн ангилах хэлбэрээс цагаан толгойн ангилах хэлбэрлүү шилжих үйл явц зөвхөн зуун жилээс дээшгүй байна. The shift from syllabic sorting to alphabetical sorting took place only no more than a century ago.
 - η тэмдэгийг Хятад хэлээр хэвлэгдсэн толь бичиг дээрөөр ялгаатай байдлаар ангилан бичсэн байдаг. There are still discrepancies in unfolded sorting like how to arrange η in dictionaries published in China alone.
 - Хятад хэлний стандарт доторх Монгол хэлний ялгах стандарт (GB / T 30851-2014) нь өргөн цар хүрээг хамарсан хэдий ч Хятад хэлний скриптүүдийн стандартуудаас олон талаараа зээлж авсан байдаг. (GB / T 32912-2016 г.м). Илүү нарийн төдийгүй, ангилах стандартын текстийн тайлбар нь тэдгээрийн харьцуулах хүснэгтэй өөр хоорондоо мөн зөрчилддөг. The Chinese standard of Mongolian sorting (GB/T 30851-2014) contradicts widely adopted sorting practices and other Chinese standards of the script (e.g., GB/T 32912-2016) in many aspects. More ironically, the text description of the sorting standard contradicts with its collation table.
 - Хэрэглэгчид удаан хугацаагаар үсэгний ангилалаас үүсэн тогтворгүй байдал дээр дассан байна. Users have long been accustomed to the sorting instability resulting from fickle letter identification.
- Ангиллах аргыг нэгтгэх нь хэрэглэгчидэд хэрэглэж, эзэмшихэд илүү хялбар байдаг. Folded sorting is much easier for users to master.
 - Нэгтгэсэн үсэгнүүд нь цээжлэхгүй, харин шүүд бичсэн хэлбэрээс авдаг. Folded letters are not memorized but derived directly from the written form.
 - Ангиллалыг нэгтгэх арга нь үгийн гарал зүй, үг зүйн бүтэцийн хувилбаруудын хувьд ангилчлан задлах аргатай адилхан байх шаардлаггүй. Зөвхөн ерөнхий зөв бичих дүрэм нь өгөгдсөн хэлбэрээс ангилан нэгтгэх дарааллыг тодорхойлох бараг л хангалттай байдаг хэдий ч ярианы хэллэг хоёрдмол утгатай бол дуудлагаар ярих нь тустай байж болох юм.

Unlike unfolded sorting, etymological or morphological knowledge is no longer necessary in folded sorting. Basic orthography alone is almost enough to determine the folded letter sequence from a given shape, though colloquial pronunciation might be marginally helpful in case that syllabification ambiguity arises.

- Олон нийтэд зориулсан тоон текстэнд алдаа бичиж тэмдэглэх үед, хэрэглэгчид бичиж байхдаа нэгтгэсэн цагаан толгойн хэсгийг хэсэгчилсэн гэж үздэг.
- Typing errors in digital text contributed by the public shows that users have partially assumed a folded alphabet in typing.

Эцэст нь дүгнэж хэлэхэд, тодорхой бус талыг ангилан ялгахгүйн тулд эрэмбэлэх боломжгүй юм.

In the end, the conclusion is that in no definite aspect is folded sorting inferior to unfolded sorting.

2. Автомат ангилал ба бүдүүвч зураглалыг кодлох

Sorting automation and encoding schemes

Кодчлолын асуудалд эргэж орсны дараа автомат ангилалд гарч буй асуудлууд нь бидний санаа зовж буй асуудал юм. Хэдийгээр яг тодорхой бүдүүвч зураглалыг кодлоход болон автомат ангилалын хэрэгжүүлэлтийн талаар ярих боломжгүй ч, бүдүүвч зураглалыг кодлох болон автомат ангилал нь үндсэндээ хоёр өөр хэмжүүртэй хэдийч хоорондоо огтлолцдог.

As we turn back to the encoding issue, problems that arise in sorting automation is what we are concerned about. Although it is impossible to talk of implementation of automated sorting without addressing a specific encoding scheme, encoding schemes and automated sorting are essentially two dimensions and intersect each other.

Авиа зүй болон дүрсээр кодчилох аргад зориулагдсанавтоматаар ангилалын талаарх тайлбарыг 1-р хүснэгтэд үзүүлэв. Үсэгийг ангилахын тулд нөхцөл байдлаас үл хамаарсан кодлогдсон үсгийг ангилахад чиглэгддэг. Хэд хэдэн ялгах тэмдэглэгээг эс тооцвол (авиа зүйн кодчилолд FVS гэх мэт) ялгах тэмдэгтүүдийг үсгээр тэмдэглэсэн байдаг.

Remarks on various automated sorting methods for phonetic and graphetic approaches are summarized as in Table 1. Here character sorting refers to sorting encoded characters without any contextual sensitivity (expansion or contraction). Except for a few diacritical characters (such as FVSes in phonetic encoding), character sorting treats encoded characters as letters.

Хүснэгт 1. Автомат ангилалын арга × кодлох аргачилал

Table 1. Automated sorting methods × encoding approaches

	Үсгийг ангилах Character sorting	Энгийн ангилал Natural sorting	
		Ангилчлалыг задлах Unfolded sorting	Ангилчлалыг нэгтгэх Folded sorting
Авиа зүйн кодчилол Phonetic encoding	Найдваргүй Unreliable		Хэрэгжиж болохуйц Practicable
Дүрсийн кодчилол Graphetic encoding	Ашиглах боломжгүй Unusable	Боломжгүй Impossible	Хэрэгжиж болохуйц Practicable

Авиа зүйг кодлох тогтолцоо

For phonetic encoding schemes:

- Үсгийг ангилах болон ангиллалыг задлахад ихэнх тохиолдолд адил үр дүн гардаг. Өнөөгийн хамгийн өргөн хэрэглэгддэг Монгол хэлний автомат ангиллах аргын үр дүнд нь найдваргүй байна. Учир нь нэг оруулж байгаа хувилбар нь өөр өөр үсгийг илэрхийлж хэрэглэгч тус бүр өөр өөр тэмдэгтийг ашиглаж байгаатай холбоотой. Энэхүү бичилтэнд нийцэхгүй байгаа хэд хэдэн шалтгаан бий. Үүнд:
Character sorting and unfolded sorting yield identical results in most of the cases. Being the most widely used automated sorting of Mongolian today, however, the sorting results are unreliable, because one entry can be represented with different character sequences by different users. There are several reasons for this typing inconsistency:
 - Бичилтийн алдаа: Бичээч нар бичиж байх явцдаа алдаагаа анзаарахгүй байх.
Typo: Typists may not always be able to notice their typing errors.
 - Үсгийг буруугаар ашиглах: Бичээч нар үсгийн хэрэглээ өөр байгааг тодорхойлж болох талтай.
Misuse: Typists may identify letters differently.
 - Хэтрүүлэн ашигах: Бичээчид өөрсдийн зорилгодоо хүрэхийн тулд зарим буруу товчлууруудыг санаатайгаар ашигладаг. Тиймээс, янз бүрийн эх сурвалжаас бичсэн текстийг нэгтгэх үед нэг оруулгын ангилах үр дүн нь хэд хэдэн газруудад гарч ирдэг.
Abuse: Typists may favor some incorrect keystrokes deliberately for the sake of expediency. Therefore, when text from different sources is aggregated, one entry can appear in several places of the sorting result.
- Нэгтгэх ангилал нь авиа зүйг кодчилоход хэрэгжих боломжтой, яагаад гэвэл энэ нь мөн л дүрсээр кодчилоход хэрэгжих боломжтой байдаг. (доорхийг уншина уу) Авиа зүйн кодчилал нь дүрсээр кодчилохын тулд дүрсээр кодчилох бүх мэдээлэлийг дамжуулах шаардлага гардаг.
Folded sorting is practicable for phonetic encoding because it is also practicable for graphetic encoding (see below), and phonetic encoding contains all information graphetic encoding conveys and beyond.

Дүрсээр кодлох тогтолцоо

For graphetic encoding schemes:

- Шууд дүрсийн үсгийг ангилахдаа, хамгийн бага өртөгөөр хэл шинжлэлийн хувьд тодорхойгүй болон хэрэглэгчидэд ч мөн ойлгоход бэршээлтэй байдаг.
Direct graphetic character sorting, being the least costly method, is not linguistically significant, and is incomprehensible to users.
- Ангиллах аргыг задлах боломжтой, учир нь кодлогдсон текстээс хойш адил бичигддэг өөр утгатай үг болон үсгийн хэлбэрийн хооронд ялгаа байдаг.
Unfolded sorting is impossible, because the distinctions between homographic letterforms are missing from the encoded text.
- Дүрсийг кодлон ангиллах аргыг нэгтгэхэд хэрэгжүүлж болохуйц байна.
Folded sorting of graphetic encoding is practicable.
 - Энэ нь хэрэглэгчдэд ойлгомжтой болж, ингэснээр 1.2-хэсэгт хэлэлцсэн заасанчилан ашиглагддаг.
It is comprehensible to users and thus usable, as has been discussed in Section 1.2.
 - Энэ нь техникийн хувьд хэрэгжих боломжтой, учир нь үсгийн нэгтгэсэн мэдээлэл нь нөхцөл байдлаас хамааран өөрчлөгдсөн дүрсийн холбооноос сэргээн босгогддог. Үүнийг доор тайлбарласан байгаа.

It is technically implementable, because the information of folded letters is still largely reconstructible from graphetic strings through a series of contextual transformations, which I will demonstrate shortly.

Автоматаар дүрсийг кодчилон ангилах, ялангуяа G2P ангилах арга нь маш их ач холбогдолтой юм. Учир нь энэ нь Монгол хэлийг кодчилсон монгол хэл дээрх автоматаар ангилах цорын ганц арга юм.

Automated folded sorting of graphetic encoding, or G2P sorting, is particularly of great significance because it seems to be the only practical solution to automated sorting for graphetically encoded Mongolian.

3. G2P ангилалын танилцуулага

Introduction to G2P sorting

G2P ангилал арга нь хоёр үе шаттайгаар хэрэгжсэн. Үүнд:

G2P sorting is implemented in two steps:

- 1-р үе: дүрсийн үсгийн холбоосоос үсгийг нэгтгэх дараалалыг дахин сэргээн завсарладаг.
Step 1: Reconstruct folded letter sequences from graphetic character strings.
- 2-р үе: Оруулгуудыг үсгийг нэгтгэх дарааллаар нь ангилдаг.
Step 2: Sort the entries by their folded letter sequences.

G2P ангилах чадвар нь Юникодын Харьцуулалтын Алгоритмийн гол алгоритм дотор нөхцөл байдлаас хамаарсан хязгаарлагдмал нөхцөлөөс шалтгаалан багтаж чаддаггүй. (харьцуулах элементийн хүснэгтэд цөөн хэдэн товчлол нэмэхэд ажиллана, Thai/Lao CV гэх мэт хувиргаж оруулана гэсэн үг). Гэсэн хэдий ч, үсгийн дүрсийг үсгийн нэгтгэл уруу буулгахдаа хам сэдвээр нь болгоомжтой буулгана. Үүний үр дүнд, ангилах үндсэн процессоос өмнө дүрсийн нэгтгэх эгнээнд сэргээн засварлахад тусдаа өөрчлөлт хийх шаардлагатай болно.

G2P sorting cannot be accommodated within the main algorithm of the Unicode Collation Algorithm, as only limited contextual sensitivity (to the extent that adding a few contractions to the collation element table will work, like Thai/Lao CV inversion) can be handled thereby. However, the mapping from graphetic characters to folded character is highly context-sensitive. As a result, a separate step of transformation is needed to reconstruct folded strings before the main procedure of sorting.

Энэхүү хэлэлцүүлэг нь үсгийн нэгтгэх аргыг бүтээж сэргээн засварлах практик арга байгаа гэдгийг харуулах зорилготой юм. Гэсэн хэдий ч үндсэн дүрсээр кодчилох тогтолцоог үгийн дунд ордог хоёр шүд (ᠮ) Х, Г үсгийг төлөөлөх боломжтой эсвэл Г үсэг үгийн дунд орохдоо хоёр салангид (ᠨ) нэг нэг шүд хэлбэрээр бичигдэг.

Дүрсээр кодчлолын аргын дотор бүх хувилбаруудыг багтаахын тулд зөвхөн хамгийн хувирашгүй кодчиллын тогтолцооны хэрэгжилтийг энд авч үзэхэд уншигчид үүнийг бусадтай нь тохируулахад хялбар болдог. Тус хэлэлцүүлэг нь дүрсийн кодчиллын хамрах хүрээ ямар нэгэн тодорхой хувилбаруудыг дэмжиж байгаагаар тайлбарлахгүй гэдгийг онцлон тэмдэглэх нь зүйтэй.

The present discussion is aimed at showing that there are always practical solutions to the reconstruction of folded letters, however graphetically radical the encoding scheme is (as long as the medial double tooth (ᠮ) representing x or g is not broken into two single teeth (ᠨ)). In order to cover all practical variations within the graphetic encoding approach, only the implementation for the most radical graphetic encoding scheme is addressed here, and the readers should find it easy to adjust it to the rest. It should be emphasized that the present discussion shall not be construed as favoring any specific variant in the spectrum of graphetic encoding.

3.1 Ялгах элементийн хүснэгт

Collation element table

G2P-н ялгах элементийн хүснэгтийг Хүснэгт 2 харуулсан болно, Энд нэгтгэсэн үсгүүдийг мөн тэдгээрийг ангилах дарааллаар нь жагсаасан байдал, харгалзах дүрсүүд, үсгийн дүрс, ангилах аргын нэгтгэсэн болон задаргасан байдлын дараалалуудыг тус тус харуулсан. Нэгтгэсэн үсэг, дүрсийн үсгийн хоёулангынх нь онцлогыг буулгахдаа жижиг үсгээр тэмдэглэж оруулсан. Харин дүрсийн үсгээр бичсэн тэмдэгтүүд нь том үсгээр бичигдсэн. Доорхи хүснэгтийг хархад, дүрс кодлох тогтолцооны үндсэн язгуурын чанараас хамааран 9 хүртэлх тоогоор илэрхийлэгддэг график дүрс байж болдог. Тухайлбал A, I, O, U, X, G, L, W, and H гэх мэт.

The collation element table of G2P sorting is given as Table 2, where folded letters are listed by their sorting order, and corresponding glyphs, graphic characters, and unfolded letters along with their orders are given together. Folded letters and graphetic characters in biunique mapping with them are transliterated with small letters, while graphetic characters not in biunique correspondence with folded letters are transliterated with capital letters. As can be seen from the table, there might be up to 9 under determined graphetic characters depending on the radicalness of the graphetic encoding scheme, namely A, I, O, U, X, G, L, W, and H. Secondary weights of collation elements in the table are largely arbitrarily given, as there is no widely accepted secondary weighting convention for either unfolded or folded sorting.

Хүснэгт 2. Ялгах элементийн хүснэгт

Unfolded order	Unfolded letter	Collation element	Folded letter	Graphetic character	Glyph			
					IS	I	M	F
—	—	[.00.3]	ʔ	A	□	ʔ	ʔ	□
1 2 ₁	a e	[.01.1]	a	A	□	ʔ	ʔ	√/η
		[.01.2]	α = α		γ	■	■	■
2 ₂	e	[.02.1]	ø	—	■	■	-	γ
		[.02.2]	ë	W	■	□	ʔ	ʔ
3	i	[.03.1]	ij	II	□	□	π	□
		[.03.2]	i	I	∠	ʔ	ʔ	∠/η
4 5 6 ₁ 7 ₁	o u ö ₁ ü ₁	[.04.1]	o	O	■	σ	σ/σ	σ/σ
		[.04.2]	u	U	θ	■	■	θ
6 ₂ 7 ₂	ö ₂ ü ₂	[.05.1]	ö	OI	■	■	σ/σ	□
		[.05.2]	ü = ü		■	■	■	σ/σ
11	n	[.11.1]	n = n		■	ʔ	ʔ	√.
		[.11.2]	ñ	A	□	□	ʔ	√
12	η	[.12.1]	η	AG	□	□	π	η
13	b	[.13.1]	b = b		■	θ	θ	θ
14	p	[.14.1]	p = p		■	θ	θ	θ
15 ₁	x	[.15.1]	x	X	■	ʔ	π	π
16 ₁	g	[.16.1]	ġ = ġ		■	ʔ	π	π.
		[.16.2]	ġ	X	■	□	π	π
15 ₂ 16 ₂	x ₂ g ₂	[.17.1]	g	G	■	ʔ	ʔ	η
17	m	[.18.1]	m = m		■	ʔ	ʔ	ʔ
18	l	[.19.1]	l	L	■	ʔ	ʔ	ʔ
19	s	[.20.1]	s = s		■	ʔ	ʔ	π
20	š	[.21.1]	š = š		■	ʔ	ʔ	π

21 ₁ 22 ₁	$t_1 d_1$	[.22.1]	$t = t$	■ 𐑃 𐑃 𐑃
21 ₂ 22 ₂	$t_2 d_2$	[.23.1]	$d = d$	■ 𐑄 𐑄 𐑄
		[.23.2]	$\delta \quad \text{OA}$	□ ■ 𐑄 𐑄
23	\check{c}	[.24.1]	$\check{c} = \check{c}$	■ 𐑆 𐑆 𐑆
24	\check{j}	[.25.1]	$\check{j} = \check{j}$	■ ■ 𐑇 𐑇
		[.25.2]	$\check{j} \quad \text{I}$	◁ 𐑇 □ □
25	y	[.26.1]	$y = y$	■ 𐑉 𐑉 ■
		[.26.2]	$\check{y} \quad \text{I}$	□ 𐑉 𐑉 ◁
26	r	[.27.1]	$r = r$	■ 𐑊 𐑊 𐑊
27	w	[.28.1]	$w \quad \text{W}$	■ 𐑋 𐑋 𐑋
		[.28.2]	$\check{w} \quad \text{U}$	□ ■ ■ 𐑋
28	f	[.29.1]	$f = f$	■ 𐑌 𐑌 𐑌
29	k	[.30.1]	$k = k$	■ 𐑍 𐑍 𐑍
30	c	[.31.1]	$c = c$	■ 𐑎 𐑎 𐑎
31	z	[.32.1]	$z = z$	■ 𐑏 𐑏 𐑏
32	h	[.33.1]	$h \quad \text{H}$	■ □ 𐑑 𐑑
		[.33.2]	$\check{h} \quad \text{AH}$	■ 𐑑 □ □
33	\check{r}	[.34.1]	$\check{r} = \check{r}$	■ 𐑒 ■ ■
34	\check{t}	[.35.1]	$\check{t} \quad \text{LH}$	■ 𐑓 𐑓 □
35	\check{z}	[.36.1]	$\check{z} \quad \text{H}$	■ 𐑔 □ □
36	\check{c}	[.37.1]	$\check{c} \quad \text{OO}$	■ 𐑕 □ ■

3.2 График дүрслэлээс нэгтэгсэн үсгийг дахин сэргээх

Reconstructing folded letters from graphetic characters

Аз болоход, энгийн орлуулалтын програмчилал нь(хэсэгчилэхгүй, давталтгүйгээр) бидний хүсэн тэмүүлж байгаа үсгийг дахин ангилахад илүү сайн хангаж өгч байна. Мөн толь бичгийн зүйлийг агуулаагүй үг хэллэгээр орлуулсан багц орлуулалтын талаар доорхи 3-х хүснэгтэд жагсааж бичсэн. Өөрчлөлтийг бүрэн дүүрэн болгох бүх орлуулалтыг тодруулсан болно. Бүлгүүдийг (...) эсвэл тоон дайвар үгийг (*, +, ?, {m, n}, гэх мэт) энд эрэмбэлж оруулаагүйг дурдах нь зүйтэй байх.

Luckily, finite unilinear (no branching, no looping) application of regular expression substitutions will suffice for our purposes of folded letter reconstruction. A set of substitutions containing no dictionary items for content words is adopted, as is listed in Table 3. Catch-all substitutions that ensure completeness of transformation are highlighted. It is worth mentioning that devices like groups (...) or quantifiers (*, +, ?, {m,n}, etc.) are not even resorted to here.

Хүснэгт 3: Нэгтэгсэн үсгийн сэргээн засварлахад хэрэглэдэг орлуулалтын жагсаалт

Table 3 List of regex substitutions employed in folded letter reconstruction

Search	Subs.
<code>\bOOA\b</code>	ođ
<code>\bI\b</code>	i
<code>\bIIAA\b</code>	iĭań
<code>\bIIAr\b</code>	iĭar
<code>\bOOI</code>	ĉi
<code>\bH</code>	ž

	\bAH	h
	\bLH	l
	\bW	w
	\bI	i
	I (?=[α])	ı
	U (?=[α])	ü
	U	u
	H	h
	L	l
	OI\b	oi
	OII\b	öi
	(?<=\bG) OII	öi
	OII	oj
	OI	ö
	(?<=\bA) II\b	ıi
	(?<!I) I (?!I)	i
	\bO	o
	O (?!A)	o
	(?<!A) G	g
	(?<=\b[AWO]) II	ij
	I	i
	OAG (?! [AIOαüaøëiijouö])	oŋ
	(?<=\b[AIOαüaøëiijouö]) W	w
	(?<=\b[AWIαüaøëiijouö]) OA	ö
	(?<!A) G	g
	O	o
	X (?=[üAIαüaøëiijouö])	x
	X	ğ
	GW (?=[nbpmsštdčjrfkczgljwhgyř?xiüşžčhı])	gë
	G (?=[AIαüaøëiijouö])	g
	(?<=[nbpmsštdčjrfkczgljwhgyř?xiüşžčhı]) AAG	aŋ
	WAAG	waŋ
	(?<=[nbpmsštdčjrfkczgljwhgyř?xiüşžčhı]) AG	ag
	(?<=[Iαüaøëiijouö]) AG	ŋ
	WAG	ëŋ
	(?<=[nbpmsštdčjrfkczgljwhgyř?xiüşžčhı]) AA	ań
	WAA	wań
	(?<=[nbpmsštdčjrfkczgljwhgyř?xiüşžčhı]) A	a
	(?<=[Iαüaøëiijouö]) A	ń
	(?<=\bA) WA (?! [ńŋğđαüaøëiijouö])	ěń
	(?<=[nbpmsštdčjrfkczgljwhgyř?xiüşžčhı]) WA (?! [ńŋğđαüaøëiijouö])	ěń
	WA	wa
	AAAG	?aŋ
	AAG	?ğŋ
	AG	?ğg
	G	g
	AAA	?ań
	AA	?a
	A (?=[Wαüaøëiijouö])	?
	A	?ğ
	W (?=[ı] \b)	w
	W (?=[ijńŋğđ])	ë

	W (?=[u])	ë
	(?<=[ʔhɪ]) W	ë
	(?<=[Iɑüaøëiijjouöńǵǫ]) W	w
	WW	ëw
	W	ë

4. G2P ангилалыг туршиж үзэх.

Testing G2P sorting

Өмнөх хэсэгт тодорхойлсон G2P ангилах аргыг 26433-үг зөв бичгийн өгөгдлийн сантай тулган шалгадаг. Үр дүн нь дараах байдалтай байна.

G2P sorting method specified in the previous section is tested against a 26433-word spelling database. The results are as follows.

4.1 Нэгтэгсэн үсгийг буруу сэргээсэн

Misreconstruction of folded letters

График дүрслэлээс нэгтэгсэн үсгийг дахин сэргээх мэдээлэл нь анхны задарсан үсэгтэй маш сайн нийцэж байгааг 2-р хүснэгтээс харж болох хэдий ч, 4-р хүснэгтэд үл хамаарах зүйлүүдийг жагсааж бичсэн.

The folded letters reconstructed from graphetically encoded data matched the original unfolded letters quite well in the correspondences defined in Table 2, except for the exceptions listed in Table 4 below.

Хүснэгт 4. Буруу сэргээгдсэн төрлүүд

Table 4 Types of misreconstruction

Төрлүүд Type	Задарсан үсэг Unfolded letter	Нэгтэгсэн үсэг Folded letter	Нэгтгэн дахин сэргээсэн үсэг Reconstructed folded letter	График үсэг Graphetic character	Ерөнхий үгийн гарал зүй Major etymological class	Тоо Count
A	e	ʔø	ʔ	A	Уугуул Native	27/26433
B	Un	oń	ǫ	OA	Уугуул/Зээлсэн Native/Loan	12/26433
C	w ë	W ë	Ë w	W	Зээлсэн Loan	6/26433
D	Алдаатай зөв бичигдсэн Ill-formed spelling				Уугуул Native	2/26433

Хүснэгтээс хархад буруу сэргээгдсэн дөрвөн үндсэн төрөл байгааг харж болно: Үүнд:

As can be seen from the table, there are four major types of misreconstruction:

- Төрөл A: *ei* эсвэл *eü* үсгээр эхэлсэн бүх үг болон *en* эхэлсэн бүх үгнүүдийн дараа эгшиг дагаагүйг харж байна. Type A: All words beginning with *ei* or *eü*, and all words beginning with *en* which is not followed by a vowel.
- Төрөл B: Бүх үг *on/un/ö n/ü n* эдгээрийг багтаасан байдаг ба эдгээр нь эгшигийн өмнө орсон байх ба эгшигийн ард дагаж бичигддэггүй. Ихэнхдээ харьяалахын тийн ялгалд хамааралтай хэдий ч, заримдаа гадаад үгнээс зээлж авсан үгэнд мөн харагддаг.

Type B: All words containing *on/un/ö n/ü n* which is preceded by a vowel but not followed by a vowel. Mostly in genitive stems of fugitive-n native words (e.g., EREÜ: nom. ereü, gen. ereü n=ü), but may also occur in loan words.

- Төрөл C: Цөөн хэдэн үгэнд *w* эсвэл *ë* орсон байдаг. Нөлөөлөлд өртсөн үгсийн жагсаалт сонгосон хөрвөх зарчмаас хамааран өөр өөр байж болдог.
Type C: A few words containing *w* or *ë*. The list of affected words may vary, depending on the specific transformational rule set chosen.
- Type D: Маш цөөн хэдхэн үг алдаатай зөв бичигдсэн байдаг байна.
A few words with ill-formed spelling.

Энэ буруу сэргээх гэдэг нь ерөнхийдөө урьдчилан таамаглах боломжтой гэсэн үг юм. 5-р хүснэгтэнд буруу сэргээгдсэн төрлүүдийн талаарх жишээнүүдийг харуулсан байгаа.

This means that cases of misreconstruction are predictable in general. Table 5 gives some examples of misreconstruction of each type.

Хүснэгт 5.

Table 5

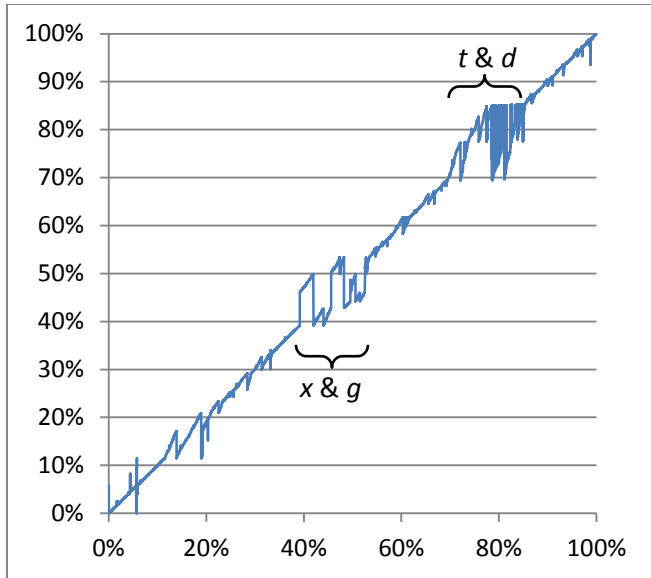
Төрлүүд Type	Задарсан үсэг Unfolded letter	Нэгтгэсэн үсэг Folded letter	Нэгтгэн дахин сэргээсэн үсэг Reconstructed folded letter	График үсэг Graphetic character
A	ei ein eü ende	ʔɛi ʔɛijn ʔɛu ʔɛnda	ʔi ʔiin ʔu ʔada	AI AIIA AU AAAdA
B	ereün ondoön	ʔɛraoń ʔońdooń	ʔɛrað ʔońdoð	ArAOA AOAdOOA
C	niswanis nirwalaxu burwasad yêlwi	niswanis nirwalaxu borwasað yêlwi	nisêńnis nirêńlaxu borêńsað yêlêi	nIsWANIs nIrWALAXU bOrWAsAOA yWLWI
D	tɯri tɯrilig	tɯri tɯrilig	tagri tagrilig	tAGrI tAGrILIG

26433-үгийн мэдээллийн санд, зөвхөн 47 буруу сэргээлтэй оролтын тохиолдол байдаг хэдий ч нийт 0.17% -ийг эзэлж байна.

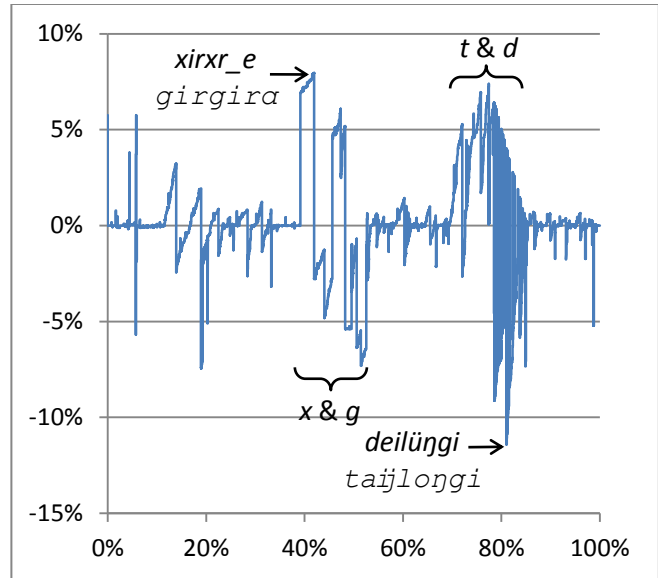
In the 26433-word database, only 47 entries of misreconstruction occur, making up 0.17% of the total.

4.2 G2P ангилалыг хийсвэр задалдаг ангилалтай харьцуулах нь Comparison of G2P sorting with ideal unfolded sorting

G2P ангилалын үр дүнг хийсвэр задалдаг ангилалтай харьцуулан хархад, Үр дүнгийн зураг 1 ба зураг 2-т харьцуулан харуулав. Үүнд:



Зураг 1: G2P ангилалын хийсвэрээр задалдаг ангилал
Figure 1 G2P sorting ~ ideal unfolded sorting



Зураг 2: G2P ангилалын чиглэл, сааталыг хийсвэр задалддаг ангилалтай харьцуулах
Figure 2 Lead and lags of G2P sorting relative to ideal unfolded sorting

Зураг 1 дээр G2P ангилалын дараалалыг хийсвэрээр задалдаг ангилалтай хувиар харуулав. Зураг 2 дээр ялгаануудыг харьцуулан харуулав. G2P-ийн хамгийн дээд эрэмбийн ялгаа нь -11.41% -иас ялимгүй ялгаатай, *deilüngi*-н хувьд ангилахдаа *taijlongi* хэлбэрээр хамгийн их хоцролт + 7.96% ангилсан байгааг харж байна. Харин *xirxr_e*-н хувьд ангилахдаа *girgira* хэлбэрээр ангилсан байна. Бүх 26433 оруулгуудын ялгах дарааллын стандарт хазайлт (σ) нь 2.84% байдаг.

Figure 1 shows the relative order of G2P sorting against ideal unfolded sorting in percentage, and Figure 2 shows their difference. The maximum lead of G2P sorting to ideal unfolded sorting is -11.41%, occurring at *deilüngi* sorted as *taijlongi*, while the maximum lag is +7.96%, occurring at *xirxr_e* sorted as *girgira*. The standard deviation (σ) of the sorting order differences of all 26433 entries is 2.84%.

Ангилалын ялгааны гол хувь нэмэр нь эхний үсэгний *x / g* ба *t / d* хосуудаас бүрдэх ба эхний эгшигний */ e, o / u* ба *ö / ü* хосуудаас ихээхэн хамаардаг. Үсгийг задлан ангилах аргачлалд дассан хэрэглэгч нь G2P-ыг ямар нэгэн сургалтгүйгээр дасан зохицоход хялбар байдаг, зөвхөн энэ орооцолдсон хос бичсэн хэлбэрээрээ тохируулахыг санаарай.

The major contributions to the sorting differences come from *x/g* and *t/d* pairs of the first letter, and less significantly from *a/e*, *o/u* and *ö /ü* pairs of the first vowel. A user accustomed to ideal unfolded sorting should find it easy to adapt to G2P sorting without much training, just to remember to make some adjustment to these tangled pairs according to their written shapes.

Адил төстэй байдлыг хялбархан ангилах үүднээс ялгах аргын үр дүнг 6-р хүснэгтэд үзүүлэв. G2P ангилалын мэдээлэлийн санд G2 үсэг задарч байгаа тэр хэсэгт жигд хэлбэрээр явагддаг. Нийт 52 өөр оруулгаас зөвхөн 5 (тодруулсан) оруулга боломжтой гэсэн үр дүн гарч байна. Эцэст нь хэлхэд, G2P ангилал нь ямар нэгэн байдлаар хазаахгүй хамгийн тохиромжтой нь гэж үзэж болно

To give an intuitive view of the sorting similarity, an extract of the results of both sorting methods are given in Table 6, where the items are extracted in a uniformly spaced manner from the G2P-sorted database. The results show that only 5 entries (highlighted) out of total 52 differ. In conclusion, G2P sorting does not deviate as much from ideal unfolded sorting as one might expect.

6-р хүснэгт: G2P ангилалын задаргаа

Table 6 Extract of G2P sorting and ideal unfolded sorting

Rel. delta sort. ord.	Unfolded sort. ord	Unfolded letter	G2P sort. ord.	Reconstructed folded letter	Graphetic character
0.06%	517	agawa	500	᠎agawa	AAgAWA
0.00%	1000	asagtuxu	1000	᠎asa᠑doxu	AAsAXdOXU
-0.05%	1488	arčigur	1500	᠎arčigor	AArčI᠑Or
0.06%	2015	elegdexü	2000	᠎elagdago	ALAGdAGO
0.04%	2511	iin	2500	᠎iin	AIIA
0.01%	3003	irgagalaxu	3000	᠎irgagalaxu	AIrgAgALAXU
-0.91%	3259	obosxixü	3500	᠎obosgigo	AObOsGIGO
0.62%	4164	usadxagči	4000	᠎osa᠔xa᠑či	AOSAOAXAXčI
-0.14%	4462	ursigtai	4500	᠎orsig᠑dai	AOrsIXdAI
	(5000	örgedxel)			
1.04%	5276	ülemjidexü	5000	᠎ölamjidago	AOILAmjIdAGO
-1.89%	5000	örgedxel	5500	᠎örga᠔gal	AOIrGAOAGAL
-0.63%	5834	namči	6000	namči	nAmčI
-0.09%	6475	noxaituxu	6500	noxaijdoxu	nOXAIIdOXU
-0.01%	6997	bagaturčud	7000	ba᠑adorčo᠔	bAgAdOrčOOA
0.71%	7689	beyeči	7500	bayači	bAyAčI
1.23%	8326	buurji	8000	boorji	bOOrjI
0.25%	8565	bujigirtuxu	8500	bojigirdoxu	bOjIGIrdOXU
0.02%	9005	bürxüixü	9000	börgoi᠑go	bOIrGOIIGO
0.00%	9500	xabursil	9500	xaborsil	XAbOrsIL
0.00%	10000	xasumal	10000	xasomal	XAsOmAL
2.56%	11179	xoordalg_a	10500	xoordal᠑a	XOOrdAL᠑a
1.82%	11482	xo᠑o᠑čilaxu	11000	xo᠑o᠑čilaxu	XO᠑OAGčilAXU
	(12250	xög)			
5.39%	12928	gagaglaxu	11500	᠑a᠑a᠑laxu	᠑AgAXLAXU
6.75%	13788	gutugaxu	12000	᠑odo᠑axu	᠑OdO᠑AXU
2.63%	13198	gemsixü	12500	gamsigo	GAmSIIGO
1.20%	13319	gilaljam_a	13000	gilal᠑ama	GILALjAma
	(13788	gutugaxu)			
-4.72%	12250	xög	13500	᠑ög	GOIG
0.45%	14120	güyüldüxü	14000	᠑öyoldogo	GOlyOLdOGO
0.34%	14590	mesil	14500	masil	mAsIL
-0.01%	14997	mönxexlexü	15000	mön᠑alago	mOIAGGALAGO
-0.11%	15471	sanagatai	15500	sana᠑adai	sAnAgAdAI
-0.80%	15789	salbarxai	16000	salbarxai	sALbArXAI
0.05%	16514	sibeg	16500	sibag	sIbAG
-0.03%	16993	siratuxu	17000	siradoxu	sIrAdOXU
-0.82%	17283	soyoxai	17500	soyoxai	sOyOXAI
-0.02%	17995	šaliyatuxu	18000	šaliyadoxu	šALIyAdOXU
-0.23	18440	taitarar	18500	taijda᠑ar	tAIIdAgAr
0.86%	19227	tebxe	19000	tabga	tAbGA
6.16%	21133	dagir	19500	tagir	tAGIr
-4.21%	18884	tasixu	20000	tasixu	tAsIXU

-5.30%	19097	tawarčilaxu	20500	tawarčilaxu	tAWArčILAXU
	(19227	tebxe)			
-4.68%	19761	togonočaxu	21000	toġonočaxu	tOġOnOčAXU
-4.62%	20275	tulgagurtai	21500	tolġaġordai	tOLġAġOrdAI
	(21133	dagir)			
2.03%	22537	dünjsüixü	22000	tönsoijgo	tOIAGsOIIGO
0.39%	22604	dürsütei	22500	törsodai	tOIRSodAI
-0.48%	22874	časuraxu	23000	časoraxu	čAsOrAXU
0.00%	23500	čisorxau	23500	čisorxau	čIsOrXAU
0.35%	24092	čüüreljexü	24000	čöoraļjago	čOIOrALjAGO
0.87%	24730	ĵegülte	24500	ĵagolda	IAGOLdA
-0.01%	24998	ĵilabči	25000	ĵilabči	IILAbči
-0.46%	25378	ĵoltai	25500	ĵoldai	IOLdAI
0.43%	26114	yexemsüg	26000	yagamsog	yAGAmSOG

5 ДҮГНЭЛТ

5 Summary

Тус баримт бичигт тусгагдсан гол чухал асуудлууд нь:

The key points of this document are:

- Нэгтгэн ангилах аргачилал нь шинэ нээлт биш, мөн энэ нь задлан ангилалд муу үр дагавар үзүүлдэггүй.
Folded sorting is not a new invention, nor is it inferior to unfolded sorting.
- G2P ангилах аргачилал, нэгтгэн ангилах хувилбар болох нь графикаар кодчилох шаардлагатай байгаа Монгол хэллэгийг автоматаар ангилах хамгийн сайн шийдэл байж болох юм.
G2P sorting, as a variant of folded sorting, can be a good solution when automated sorting of graphetically encoded Mongolian is needed.
- G2P-ийн ялгах ангилалын үр дүн нь өргөн хэрэглэгддэг задаргаатай ангиллуудын хувьд ялгаатай биш бөгөөд хэрэглэгчид үүнийг хялбархан ашиглах боломжтой юм.
The results of G2P sorting are not dissimilar to those of widely used unfolded sorting, and users will readily get used to it.

A. Толь бичгийн толгойн дарааллыг эрэмбэлэх аргачилал болон G2P ангилах аргачилал

A. Ordering of dictionary headings in ideal unfolded sorting and G2P sorting

7-р хүснэгтэд орчин үеийн Хятад дахь хамгийн өргөн дэлгэрсэн хувилбар болох толь бичгийн гарчигыг нэлээд задлан ялгаж харуулав. Доорхи хүснэгтэнд байгаа нэг эгнээнд онцлон тэмдэглэсэн үсэгнүүд нь адил бичигддэг өөр утгатай үсэгнүүд юм.

Table 7 shows an ideal unfolded sorting of dictionary headings, which is the most widely accepted version in modern China. Headings in adjacent cells with the same underline patterns are homographic.

Хүснэгт 7 Толь бичгийн гарчигийг ангилах

Table 7 Ideal unfolded sorting of dictionary headings

	<i>a</i>	< ₁	<i>e</i> < ₂ <i>ě</i>	< ₁	<i>i</i>	< ₁	<i>o</i>	< ₁	<i>u</i>	< ₁	<i>ö</i>	< ₁	<i>ü</i>
< ₁	<i>na</i>	< ₁	<i>ne</i> < ₂ <i>ně</i>	< ₁	<i>ni</i>	< ₁	<i>no</i>	< ₁	<i>nu</i>	< ₁	<i>nö</i>	< ₁	<i>nü</i>
< ₁	<i>ba</i>	< ₁	<i>be</i> < ₂ <i>bě</i>	< ₁	<i>bi</i>	< ₁	<i>bo</i>	< ₁	<i>bu</i>	< ₁	<i>bö</i>	< ₁	<i>bü</i>
< ₁	<i>xa</i>	< ₁	<i>xe</i> < ₂ <i>xě</i>	< ₁	<i>xi</i>	< ₁	<i>xo</i>	< ₁	<i>xu</i>	< ₁	<i>xö</i>	< ₁	<i>xü</i>
< ₁	<i>ga</i>	< ₁	<i>ge</i> < ₂ <i>gě</i>	< ₁	<i>gi</i>	< ₁	<i>go</i>	< ₁	<i>gu</i>	< ₁	<i>gö</i>	< ₁	<i>gü</i>
< ₁	<i>ta</i>	< ₁	<i>te</i> < ₂ <i>tě</i>	< ₁	<i>ti</i>	< ₁	<i>to</i>	< ₁	<i>tu</i>	< ₁	<i>tö</i>	< ₁	<i>tü</i>
< ₁	<i>d'a</i>	< ₁	<i>d'e</i> < ₂ <i>d'ě</i>	< ₁	<i>d'i</i>	< ₁	<i>d'o</i>	< ₁	<i>d'u</i>	< ₁	<i>d'ö</i>	< ₁	<i>d'ü</i>
< ₁	< ₂ <i>da</i>	< ₁	< ₂ <i>de</i> < ₂ <i>dě</i>	< ₁	< ₂ <i>di</i>	< ₁	< ₂ <i>do</i>	< ₁	< ₂ <i>du</i>	< ₁	< ₂ <i>dö</i>	< ₁	< ₂ <i>dü</i>

8-р хүснэгтэд G2P ангилалын толь бичгийн гарчгийг үзүүлэв. "≡" -н хоёр тал дээрх гарчиг нь адил бичигддэг өөр утгатай үсэг. "⊃" гэдэг тэмдэглэгээ нь нэмэлт таамаглалыг илэрхийлнэ. Бага хаалтад байгаа тэмдэглэгээнүүд хийсвэр нэгтгэсэн ангилалыг илэрхийлнэ. Жнь: **en** гэдэг тэмдэг G2P ангилалын оронд хийсвэр нэгтгэсэн ангилалын **e** тэмдэглэгээний дотор гарчиглан ангилагдаж орно.

Table 8 shows G2P sorting of dictionary headings. Headings on both sides of "≡" s are homographic. "⊃" s indicate additional subsumptions, and the parenthesized terms should be in their normal places in ideal folded sorting. For example, **en** should be sorted under the heading of **e** in ideal folded sorting, instead of under **a** in G2P sorting.

Хүснэгт 8: G2P ангилалын толь бичигийн гарчиг аргачилал

Table 8 G2P sorting of dictionary headings

$a (\supset en)$	$<_1$	$e <_2 \ddot{e}$	$<_1$	$i (\supset ei)$	$<_1$	$o \equiv u (\supset \ddot{u}i)$	$<_1$	$\ddot{o} \equiv \ddot{u}$
$<_1 na \equiv ne$	$<_1$	$n\ddot{e}$	$<_1$	ni	$<_1$	$no \equiv nu (\supset n\ddot{u}i)$	$<_1$	$n\ddot{o} \equiv n\ddot{u}$
$<_1 ba \equiv be$	$<_1$	$b\ddot{e}$	$<_1$	bi	$<_1$	$bo \equiv bu (\supset b\ddot{u}i)$	$<_1$	$b\ddot{o} \equiv b\ddot{u}$
⋮								
$<_1 xa$					$<_1$	$xo \equiv xu$		
$<_1 ga$					$<_1$	$go \equiv gu$		
$<_1 xe$	$<_1$	$g\ddot{e}$	$<_1$	xi				$<_1$
$\equiv ge$				$\equiv gi$				$\equiv x\ddot{o} \equiv x\ddot{u}$
⋮								
$<_1 ta \equiv te$	$<_1$	$t\ddot{e}$	$<_1$	ti	$<_1$	$to \equiv tu (\supset t\ddot{u}i)$	$<_1$	$t\ddot{o} \equiv t\ddot{u}$
$\equiv da \equiv de$		$\equiv d\ddot{e}$		$\equiv di$		$\equiv do \equiv du (\supset d\ddot{u}i)$		$\equiv d\ddot{o} \equiv d\ddot{u}$
$<_1 d'a \equiv d'e$	$<_1$	$d'\ddot{e}$	$<_1$	$d'i$	$<_1$	$d'o \equiv d'u (\supset d'\ddot{u}i)$	$<_1$	$d'\ddot{o} \equiv d'\ddot{u}$
⋮								

Толь бичгийн гарчигийг Худам бичгийн задарсан болон нэгтгэгсэн үсэгний аль алинаар жагсаасан хүснэгтийг энд тайлбар өгч ойлгомжтой байхын тулд харууллаа .

Tables of dictionary headings in Hudum listed by both unfolded and folded letters are given here for the convenience of reference.

Хүснэгт 9: Худам толь бичгийн гарчиг

Table 9 Dictionary headings in Hudum

(a) Задарсан үсгээр жагсаасан хэлбэр

(a) Listed by unfolded letters

	-a	-e	-ë	-i	-o/u	-ö/ü
-	ᠠ	ᠡ	ᠢ	ᠣ	ᠤ	ᠥ
n-	ᠨᠠ	ᠨᠡ	ᠨᠢ	ᠨᠣ	ᠨᠤ	ᠨᠥ
b-	ᠪᠠ	ᠪᠡ	ᠪᠢ	ᠪᠣ	ᠪᠤ	ᠪᠥ
⋮	⋮					
x-	ᠬᠠ	ᠬᠡ	ᠬᠢ	ᠬᠣ	ᠬᠤ	ᠬᠥ
g-	ᠭᠠ	ᠭᠡ	ᠭᠢ	ᠭᠣ	ᠭᠤ	ᠭᠥ
⋮	⋮					
t/d-	ᠲᠠ	ᠲᠡ	ᠲᠢ	ᠲᠣ	ᠲᠤ	ᠲᠥ
d'-	ᠳᠠ	ᠳᠡ	ᠳᠢ	ᠳᠣ	ᠳᠤ	ᠳᠥ
⋮	⋮					

(б) Нэгтгэгсэн үсгээр жагсаасан хэлбэр

(b) Listed by folded letters

	-a	-ø	-ë	-i	-o	-ö
ᠲ-	ᠠ	ᠡ	ᠢ	ᠣ	ᠤ	ᠥ
n-	ᠨᠠ	■	ᠨᠢ	ᠨᠣ	ᠨᠤ	ᠨᠥ
b-	ᠪᠠ	■	ᠪᠢ	ᠪᠣ	ᠪᠤ	ᠪᠥ
⋮	⋮					
x-	ᠬᠠ	■	ᠬᠢ	ᠬᠣ	ᠬᠤ	ᠬᠥ
ᠭ-	ᠭᠠ	■	ᠭᠢ	ᠭᠣ	ᠭᠤ	ᠭᠥ
g-	ᠭᠠ	■	ᠭᠢ	ᠭᠣ	ᠭᠤ	ᠭᠥ
⋮	⋮					
t-	ᠲᠠ	■	ᠲᠢ	ᠲᠣ	ᠲᠤ	ᠲᠥ
d-	ᠳᠠ	■	ᠳᠢ	ᠳᠣ	ᠳᠤ	ᠳᠥ
⋮	⋮					

Б. График үсэг болон нэгтгэгсэн үсгийн хоорондын уялдаа холбоо

B. Correspondence between graphetic characters and folded letters

График болон нэгтгэгсэн үсгүүдийн хоорондын уялдаа холбоог ойлгоход хялбар болгох зорилгоор Хүснэгт 10 орууллаа . Нэг үсэг болон үсгийн нийлбэрүүдийг хүснэгтэд эгнүүлж бичиж оруулсан байна.

Correspondence between graphetic characters and folded letters is given as Table 10 for the convenience of reference. Single characters and character combinations are aligned in the table.

Хүснэгт 10: График болон нэгтэгсэн үсгийн хоорондын үчлдаа холбоо

Table 10 Correspondence between graphetic characters and folded letter

Glyph				Graphetic character	Folded letter	Unfolded letter	Glyph				Graphetic character	Folded letter	Unfolded letter
IS	I	M	F				IS	I	M	F			
ᠬ	ᠰ	ᠷ	ᠠ/ᠨ	A	ᠠ	—	□	□	ᠠ	ᠨ	AG	ᠠ	ᠨ
					ᠠᠶ	<i>e</i>	■	ᠠᠶ	□	□	AH	ᠠ	<i>h</i>
					ᠠ	<i>a</i>	□	■	ᠠ	ᠠ	OA	ᠠ	<i>d</i>
					ᠠ	<i>n</i>	□	□	ᠠ	□	II	ᠠ	<i>i</i>
ᠢ	ᠰ	ᠷ	ᠠ/ᠨ	I	ᠢ	<i>i</i>	■	■	ᠢ/ᠢ	□	OI	ᠢ	<i>ö/ü</i>
					ᠢ	<i>j</i>	□	■	ᠢ	ᠠ	OA	ᠢ	<i>d</i>
					ᠢ	<i>y</i>	■	■	ᠢ/ᠢ	□	OI	ᠢ	<i>ö/ü</i>
■	ᠢ	ᠢ/ᠢ	ᠠ/ᠨ	O	ᠠ	<i>o/u/ö/ü</i>	□	■	ᠢ	ᠠ	OO	ᠢ	<i>ĉ</i>
ᠢ	■	■	ᠢ	U	ᠢ	<i>u</i>							
					ᠢ	<i>ü</i>							
					ᠢ	<i>w</i>							
■	ᠠ	ᠠ	ᠠ	X	ᠠ	<i>x</i>							
					ᠠ	<i>g</i>							
■	ᠢ	ᠢ	ᠠ	G	ᠢ	<i>g</i>	□	□	ᠢ	ᠠ	AG	ᠢ	ᠠ
■	ᠠ	ᠠ	ᠠ	L	ᠠ	<i>l</i>	■	ᠠᠶ	ᠠᠶ	□	LH	ᠠ	ᠠ
■	ᠢ	ᠢ	ᠠ	W	ᠢ	<i>e</i>							
					ᠢ	<i>w</i>							
■	ᠠ	ᠠ	ᠠ	H	ᠠ	<i>h</i>	■	ᠠᠶ	□	□	AH	ᠠ	<i>h</i>
					ᠠ	<i>z</i>	■	ᠠᠶ	ᠠᠶ	□	LH	ᠠ	ᠠ