# The order of components in Khmer orthographical syllables[1]

In Unicode, letters or 'characters' should be typed in the order of their pronunciation in a word or syllable, regardless of the visual order of elements. The reordering and rearranging should be done by code in the font, the software and/or the operating system. Any one syllable[2] should be typed in only one single way, otherwise a text may not be found using search-functions or words may be sorted unexpectedly or wrongly by software. In order to assist the user to type syllables always the same way, the Windows rendering system (the program that assembles the text input into its visual result) will insert dotted circles in text if elements/characters are typed in non-standard order. However, other systems using the same fonts (e.g. Adobe software or Linux) may implement the insertion of dotted circles either not at all (Adobe[3]) or more or less differently from Windows[4]. The Apple rendering system does not insert dotted circles at all, but leaves this to the font developers to implement. Every font developer can decide to implement such an insertion or not, and if yes, how to implement it. Every font can differ in the rules applied to the insertion of dotted circles and in consequence a text typed with one font may no longer display correctly in another. This is contrary to the goals of Unicode. Unfortunately some OS X fonts do not help the user to type Khmer Unicode text by inserting dotted circles at all. Text typed in systems not inserting dotted circles will most likely by difficult to read in systems that do or the text may even become nearly incomprehensible.

The Mondulkiri fonts do try to help the user by inserting dotted circles whenever text does not comply with Khmer text rules in OS X. Therefore text typed with Mondulkiri fonts in OS X and displaying correctly, will also display correctly in Windows. However, Windows does permit some character sequences that should not be permitted[5]. This may lead to some dotted circles being shown in text typed in Windows and viewed in OS X with Mondulkiri fonts. In software that does support typographical features in OS X, the insertion of dotted circles can be switched off with Mondulkiri fonts[6].

This order of elements in a syllable should be the same in all implementations of the Khmer script. However, this is not the case. The order of elements as it should be is described in 'The Unicode Standard, Version 6.1 (short: TUS) Core Specifications, chapter

---

[1] Technical terms in this document mostly follow "The Unicode Standard 6.1" (TUS). However, what is called 'consonant shifter' in TUS is called 'register shifter' (a well established linguistic term) in this document.

[2] 'Syllable' here always refers to the orthographical syllable as defined in the formula below. It differs from the linguistic (or phonological) syllable.

[3] Adobe software does use the OpenType rendering rules that are part of a font, but their rendering engine differs from the normal Windows rendering engine and does not insert dotted circles at all. Therefore text typed in Adobe suite applications may not display properly if exported to Windows programs. However, text typed in Windows programs other than Adobe suite software will mostly display correctly when exported to an Adobe suite application. For limitations to this in Adobe CS5 and CS5.5 and ways to overcome these limitations, please see the document on Stylistic Sets in Mondulkiri fonts.

[4] Please see also the document 'Differences in the Implementation of Khmer Unicode'.

[5] Microsoft gives some information on 'Developing OpenType Fonts for Khmer Script' and also gives a definition of the Khmer syllable (http://www.microsoft.com/typography/otfntdev/khmerot/shaping.aspx), but that definition does not quite agree with TUS. It also does not correctly reflect any real OpenType implementation on Windows, at least not after 2003.

[6] This can be done using 'Typography options.' Please see document 'Typography options in Mondulkiri fonts.'

11.4[7]'. However, even this standard is lacking in a few details in order to provide the best possible support for font and software developers. Please find below a description of the orthographical syllable in Mondulkiri fonts and its deviation from TUS.

## Definition of the orthographical syllable in Mondulkiri fonts

## B {{{Z$_1$} S} or R} C {{Z$_1$} S} C {{Z$_1$} S} {{Z$_2$} V} NS C SS

**Whereby:**

B  = base-character. The base character is obligatory and it is the only obligatory component of the syllable. It can be any consonant or independent vowel.

Z$_1$ = zero-width-non-joiner (U+200C, short: zwnj) **or** zero-width-joiner (U+200D, short: zwj)

S  = register shifter (ㅤ ㅤ)

R  = robat (ㅤ)

C  = coeng (or 'subscript consonant'), i.e. a sequence of one coeng-character (ㅤ, U+17D2) and one base-character.

Z$_2$ = zero-width-non-joiner **or** zero-width-joiner

V  = vowel (ា ិ ី ឹ ឺ ុ ូ ួ ើ ឿ ៀ េ ែ ៃ ោ ៅ).

NS = non-spacing-symbol (ំ ់ ៉ ៊ ៍ ៎ ៏ ័)

SS = spacing-symbol (ះ ៈ)

**Additional rules:**

1. There may be only one of each element in the formula per occurrence[8].

2. ប៉ុ should be typed as ប ㅤ ㅤ and not as ប ㅤ ㅤ.

3. There may only be one register-shifter per syllable.

4. Theoretically there could be up to three coengs in a syllable, two before and one after a vowel. Syllables without a vowel can have a maximum of 2 coengs.

5. There may not be more than two zwnj or zwj in a syllable (one relating to a register shifter and one to a vowel).

6. None of the elements apart from the base-character can stand on its own without a preceding base-character.

This means also that Robat is only permitted directly following the base character and cannot occur in the same syllable with a register-shifter.

---

[7] http://www.unicode.org/versions/Unicode6.1.0/ch11.pdf

[8] In the list of Khmer syllables no syllable has more than one non-spacing symbol. However, in the body of the Chuon Nath dictionary there are examples of up to two Khmer non-spacing symbols and a few examples of 2 Khmer non-spacing symbols plus one Thai non-spacing symbol, making a total of three superscript symbols in one syllable. There is also a combination of Toandakhiat with a superscript vowel above it (usually only permitted in the opposite order) and a double E vowel related to Thai based words using Khmer base consonants. In the OS X implementation of this font currently only one non-spacing symbol can be used per syllable. The Windows implementation supports an unlimited number of such symbols.

## Notes

**Register shifters** should be placed **directly after** the base-character or coeng that they are to modify (but see further discussion below). Only then can they be correctly rendered as subscript (◌្) according to the context. Please see below for a discussion of the rules for the rendering as subscript. The rules could possibly be made more stringent by not permitting register shifters to occur before coengs with the exception of Muusikatoan (and possibly also Triisap in the same context for the Tampuan language) if preceded by Ba and followed by a spacing coeng.

**Zero-width-non-joiner (zwnj)** can be inserted directly **before register-shifters** in order to prevent them from being rendered as subscript (◌្). This is necessary for some cases where the spelling deviates from the general rule, sometimes also for personal preference because both variants with the register-shifter as superscript or as subscript may be permissible.

**Zero-width-non-joiner** can also be inserted directly **before vowels** in order to prevent the formation of ligatures between the base-character and the vowel (TUS p382, 'Ligature Control'). While this is mentioned in TUS only for superscript vowels in 'aksaa muul' fonts, it should also be permitted before the vowels Aa, Oo and Au (◌ា ◌េ◌ា ◌េ◌ៅ) in order to prevent the formation of a ligature with the base character. Unfortunately newer implementations of Khmer Unicode in Windows permit zwnj only before superscript vowels and no longer permit the placement of zwnj before Aa, Oo and Au. The Mondulkiri fonts will break these ligatures if there is a zwnj before those vowels, provided the OpenType implementation permits it (e.g. in Adobe CS or OS X). Zwnj has also proven useful before superscript vowels to force wide vowels over narrow consonants to become narrow or wide in places where rules to do this cannot be applied due to limitations of software implementations. A particular limitation is often found with zero-width-space, which is very useful in Khmer text to ensure correct line breaking, but which is not handed to the rules by the software in some cases.

**Zero-width-joiner (zwj)** can also be used before vowels to force a ligature between superscript vowels (◌ិ ◌ី ◌ឹ ◌ឺ ◌េ◌ី) and certain consonants (TUS p382). Since both zero-width-joiner and zero-width-non-joiner are mentioned in this context, it is probably thought that a font could have either ligatures or non-ligatures by default, but should be able in conjunction with zwj and zwnj to cause the respectively other behavior. The Mondulkiri fonts will also render register shifters as subscripts without regard for the preceding consonant if they are preceded by zwj and followed by sequences that would normally cause the rendering as subscript.

TUS provides for the **placement of coengs also after vowels** (TUS p375) to accommodate the 'uncommon' (TUS) practice of writing final consonants as subscript consonants and gives as examples �washing and ប៉ើៀ. This requires the correct placement of coengs and/or vowels in relationship the base-character in order for the reader to identify the function of the coeng correctly. In order to help the user to type syllable components in the correct Unicode order, the following rules should be applied to the placement of subscript consonants after vowels or symbols:

a) Coeng-Ro (◌្រ) should never be permitted after vowels or symbols because it's position in the pronunciation order cannot be established from the visual rendering.

b) No coeng should be permitted after Po + Aa (កា, with coeng e.g. កា្ក), because the visual result would be virtually indistinguishable from Nho with coeng (e.g. ញ្ញ). It is very useful to forbid coengs in this context because in pre-Unicode fonts Nho with coeng was typed as Po + Aa + coeng. This has become a very common typo in Unicode because it is still permitted in the Windows implementation.

c) Non-spacing coengs (e.g. ្ក) should only be allowed after vowels with spacing to the right, in particular the vowels Aa, Oo and Au (ា  ោ  ៅ), but possibly also after the vowels Ya and Ie (ឿ and ៀ) in order to not unnecessarily restrain the user (Most likely, however, few fonts will ever place a coeng correctly after Ya or Ie.). Non-spacing coengs also need to be allowed[9] after the vowel sign Aam (ាំ, named in TUS p378). They should not be allowed after any other vowel or symbol. The coeng's placement must make clear its position in the syllable (e.g. ខ្ញុំ and ខុំ្ញ would be two different words if the second word existed).

d) Coengs spacing to the right (e.g. ្យ) could be permitted after vowels with spacing to the right, but according to the example provided in TUS, they are also permitted after superscript vowels, but not after preposed vowels (like េ) or below-vowels (like ុ). However, the placement of the vowels must make clear the coeng's position in the syllable (e.g. ប៉ើ្យ vs ប៉ើ្យ) in relation to the vowel.

e) No coeng apart from coeng-Vo should be permitted to occur after coeng-Ro. In Khmer no coeng ever follows coeng-Ro in pronunciation order (this is most likely a phonological restraint). In the Tampuan language coeng-Vo does occur after coeng-Ro but no other coeng. Alternatively coeng following coeng-Ro could be placed in the second level below the baseline to indicate visually that it follows coeng-Ro. Unfortunately it is currently impossible for the lookups in OpenType to know in which position coeng-Ro occurred in the character string (also called 'coded characters' in TUS).

## Notes on the Unicode Standard

According to TUS "the sequence (of U+17D2 and a base-character) functions as if it had been encoded as a single character" (TUS p377). Consequently a word processor should not permit to place the cursor in between them.

The Unicode Standard says: "U+17DD khmer sign atthacan is a rarely used sign that denotes that the base consonant character keeps its inherent vowel sound" (TUS  p379). It also functions as a vowel in the Bunong language.

Contrary to the formula given under the heading 'Ordering of Syllable Components' in TUS, Robat should probably not occur after coengs (or subscripts, 'S' in TUS) and there should only ever be one per syllable. This is also how the Khmer script is implemented by Microsoft. According to TUS, Robat could also be placed after coengs and there could be two of them in a syllable.

---

[9] This is an 'uncommon' way to write final consonants. Please see TUS p375.

$Z_1$ of the syllable component formula above is described as part of the 'Ordering of Syllable Components' (p381) in TUS, but its function in this position is described under 'Consonant Shifters' (p382).

The Unicode Standard says "the consonant shifter (or 'register shifter') should always be encoded immediately following the base consonant" (TUS p382), but this is problematic, because in syllables with coengs it is sometimes clear that the consonant shifter actually changes the pronunciation of the coeng and not that of the base-character. But whether the consonant shifter is to be rendered as subscript (in the shape ្ិ) or not, depends on the character it modifies. Please see note above and discussion below under 'The rules for the dropping of register-shifters'. It seems better if the consonant-shifter is to follow the element that it is to modify or be placed after all coengs if occurring in a syllable with coengs.

The example 'ប៊ីយ៉ែរ ba + zwnj + triisap + ii + yo + ae + ro "beer"' (TUS p382) is incorrect.

Triisap should never be turned into its subscript form following Ba, the zwnj is therefore not necessary between Ba and Triisap. Please see discussion under 'Rules for the alternate rendering of register-shifters as subscripts'.

Additional note to the section 'Spacing' of the Unicode Standard (TUS p383): Though Khmer does not use whitespace between words, some languages using the Khmer script do use whitespace between words and an extra wide whitespace between clauses (e.g. Krung and Tampuan). These languages prefer a more narrow space (e.g. U+2006) between words than Khmer uses as clause separator and a wider space than the normal Khmer space (U+0020) between clauses (e.g. U+2003). Therefore fonts for Khmer script should include the full complement of whitespace characters (U+2002 - U+200A) for the user to chose from.

The shape of the ligature of Sso[10] (U+179E) and Aa (U+17B6) (ឞ and ា, ឞា) should not follow the ligature of Ba (U+1794) and Aa (U+17B6) (បា), because it cannot be confused with U+17A0 (ហ).

The shape of coeng La (U+17D2 plus U+17A1) seems unclear. TUS has it as a variant of coeng Ba (្ប), but most fonts implement it as a small version of La (្ឡ).

Tampuan uses Samyok followed by Reahmuk to mark sounds that do not exist in Khmer. In this context the Samyok should be placed right above the Reahmuk (e.g. ទ៎ះ), even if Reahmuk is preceded by Muusikatoan or spacing vowels like ា (e.g. គៀ៎ះ). In Tampuan a zero-width-non-joiner will precede Muusikatoan if it is directly followed by Samyok to prevent it from being rendered as subscript (e.g. ម៊៎ះ). There may be vowels between Muusikatoan and Samyok. In those cases Muusikatoan should never be rendered as subscript, even without the addition of zero-width-non-joiner (e.g. ប៊ី៎ះ, ប្រ៊ី៎ះ).

---

[10] For the names of Khmer letters and their Unicode 'numbers' please see the Unicode chart for Khmer at http://www.unicode.org/charts/PDF/U1780.pdf

# The rules for the alternate rendering of register-shifters as subscripts

The following rules suggested and much of the in common use, but they are not officially established in any Unicode document, some parts may be open for discussion.

The register shifters (called 'consonant shifters' in TUS) Muusikatoan (U+17C9 ់) and Triisap (U+17CA ៉) are sometimes rendered as subscript in the shape of the vowel U (ុ). The general rules for this alternate rendering are:

1. If Muusikatoan is preceded by (ង គ្ ញ ជ្ ន ដ្ ប ព្ ម ទ្ យ ៗ រ ល ឡ្ វ ឝ គ ឝ្) (note: coeng Ro is not is this list) and is followed by a superscript vowel (ិ ី ឹ ឺ េី), Samyok (់) or the vowel Aam (ាំ), then Muusikatoan should be rendered as subscript (ុ) under the base consonant or under the spacing coeng that precedes it. No and Lo do not usually carry Muusikatoan, there are only two Khmer word in the Chuon Nath dictionary where they do carry it (សិតអន្លែត as alternative spelling of សិតអន្លិត and អន្លាយ as alternative spelling of អន្លាយ). However, they regularly carry Muusikatoan in the pronunciation guides of the Choun Nath dictionary. Sha (ឝ) and its coeng carry Muusikatoan following the above rule in the Krung language.

2. If Triisap is preceded by (ផ ឡ្ ស ហ អ) (note: not all coengs of the base consonants in this list are included) and followed by a superscript vowel or Aam, then the Trisap should be rendered as subscript (ុ) under the base consonant or the spacing coeng that precedes it (e.g. ស៊ី សា៉). In the case of coeng-Sa (ស្), Triisap should only be rendered as subscript if followed by a superscript vowel, but not if followed by Aam (e.g. បន្ស៊ី or ផ្សាំ[11]). Triisap always remains in its original shape as superscript after Ba (ប), even if it is followed by a superscript vowel in order to enable the distinction between បី and ប៊ី (the latter using Muusikatoan). It also always remains in its superscript shape if it occurs in a syllable with coeng Ha (្ហ) or coeng Qa (្អ). Sso (ផ) and its coeng carry Triisap following the above rule in the Krung language.

These rules are well established as they relate to the preceding base characters, but not well known in regards to Sha and Sso.

In some cases the 'authoritative Chuon Nath dictionary' (TUS p379) explicitly permits both, the rendering as superscript or as subscript (see appendix) as a user choice and some Khmer names are spelt with Triisap as superscript even though the syllable contains a superscript vowel (e.g. រ៉ឺ). Therefore a font should permit the use of zwnj preceding the

---

[11] ផ្សាំ is the only word found in the Chuon Nath dictionary with this sequence. The following word has been seen in text, but not in any dictionary: បន្សាំ.

register shifter to achieve alternative rendering. No application should strip zwnj (nor zwj or zero-width-space (short: zws)[12]) from a text string in the copy and paste process[13].

In words containing coeng-Ro, the register shifter could be placed after the main consonant or after coeng-Ro. ស្រ្គ៉ុប ('sound produced by a heavy solid object falling') can also be spelt ស្រ៊ុប. This should be achieved by inserting a zwnj before the Triisap[14]. Because OpenType does not convey the position of the Triisap in relation to the coeng Ro to the font, both spellings should be expected in texts in other rendering systems (e.g. ស្រ៊ុ (Triisap before coeng) and ស្រ៊ុ (Triisap after coeng)). The Mondulkiri implementation will render the Triisap as superscript if Triisap follows Sa (ស) and is preceded by zwnj, or if Triisap follows the coeng-Ro.

In words with coengs with spacing to the right, the register shifter should always be placed after the coeng. The exception for this rule could be the word ប្ប៉ាណូ ('piano' found in this way in some texts), where the Muusikatoan affects the pronunciation of the base consonant as well as that of the vowel. Unlike Khmer, Tampuan has both voiced as well as voiceless 'ប' before coengs, therefore it does need to mark 'ប' before coengs with Muusikatoan as well as Triisap.

In words with non-spacing coengs the register shifter is probably also better be placed after the coeng - in contrast to the statement in TUS (p368). It should not be rendered as subscript (ុ) in conjunction with coeng Ha and coeng Qa, as well as in words where the base consonant can carry Triisap, but the coeng does not. The Mondulkiri implementation will render ហ្គ៊ើ with the Triisap as superscript if Triisap follows Ha and is preceded by zwnj, or if Triisap follows the coeng-Vo. In order to achieve text that can easily be typed and consistently be searched for and sorted, it might be best to establish a general rule to place register shifters always after coeng and leave it to the font to decide where to place it in the final image. Whether the register shifter gets rendered as subscript could then be covered by the above rules 1 and 2 plus with the exception of the rendering of words like ស្រ្គ៉ុប that in the strict sense of those rules would be rendered as ស្រ៊ុប because Triisap would not become a subscript after coeng-Ro with no means of forcing the Triisap to render as U. Theoretically, register shifters could be forced to render as subscript by insertion of zwj before the register shifter. But this is currently not provided for in TUS, though implemented in the Mondulkiri fonts.

---

[12] MS Word on Windows (Word 2010 on Windows XP) does strip zws in text strings that are copied from Word and inserted into some other applications.

[13] A text created in Windows XP run in VMware on OS X 10.6 will paste zwnj correctly into TextEdit and from Text Edit into Pages (OS X 10.6.8), but not from Windows into Pages. Nor does Pages seem to permit the insertion of zwj or zwnj from the keyboard.

[14] But it can in some fonts (including the Mondulkiri fonts) also be achieved by placing the Triisap after the coeng Ro. An Opentype font has unfortunately no means to test where the Triisap is placed in the underlying text.

# Recommendations for the implementations of Khmer Unicode

- Implementations should insert dotted circles consistently between fonts and systems to help the user to type cross-font, cross-software and cross-platform compatible text.

- Fonts should include the full complement of spacing characters.

- Fonts should not permit the 'overstriking' of characters, e.g. two bontoks or twice the same coeng in exactly the same place.

- Text-editing software and rendering systems should make zws accessible to the rules in a font in order to facilitate good kerning across syllables separated by zws. Khmer needs this in particular because some diacritics can easily overlap with components of neighboring syllables or words.

- For increasing spacing between letters, systems should insert space also between 'before'-vowels, coeng-ro  and the base character (e.g. in កើ ក្រើ ក្រ្បៀ) and also before and after other spacing coengs (e.g. ក្ប កើ្ប). Space should also be increased between spacing vowels and the base character (e.g. កា) if the ligature formation is turned of either due to typographical options or due to the insertion of zwnj.

- There is one Khmer word (and only one) with three superscript components among the main entires in the Chuon Nath dictionary that needs to be rendered correctly: អ៊ីះ! 'an exclamation'.

- Text editing software should permit the typing of zwj and zwnj and preserve them, as well as zws, faithfully.

- Text editing software should expect zero-width-space in the text and treat them like an ordinary space for the purpose of line breaking. A line of text should never be broken after the coeng-character (U+17D2) unless it is not followed by a base character (i.e. in incomplete or wrong text). Ideally it should not be possible to place the cursor between the coeng-character and the following base character, with the exception of when coeng is entered in between two already existing base characters.

- Zws, zwnj and zwj: there should never be one of these following any other or two of the same.

# List of syllables with register shifters and coeng in the same syllable

The following list was extracted from an extensive list of words found on the internet.

**RS with spacing coengs:**

ប៉្រា(in ប៉្រាណ្ឌ).ម្ប៉ា(in ចម្ប៉ា).

ស្ប៉ឺ(in ស្ប៉ឺម).

ស្ប៉ី (found in word list, but most likely misspelling of ស្ប៉ឺម in a font that renders ស〇្ប〇̃〇̃with 〇̣).

ស្ប៉ុ.រ៉ូ.ល្ប៉ូ.ម្ប៉ុ.ម្ប៉ា.ម៉្ប៉ា.ស្ប៉ា.ស្ប៉ុ.ស្ប៉ូ.ហេ្ប្រ.

ម្ប៉ា (found in word list, but considered misspelling)

ហ្ប៉ុ.ហ្ប៉ើ.ហ្ប៉ែ.ហេ្ប៉ែ.ន្ប៉ី.ស្ប៉ី.

**RS with coeng Ro:**

ម្រ៉ុ.ម្រ៉ា.ស្រ៉ា.ស្រ៉ិ.ស្រ៉ី.ស្រ៉ុ.ស្រ៉ុ.ស្រ៉ូ.ស្រ៉ែ.

**RS with non-spacing coeng:**

Muusikatoan: ង្ប៉ុ.ធ្ន៉ុ.ម្ព៉ុ.ម្ព៉ែ.ម្ព៉ូ.ម្ប៉ុ.ម្ប៉ិ.ម្ប៉ុ.ម្ប៉ូ.ម្ព៉ូ.ម្ព៉ុ.ម្ព៉ែ.ម្ព៉ូ.ម្ព៉ូ.ម្ប៉ឺ.ម្ព៉ែ.មេ្ព៉ះ.រ្ប៉ា.វ្ល៉ៃ.

ស្ព៉ូ (in ស្ព៉ប in word list, not attested in dictionary)

Ba-Muusikatoan: ប្ល៉ែះ.

Ba-Triisap: ប្ល៉ិះ.ប្ល៉ីះ.ប្ល៉ីះ.ប្ល៉ី.

Triisap: ស្ឡ៉ុ (in ស្ឡ៉កគ្រោក).ស្ល៉ូ (in ស្ល៉ត in the word list, not attested in dictionary).

ស្ល៉ុ.ស្ល៉ីះ.ស្ល៉ុ.ស្ល៉ុះ.ស្ល៉ុ.ហ្ល៉ុ.ហ្ល៉ី.ហ្ល៉ុ.ហ្ល៉ុ.

ម្ល៉ីះ.ម្ល៉ុ.អ្ន៉ីះ.អ្ន៉ី.អ្ន៉ីះ.អ្ន៉ះ.អ្ន៉ែះ.អ្ន៉ែះ.អ្ន៉ះ.អ្ន៉ះ.

**RS with coeng where the Chuon Nath dictionary permits both spellings (some examples only):**

ហ្ប៉ិះ/ហ្ប៉ិះ ហ្ប៉ុយ/ហ្ប៉ុយ ស្រ៉ុប/ស្រ៉ុប អ្ន៉ុម/អ្ន៉ុម

**Examples from the Tampuan language:**

Please note that in some examples the register shifter is placed between base character and coeng, in some it is placed after. This represents no difference in meaning, but is only due to inconsistent typing.

ហ្ប៉ា.ប្ប៉ា.ហ្ប៉ុ.ហ្ប៉ុ.ហ្ប៉ុ.ប្ប៉ុ.ហ្ប៉ៀះ.ហ្ប៉ៀះ.ហ្ប៉ា.ប្ប៉ា.ប្រ៉ុ.ប្រ៉ា.ប្រ៉ុ.ប្រ៉ី.ប្រ៉ឺ.ប្រ៉ុ.ប្រ៉ែ.ប្រ៉ូ.ប្រ៉ឺ.ប្រ៉ើះ.ប្រ៉ុ.ប្រ៉ុះ.

ប្រ៉ុះ.ប្រ៉ុះ.ប្រ៉ែ.ប្រ៉ៅ.ប្រ៉ៅះ.ប្រ៉ៅះ.ប្រ៉ៅះ.ប្រ៉ៅះ.ប្រ៉ុំ.ប្រ៉ុះ.ប្រ៉ុះ.ប្រ៉ុះ.

ប្រ៉ុះ.ប្ល៉ុ.ប្ល៉ុ.ប្ល៉ិ.ប្ល៉ី.ប្ល៉ុ.ប្ល៉ុ.ប្ល៉ៀះ.ប្ល៉ែ.ប្ល៉ែ.ប្ល៉ែ.ប្ល៉ៃ.ប្ល៉ៅ.ប្ល៉ៀះ.ប្ល៉ះ.ប្ល៉ះ.

ប្ល៉ុះ.ប្ល៉ុ.ប្ល៉ុ.ប្ល៉ុ.ប្ល៉ុះ.ប្ល៉ុ.ស្រ៉ុ.ស្រ៉ែ.ស្ល៉ែ.ស្ល៉ុ.ហ្ល៉ុ.ហ្ល៉ុ.ហ្រ៉ា.ហ្រ៉ា.ហ្រ៉ុ.អ្រ៉ុ.អ្រ៉ុ.អ្រ៉ើ.អ្រ៉ៀ. គ្រ៉ៀ.គ្រ៉ៀ.អ្រ៉ៀ.គ្រ៉.អ្ន.ម្ន៉ែ.