Structuring Tai Tham Unicode

M. Hosken, SIL International, Payap University Revision 7

Table of Contents

1 Introduction2	
2 Syllable Structure	ŀ
2.1 Consonants4	ŀ
2.2 Medials	ŀ
2.3 Vowels5	j
2.4 Tones	,
2.5 Finals5	j
2.6 Spacing Vowels	j
2.7 Syllable Chaining6	;
2.8 Mai Kang Lai7	1
2.9 Mai Sam7	,
2.10 Mai Kang	}
2.11 Ra Haam	3
2.12 Cryptogrammic Dot	}
2.13 Digits)
2.14 Extra Consonants9)
2.15 Summary9)
3 Visually Motivated Ordering11	_
3.1 Encoding Order	_
3.1.1 Examples	;
3.2 Under Differentiation14	ŀ
4 Refined Visual Encoding16	j
4.1 Issues	5
5 Implementation)
5.1 Shaping)
5.1.1 USE Structure)
5.1.2 OpenType Implementation19)
5.2 Normalization20)
5.2.1 Canonical Ordering20)
5.2.2 Sorted Ordering20)
6 Acknowledgements23	5
7 Bibliography24	ŀ

1 Introduction

This paper is a draft and is liable to change significantly. As such it also falls fowl of: "Why did you not talk about *x*?" Probably because I didn't think of it; or "Why did you not describe *y* more clearly?" Probably because my brain isn't the same as yours. In either case I am happy to try to improve things to make things clearer and more accurate. But my primary concern is that we have a good encoding and not necessarily a perfect paper describing it.

The Tai Tham script is one of the more complex scripts in Unicode. In many respects it behaves like the Myanmar script, but with added complexity. It uses subjoined consonants not just for syllable chaining but for initial clusters, ala Devanagari, for final consonants ala Old Khmer, and even for vowels. In addition, Tai Tham based orthographies typically have no word spaces¹ and, missing any phonetic syllable separators, therefore, has the added ambiguity of whether a base consonant is a final or initial. As a result, Tai Tham has some inherent visual ambiguities. Having said that, people are still able to read, and there are a number of reading principles that, certainly for modern spelling conventions, allow readers to minimise the ambiguity. The problem is that these principles are very hard to codify and are too complex for the generic computing models into which the script may be desired to fit.

The upshot of all this is that when the script was added to the standard, a relatively laissez faire attitude was taken to encoding order. The expectation was that users would type and store data in reading order. In effect, the reading ambiguities would be removed in stored text. The reason for this is that there has been considerable variation in spelling and types of syllable structure in how people have written using the Tai Tham script, and if the addition of Tai Tham to Unicode had had to wait until a full analysis had been done, Tai Tham would still not be in Unicode. But the resulting encoding order approach is problematic since it results in over differentiation. Two different, but correctly encoded, strings may render identically and there is no way to distinguish them. Tai Tham is unique in Unicode in the amount of over differentiation that it supports.

If over differentiation is bad for encodings, under differentiation is even worse. Under differentiation is where it is not possible to encode a necessary visual difference and is considered a technical deficiency of an encoding that is sufficient justification to change an encoding in Unicode. It is where two different visual representations require the same underlying string (and therefore one of the representations is lost, or out of band information is necessary to resolve the rendering ambiguity). A good encoding design, therefore, requires that it not contain over differentiation or under differentiation. For Tai Tham, due to its visual complexity, this is very hard.

Another aspect to encoding is whether an encoding is primarily visual or logical. Does it store the surface rendering form or the underlying reading form. The current Unicode standard of Tai Tham is entirely a logical encoding. But this leads to huge over differentiation. The problem of going to an entirely visual encoding is that it makes analysis really hard. For example, culturally appropriate collation becomes nearly impossible. Tai Tham collation is particularly problematic since most sort orders require the user to be able to read the word and to use the phonemic representation for collation. This goes against all the principles of computer collation. The complexity of converting from a visual form to a logical one is huge and beyond what most sorting systems can handle. The question this paper examines is whether there is some compromise position for an encoding that is sufficiently logical to allow analysis while also straightforward to render.

A further dimension to the question comes from input. The current encoding of Tai Tham requires that a typist be able to read what they type correctly. They cannot merely look at a text and type it. They have to, in effect, read it and then type the logical representation of that reading. Along with this, they have to get all the diacritics in the right relative order. If they get it wrong, at best they may see what looks like a perfectly good word being marked as a spelling error. But more often, they just end up with inconsistently encoded text which makes searching and analysis nearly impossible.

¹ Modern Pali being an exception.

The purpose of this paper is to propose a new encoding order for Tai Tham that uses all the same codepoints in Unicode but that reduces over differentiation, preferably removing it altogether. In addition, it aims to be one that lends itself to analysis by, for example keeping vowel sequences together as much as possible and storing tone after vowels, etc. Obviously, the encoding aims to not be under differentiated with respect to the needs of its user community. Tai Tham script usage is poised for much needed deployment and wide usage. Currently there is relatively little data stored and so there is a window of opportunity to make the changes needed to provide a good basis for future computing in the script.

The paper starts with a description of the current encoding and then moves into proposing what changes can be made to reduce the structural complexity of the encoding such that it can be rendered and keyed. From here we examine implementation details, particularly with regard to the Universal Shaping Engine for OpenType support and how to normalize text into a canonical form of particular relevance to keyboard entry.

A number of examples are given in tables throughout the paper. Each example shows a Lanna style rendering, a Khuen style rendering and then the underlying stored codes. Two renderings are shown since they can be quite different.

2 Syllable Structure

In this section we look at the structure of an orthographic syllable. This may well differ from a linguistic syllable in that it is primarily a graphical construct, that while motivated by the phonemics it conveys, is not governed by them. Tai Tham loosely follows the abugida structure of syllables, but stretches it beyond what many may consider the bounds of good order! The description here is introductory and more thorough treatments may be found in the literature. In the various subsections of this section we build up a regular expression description of the orthographic syllable².

The basic structure we examine is:

 $\mathsf{S}=\mathsf{C}\;\mathsf{M}\;\mathsf{V}\;\mathsf{T}\;\mathsf{F}$

C – Consonant and required; M – Medial; V – Vowel; T – Tone; F – Final. Only C is required and any of the components may consist of a sequence of Unicode characters. This is not the only structure as we will see.

2.1 Consonants

A consonant consists of any character in the range U+1A20 – U+1A4C. There are also independent vowels that may take this initial position in the range U+1A4D – U+1A52. There are two ligatures that may also start an orthographic syllable: U+1A53 (\Im LAE), U+1A54 (\Im GREAT SA).

One situation in which the initial consonant is a sequence is when considering the tone class prefix h (U+1A49) as in this table which shows a Lanna style rendering in the first column, a Khuen style rendering in the second and the sequence itself in the third.

හිුුු	ဟိုရ	1A49 (లు на) 1A60 sakot 1A36 (్రై NA) 1A62 (్ MAI SAT) 1A26 (్రై NGA)
-------	------	---

Thus in regular expression terms we can say:

C = [\u1A20-\u1A54] | \u1A49 \u1A60 [\u1A20-\u1A4A]

2.2 Medials

The most common expression of a medial is the use one of the medial characters: U+1A55 (medial RA), U+1A56 (medial LA), medial ya encoded as U+1A60 SAKOT U+1A3F (Medial Wa encoded as U+1A60 SAKOT U+1A45 (Medial Wa). Notice the number of places that this sequence crops up as a valid sequence.

There is some ambiguity over the use of medial la, U+1A56. Some fonts do not contrast the visual presentation of the medial la over subjoined la (U+1A60 U+1A43). This is unfortunate since it is important that the distinction is made. If there is no contrast in a font, then someone may not type the contrast and then if the text is rendered in a font with the contrast, the resulting text may be wrong. In addition, it means the typist may be introducing spelling errors that they cannot see in the font with no contrast. Notice that in the h prefix context, the medial la is used:

బ్	ഗു	1А49 (లు на) 1А60 sakot 1А43 (్లై LA)	
ಬ್ಬ	Ś	A49 (లా на) 1A56 (్ట్ medial La)	
ల్ట్ ఉ	တ္လ်ရ	.A49 (లా на) 1A56 (్ట్ల medial la) 1A62 (్ mai sat) 1A26 (ధ nga)	
ဂ်ု	્રા	1А23 (стка) 1А60 SAKOT 1АЗҒ (су YA) 1А69 (су U) 1А76 (сб толе-2)	
હેગ્ર	ଚୁକ୍ର	A23 (೧ ка) 1A60 SAKOT 1A45 (ႏ wa) 1A63 (୨ АА) 1A60 SAKOT 1A3F (ເງ YA)	

² For those turning this into a full regular expression, there is an implied parenthesis around each defined subexpression and spaces may be removed.

Medial ya and medial wa may also act as vowels. There is a visual ambiguity here. Rather than carry that ambiguity into the regular expression, we allow for a miscategorisation of a vowel as a medial, and then allow multiple medials. This is not a problem since the primary concern with such string processing is ensuring a consistent order rather than precise categorisation of components. There is only visual confusion over the relative order of U+1A55 (\therefore MEDIAL RA). The vertical stacking of medials/vowels is achieved given the order listed. Removing the medials from what may occur subjoined after an h, allows us to simplify and reduce the list.

 $C = [\u1A20-\u1A54] | \u1A49 \u1A60 [\u1A20-\u1A3E]$ $M = \u1A55? \u1A56? (\u1A60\u1A45)? (\u1A60\u1A3F)?$

2.3 Vowels

Unlike most Indic scripts, Tai Tham is used for languages that use much richer sequences of vowel characters. For consistent ordering, this means agreeing how compound vowel sequences are to be stored. The general principle is that vowels are stored as:

 $V = V_p? V_b? V_a?$

While it is clear that V_p must come first, the other two are somewhat arbitrarily ordered. We follow the order in the original proposal, here.

We will discuss spacing vowels later. All vowel classes are optional.

 $\label{eq:V_p} V_p = [\ulA6E-\ulA72] \\ V_b = [\ulA69-\ulA6A\ulA6C] \\ V_a = [\ulA65-\ulA68\ulA6B\ulA73] [\ulA62\ulA74]? \\ \label{eq:V_p}$

Allowing U+1A62 (Si MAI SAT) in conjunction with other upper vowels is for the case where it acts as a final -k. U+1A74 (Si MAI KANG) has all kinds of functions and may combine with any vowel, but occurs before a tone mark in this context. Some special case sequences need consideration, in particular in relation to the medials that can also work as vowels:

ිට ා -	ູ່ ເວັ	1A26 () NGA) 1A60 SAKOT 1A45 (WA) 1A6B () 0) 1A75 (TONE-1)
ń	က်ွံ့	1A20 (ເກ KA) 1A60 SAKOT 1A45 (ູ່ WA) 1A76 (໌ TONE-2) 1A60 SAKOT 1A3F (ູ່ງ YA)

The tone marks help in identifying whether something is a medial, vowel or final, but even with these, the order is not necessarily cut and dried. Is the U+1A60 sakot U+1A45 (\Im WA) here acting as a medial or a vowel? To a large extent, it is not important since we are only interested in order and not categorisation.

2.4 Tones

Tone marks are stored after the vowel sequence. There can only ever be one if present. The only complexity with tones lies in their interaction with spacing vowels, to which we will return.

 $T = [\u1A75 - \u1A79]$

2.5 Finals

There are two kinds of finals: spacing and diacritical. Spacing finals are consonants that may in turn be modified through syllable chaining. Diacritical consonants are primarily subjoined consonants or a U+1A58 (: MAI KANG LAI).

 $F_s = [\u1A4B\u1A54] | \u1A60? [\u1A20-\u1A4A] \\ F = F_s? \u1A58?$

The restricted list of consonants is used because certain consonants are never used for finals. For example, letter A U+1A4B (e_A) is used as the spacing component in a compound vowel and not as a final consonant. This list could be further constrained, but for little value.

2.6 Spacing Vowels

In addition to sequences of non spacing vowels, there may are also be spacing vowels. These can occur as part of a vowel sequence, but the overall sequence my be interrupted by a tone mark. Thus an alternative syllable pattern is:

 $C M V T V_s F$

U+1A63 (9 AA) is a common spacing vowel. U+1A64 (1 TALL AA) is a variant form of the same vowel for use where there may be confusion over whether U+1A63 (9 A) is a separate vowel or part of the preceding consonant, visually, to create another consonant and so cause reading confusion. U+1A61 (1 A) is a short -a that can be used to shorten a vowel while also acting as a final consonant. U+1A4B (1 A) can combine in a vowel sequence to be part of a complex vowel. Thus we get complex vowels like:

င့်ကို့အု	လို့သး	1А20 (ເດ ка) 1А6Е (ເວີ E) 1А6С (ເວັ OA) 1А65 (ເວົ I) 1А75 (ເວົ້ TONE-1) 1А4В (ເຊ A) 1А61 (ເ A)
ငရာ၇ႜ	ကော	1А20 (の ка) 1А6Е (の:: E) 1А63 (つ АА) 1А61 (ล А)

 $V_s = [\u1A63\u1A64] \u1A74? | [\u1A4B\u1A63]? \u1A61$

The U+1A4B (a_A) and U+1A61 (a_A), in this context, are never followed by a final. U+1A4B (a_A) may occur in a syllable chaining context (see below). One refactoring is to integrate U+1A61 (a_A) into F_s and simplify the regular expression. This then allows for finals following a spacing vowel. For example:

භුදි

Sara am is encoded U+1A63 (\circ A) U+1A74 (\therefore MAI KANG) and takes no following subscript consonants.

As a result, the regular expression is restructured to accentuate the cluster structure of a syllable:

$$\begin{split} S &= CC \; (V_s \mid V_f)? \\ CC &= C \; M? \; V? \; T? \; F? \\ C &= [\ulA20-\ulA54] [\ulA5A\ulA5B]? | \\ \ulA49 \; \ulA60 \; [\ulA20-\ulA3E] \\ V_s &= AV \; F? \\ AV &= [\ulA63\ulA64] \; \ulA74? \\ V_f &= \ulA4B? \; \ulA61 \end{split}$$

Here we are saying a syllable consists of a consonant cluster (CC) followed optionally by either a spacing vowel (that can take a final) or a vowel that is itself a final. The CC is a single cluster separated from any following clusters, except for the case of sara am where the sequence U+1A63/4 (? AA) U+1A74 (: MAI KANG) has the U+1A74 (: MAI KANG) rendered as part of the preceding CC cluster, in Lanna. There are other visual ligation and rendering that may also occur, for example na (: U+1A36 (: NA) U+1A63 (? AA).

ໍ້ຄາ	б э	1А23 (С КА) 1А76 (С толе-2) 1А63 (Э АА) 1А74 (С маі Калд)
------	------------	---

2.7 Syllable Chaining

In addition to having subjoined finals, Tai Tham allows syllable chaining as found in the Myanmar script, along with others. Here a spacing final is followed by a subjoined consonant that is the start of the next syllable in the word. The result is some sequences that would otherwise be rarely seen.

ငက္တ	တ္တေ	1A20 (ຄາ ка) 1A6E (ຄະອິ E) 1A60 SAKOT 1A32 (နဲ့ та) 1A69 (ອີບ)
్ ట్ర ీ బ్రీ	_{င်} မျင်တဲ့	1A36 (్లై NA) 1A6C (్లై OA) 1A3E (ఆ MA) 1A60 SAKOT 1A3E (్లై MA) 1A6A (్లై UU) 1A36 (్లై NA) 1A65 (ే I) 1A3F (లు YA) 1A74 (ం) MAI KANG) 1A60 SAKOT 1A32 (్లై TA)

The way to include syllable chaining is through recursion in the regular expression. As a result, though, the expression is no longer truly regular. This is addressed in a later section.

 $F = F_s? [\u1A58]? (\u1A60 S)?$

Including syllable chaining in the definition of a syllable can result in what are linguistically polysyllabic syllables. For example, വല്ലാമും This is undesirable and a more refined definition of an orthographic syllable would address this issue. One solution is to break the chain by treating a spacing final as the start of a new syllable. This is not linguistically correct, but it does reduce the complexity of a sequence.

 $F_s = [\u1A4B\u1A54] | \u1A60 [\u1A20-\u1A4A]$

At this point our orthographic cluster starts to look a lot like and extended grapheme cluster.

2.8 Mai Kang Lai

U+1A58 (: MAI KANG LAI) represents a syllable final -ng. It may be placed in different places according to font preference. It may occur over the top right of the cluster it follows (as in the Khuen example), or it may occur between the cluster and the next cluster or it may occur over the start of the next cluster, like a kinzi in Burmese. This can even extend to being over a prevowel, as shown here from

a Northern Thai Dictionary (Udom 2004). But this could well be a misrender.

In Khuen the convention is to render the U+1A58 (SMAI KANG LAI) over the right hand side of the base. Which is somewhat in the space between the clusters. (Tengok).

The key thing here is that the position of U+1A58 (: MAI KANG LAI) is font specific.

2.9 Mai Sam

U+1A7B (SMAI SAM) has a number of functions. One of these is to indicate that a following subjoined character is the start of a new syllable. In the original Unicode proposal, the mai sam is stored immediately after the new syllable initial. For example:

သ္လ်ာ သို႔ 1A48 (သ SA) 1A60 SAKOT 1A43 (္လ္လ္လ္လ္လု LA) 1A7B (်က္ကန	I SAM) 1A6A (∷ UU)
---	----------------------------

One way of adding this to our regular expression is through recursion as part of the initial consonant sequence:

 $C = [\u1A20-\u1A54] \u1A7B? | \u1A49 \u1A60 [\u1A20-\u1A3E]$

It is noticeable that different fonts position U+1A7B (: MAI SAM) in different positions depending on the presence of another upper diacritic. For example, it can occur before the diacritic or even at the end of the syllable: స్పేద vs స్పేద in two editions of the same dictionary.

Another use of mai sam is as a reduplication mark. For this, the character is stored at the end of the syllable.

The final use of mai sam is in the place of a double acting consonant, as such it is a kind of final. Both of these cases are handled by:

 $F = F_s? [\u1A58\u1A7B]? (\u1A60 S)?$

The net effect of both of these positions is that U+1A7B (SMAI SAM) immediately follows the consonant.

2.10 Mai Kang

U+1A74 MAI KANG has a number of functions, reflecting the use of U+0E4D (் NIKAHIT) in Thai script.

- 2. In combination with U+1A6C (ු oa) as an open final to the vowel as in ද් U+1A21 (ස кна) 1A6C (ු oa) 1A74 (් маг кало) U+1A75 (් толе-1).
- 3. In Khuen it can act as a superjoined wa, in place of the medial: $\frac{1}{3}$ U+1A21 (\approx KHA) U+1A74 (\approx MAI KANG) U+1A60 SAKOT U+1A3F (\approx YA). We encode it as a vowel in this context.

Since, in Lanna, MAI KANG in 'am' is rendered over the consonant, MAI KANG is not used over a consonant before a U+1A63 AA vowel. In some older spellings, this does occur.

2.11 Ra Haam

U+1A7A (3 RA HAAM) acts as a silencer and is placed over anything that is to be silenced, including misspellings. In effect it can be used to mark a final consonant by silencing the inherent vowel on the consonant. It is stored after everything else but before a U+1A7B (3 MAI SAM):

 $SM = [\u1A7A\u1A7C]? \u1A7B?$

To support this, we add U+1A7B (3 MAI SAM) to the syllable modifiers. Notice that U+1A7C KHUEN-LUE KARAN has the same function in Khuen.

There is some confusion over why there are separate codes for the same function in Lanna versus Khuen. Lanna uses U+1A58 MAI KANG LAI and U+1A7A RA HAAM, while Kuen uses U+1A58 MAI KANG LAI and U+1A7C KHUEN-LU KARAN. The reason that a common code couldn't be used and styled appropriately, like so many other characters in Tai Tham is that they style of U+1A7C KHUEN-LUE KARAN is very similar to the style of U+1A58 MAI KANG LAI in Lanna. Thus if the wrong style were used, there would be confusion over which character was being referred to. To allow for more unified fonts, it was decided to keep U+1A7C KUEN-LUE KARAN only for use in Kheun and U+1A7A RA HAAM only for use in Lanna. Kheun does make use of U+1A7A RA HAAM for an archaic final ra. U+1A58 MAI KANG LAI is styled differently for Lanna versus Kheun.

2.12 Cryptogrammic Dot

The U+1A7F (CRYPTOGRAMMIC DOT) is used to modify consonants in a standard cryptographic pattern where the number of dots says how many characters to advance the base character through the sequence to find the actual letter this character represents. There are often multiple dots, and they can occur with other diacritics:



This calls for the cryptogrammic dot to be a consonant modifier.

```
CC = C CM* M? V? T? F?
CM = \u1A7F
```

2.13 Digits

Digits can act as consonants and are used in contractions or in cryptogrammic text:



As such it makes sense to add both sets of digits to the consonants class:

 $C = [\u1A20-\u1A54\u1A80-\u1A89\u1A90\u1A99] [\u1A5A\u1A5B]? | \u1A49\u1A60 [\u1A20-\u1A3E]] \\$

2.14 Extra Consonants

There are various extra consonants in the range U+1A57, U+1A59 – U+1A5E that need to be accounted for. In addition U+1A6D still requires a home.

The characters U+1A57 (Ising La Tang Lai) (strictly a ligature of subjoined la and subjoined nga), U+1A5C (Ising MA), U+1A5D (Ising MA) and U+1A5E (Ising SIGN SA) are merely alternative forms of subjoined characters to be used when they contrast with the normal subjoined forms, and therefore can be used wherever such subjoined consonants go. Although, in these cases they do not occur following a U+1A49 (ID HA), so are limited to finals. U+1A59 (ISING FINAL NGA) is also a final. U+1A6D (ISIGN OY) is a combination of vowel and final so we add it as a final.

 $F_s = [\u1A4B\u1A54\u1A6D] | \u1A60 [\u1A20\u1A4A]$

 $F = F_{s}? [\u1A57-\u1A59\u1A5C-\u1A5E\u1A74\u1A7B]? (\u1A60 S)?$

U+1A5B (:: HIGH RATHA OR LOW PA) is used in conjunction with another consonant even when the consonant is subjoined. Thus we can add it to the consonant specification. Likewise, U+1A5A (:: SIGN LOW PA) only occurs in one word in Tai Lue. And while it is stored after a base consonant, it is sounded before it. Therefore we add it to the initial.

 $\label{eq:C} C = ([\u1A20-\u1A54\u1A80-\u1A89\u1A90-\u1A99] [\u1A5A\u1A7B]? | \u1A49 \u1A60 [\u1A20-\u1A3E]) \u1A5B?$

2.15 Summary

In summary our syllable description is:

```
\begin{split} S &= CC (V_{s} \mid V_{f}) SM? \\ CC &= C CM* M? V? T? F? \\ C &= ([\u1A20-\u1A54\u1A80-\u1A89\u1A90-\u1A99] [\u1A5A\u1A7B]? | \u1A49 \u1A60 [\u1A20-\u1A3E]) \u1A5B? \\ u1A5B? \\ CM &= (u1A7F \\ M &= (u1A55? (u1A56? ((u1A60)(u1A45)? ((u1A60)(u1A3F)? \\ V &= V_{p}? V_{b}? V_{a}? \\ V_{p} &= [(u1A6E-(u1A72] \\ V_{b} &= [(u1A6E-(u1A72] \\ V_{b} &= [(u1A65-(u1A68)(u1A6B)(u1A73)] [(u1A62)(u1A74]? \\ T &= [(u1A75-(u1A79] \\ \end{split}
```

```
F_{s} = [\u1A4B\u1A54\u1A6D] | \u1A60 [\u1A20-\u1A4A]
```

```
V_s = AV F_s?
```

 $V_f = \ula4B? \ula61 \\ F = F_s? [\ula57-\ula59\ula5C-\ula5E\ula74\ula7B]? (\ula60 S)? \\ SM = [\ula7A\ula7C]? \ula7F?$

Notice that we have not tried to ensure this regular expression works with normalized text since this is purely for descriptive purposes.

3 Visually Motivated Ordering

The problem with the above scheme, even fully filled out to describe every situation, is that it requires someone to be able to read a word before they can type it. The prevalence of mistyping is huge and the mistypings are usually not visible to the typist and so cannot be caught, except by spell checking. And then all the poor user sees a red squiggly line or whatever indicates a spelling error, while the word looks perfectly OK to them, and is visually correct. Consider this little gem:

ທີ່ພຸດາຄ	လိုဖက႔	1А49 (ഗ на) 1А60 SAKOT 1А36 (ු NA) 1А65 (ි I) 1АЗЕ (ය MA) 1А20 (ග КА) 1А41 (න RA)
		1А49 (ഗ на) 1А65 (ි I) 1А60 SAKOT 1А36 (ු NA) 1АЗЕ (ප ма) 1А20 (ന ка) 1А41 (ာ RA)

With | separating the linguistic syllables. Both readings are correct. Also, by analagous environments:

Id	Lanna	Khuen	Thai	Encoding	Encoding	Thai	Khuen	Lanna
1.	(S2	ى گە	หนี	1A49 (ເວ на) 1A60 sakot 1A36 (ເຼັ NA) 1A66 (ເະີ້ II)	1A49 (లు на) 1A65 (్రీ I) 1A60 ѕакот 1A36 (్రై NA)	หิน	S	(S2
2.	ငက္မ	တွေ	เกว	1A20 (ග KA) 1A6E (E) 1A60 SAKOT 1A45 (ු WA)	1А20 (Ср. КА) 1А60 SAKOT 1А45 (S WA) 1А6Е (СС Е) 1А41 (Б RA)	เกวน	လွှေ <i>။</i>	ငက္ပစ

But it must be born in mind that the text itself is ambiguous. By asking a user to type the correct reading of the text, we are asking them to enter more information than is in the text itself. Tai Tham is not unique in this respect, but it is recognised as a problem in all areas where it occurs. Why should we expect a typist to disambiguate the ambiguous spelling? Can we find an encoding scheme that truly represents the text as it is presented rather than trying to encode some underlying interpretation of the text? This is a key principle:

Unicode encodes text and not underlying morphophonemics

This does not preclude encoding text in such a way as to aid its processing, but where there is ambiguity in the text, we should expect to see ambiguity in the encoding. Can we therefore come up with an encoding scheme that is "mostly right" and acknowledge that there will be situations where the encoding falls on the wrong side of an ambiguity?

One option is to standardise on some of the ambiguities regarding subjoined consonants as to whether they are syllable chaining, medials or finals. If we can get the order right and there is no impact on rendering, then we would gain a consistent storage order which would make implementing keyboards and many other things much easier / doable. On the other hand we would probably mess up collation. But since all the collation schemes presented for Tai Tham based languages involve being able to read the word in order to collate it, this is no bad thing to consider collation from scratch! But it does need to be born in mind that a switch to a visual encoding will make algorithmic analysis much harder for Tai Tham script. Getting from a visual encoding to a linguistically accurate one is not easy and will probably require dictionary approaches as well as algorithmic.

3.1 Encoding Order

So far, our analysis has been in terms of linguistic syllables. But a visually based interpretation works in terms of grapheme clusters. A grapheme cluster consists of a base character surrounded by non-spacing diacritics. This is a more general definition than is strictly defined for Tai Tham in UAX 29. In Tai Tham, we stretch even this a little in many ways, but the basic principle applies. A Tai Tham linguistic syllable may consist of multiple grapheme clusters and a grapheme cluster may represent more than one linguistic syllable.

The easiest encoding order to work with is where any diacritic in a grapheme cluster has a fixed position in that order. We will do our best to achieve that, and treat as special cases those situations where Tai Tham cannot fit that model. There are three grapheme cluster types in Tai Tham:

- 1. A consonant based cluster, based around a consonant and with diacritics attaching to that base or each other.
- 2. A spacing vowel cluster with a base of U+1A63 (2 AA) or U+1A64 (2 TALL AA) with diacritics attaching to it or each other.
- 3. A simple spacing vowel U+1A61 (a A), which takes no diacritics.

Given that the encoding requirement is that two different sequences will look different, we are not concerned with ensuring that crazy sequences, such as ensuring that illegal diacritics added to U+1A61 A, are marked as erroneous. It is then possible to merge all these cluster types into a single cluster. This resolves many of the questions of what can and cannot follow a spacing vowel.

The precise order that diacritics are stored following a base can be arbitrary so long as everyone agrees on that order and that shapers and fonts know how to get from that order to the intended visual rendering. Thus, for example, subjoined consonants may be stored early in the sequence or late. We choose early because that aligns itself more easily with other scripts and the needs of the USE (Universal Shaping Engine). Linguistically, this is not ideal, but a visual encoding is not intended to be linguistically accurate.

Based on these principles we can create an initial regular expression:

- 1 S = C CM? (SUB CM?)* M? V_p ? V_a * (V_b SUB*)? T* F?
- 2 C = [\u1A20-\u1A54\u1A61\u1A63\u1A64\u1A80-\u1A89\u1A90-\u1A99]
- 3 CM = \u1A5B? \u1A7F*
- 4 SUB = \u1A60 [\u1A20-\u1A4A\u1A53\u1A54] | [\u1A56\u1A57\u1A5C-\u1A5E]
- 5 M = [\u1A55\u1A5A]
- $V_{p} = [\u1A6E-\u1A72]$
- 7 V_a = [\u1A62\u1A65-\u1A68\u1A6B\u1A73\u1A74]
- 8 $V_b = [\ulabel{eq:Vb} = [\ulabel{eq:Vb} V_b = [\ulabel{eq:Vb}$
- 9 $T = [\u1A58\u1A59\u1A75\u1A79\u1A7A\u1A7C]$
- 10 F = \u1A7B

This definition takes some explaining.

Line 1 gives the primary cluster structure. V_b and V_a have been reordered to conform to USE order. Although in general usage, this is not preferred. In the USE multiple vowel diacritics are allowed and therefore we simplify the definition of V_a in line 7 and simply allow multiple. Since below vowels and subjoined consonants can interact visually, particularly in Khuen, it is necessary to allow subjoined consonants after below vowels.

Line 3 lists the consonant modifiers. U+1A5A (SPA) occurs only in one rare word and occurs only once. Multiple occurrences are expected to stack vertically and so be evident. U+1A5B (SPATHA OR PA) also only occurs once. U+1A7F (SPATHA OR PA) also only well occur multiple times under a character and stacks downward.

Line 4 is the list of subjoined consonants which may interact visually, vertically (or in the case of Khuen and some Lanna, horizontally). Whether the interaction is visible is based on font preference. But since some fonts to require appropriate order for visual interaction, the order needs to be the order they would be if they did interact visually.

Line 5 contains the sole medial. The reason that U+1A55 (\therefore MEDIAL RA) is a medial and not a subjoined character is that a cluster can have only one.

Line 7. U+1A74 (\therefore MAI KANG) interacts with other above vowels. U+1A73 (\therefore OA ABOVE) is a vowel in contrast to U+1A77 (\therefore TONE-3) which is a tone mark which only occurs in Khuen, and may occur after an above vowel.

Line 8 describes below vowels. Due to a contrast between words like හු and හා it is necessary to allow a SUB after a lower vowel as well as before it. But this contrast is only needed because of the visual interaction between lower vowels and subjoined consonants. In addition, only true subjoined characters may be used there and not medials.

Line 9 shows that tones and other finals do interact. For example, compare ஜி U+1A49 (ம на) U+1A6C (இ оа) U+1A59 (இ FINAL NGA) U+1A76 (இ TONE-2) compared to ஜி U+1A49 (ம на) U+1A76 (இ TONE-2) U+1A60 SAKOT U+1A38 (இ РА) U+1A59 (இ FINAL NGA).

Line 10 shows the final modifier.

The recursion is removed. Instead we change the syllable structure to reflect a grapheme cluster. We do this by treating a spacing final base as the start of a new cluster and then explicitly describe the subjoined finals via SUB.

3.1.1 Examples

Based on this structure, we give some example words along with their current and proposed spellings. Renderings are given in Lanna and Khuen style.

Id	Lanna	Khun	Original encoding	Proposed encoding
1.	රි)	ဇင်	1A20 (ເລ KA) 1A65 (ີ່ I) 1A60 SAKOT 1A36 (ເຼັ NA)	1А20 (ရာ ка) 1А60 sakot 1А36 (ູ່ Na) 1А65 (ີ່ I)
2.	ວ່າ	อ่ว	1A27 (☉ CA) 1A75 (᠅ TONE-1) 1A63 (੭ AA) 1A74 (᠅ MAI KANG)	1A27 (ᢒ CA) 1A75 (᠅ TONE-1) 1A63 (୨ AA) 1A74 (᠅ MAI KANG)
3.	C.	ေက်ီး	1A41 (⊅ RA) 1A6E (✑ E) 1A6C (ஜ OA) 1A65 (☺ I) 1A75 (᠅ TONE-1) 1A60 SAKOT 1A3F (♡ YA)	1A41 (♂ RA) 1A6E (♂᠅ E) 1A65 (᠅ I) 1A6C (᠅ OA) 1A60 SAKOT 1A3F (᠅ YA) 1A75 (᠅ TONE-1)
4.	(Ĵ)	Ś	1A23 (ఐ ка) 1A55 (్లో medial ra) 1A60 sakot 1A45 (్లో wa) 1A6B (ీ o)	1A23 (ග ка) 1A60 SAKOT 1A45 (ු wa) 1A55 (ြ MEDIAL RA) 1A6B (૽ o)
5.	ଶୁଁଥ	ຢູ່ງ	1АЗ4 (@ ТА) 1А58 (ё́маі калд Lai) 1А60 sakot 1А43 (ё́LA) 1А63 (э аа) 1А60 sakot 1АЗҒ (ё́J YA)	1АЗ4 (@ та) 1А60 sakot 1А43 (ूj la) 1А58 (маі калд lai) 1А63 (э аа) 1А60 sakot 1АЗҒ (ु ya)
6.	୍କଳୁ	က္ကေ	1А20 (ఐ ка) 1А6Е (ంప е) 1А60 SAKOT 1А32 (్హా та) 1А69 (్రై U)	1A20 (ఐ ка) 1A60 ѕакот 1A32 (్రై та) 1A6E (ంస E) 1A69 (్రై Ս)
7.	ေရာ၇ႏ	ကော	1A20 (೧೧ KA) 1A6E (೧೮೫ E) 1a63 (೧ AA) 1a61 (೧ A)	1А20 (ເກ ка) 1А6Е (ເວີ E) 1а63 (ว аа) 1а61 (ສ а)
8.	ငင်မွှဲာ	ငွမ့်ခခ	1A3E (డి MA) 1A6F (దరాపి EE) 1A76 (ే TONE-2) 1A60 SAKOT 1A36 (్రి NA) 1A60 SAKOT 1A45 (్ర WA) 1A75 (ీ TONE-1) 1A63 (ం AA)	1A3E (എ MA) 1A60 SAKOT 1A36 (ⓒ NA) 1A60 SAKOT 1A45 (ⓒ WA) 1A6F (♈்⇔ EE) 1A76 (᠅ TONE-2) 1A75 (᠅ TONE-1) 1A63 (♡ AA)
9.	ဂ်ျွာ	မ်းဂ	1АЗВ (こ РА) 1А6D (こ оч) 1А60 SAKOT 1А45 (WA) 1А75 (ご толе-1) 1А63 (ว АА)	1A3B (ූ PA) 1A6D (්) OY) 1A60 SAKOT 1A45 (ු WA) 1A75 (ਂ TONE-1) 1A63 (ා AA)
10.	તુર્ણ	<u>an</u>	1А35 (Ф ТНА) 1А64 () ТАЦL АА) 1А60 SAKOT 1А32 (Ф ТА) 1А69 (Ф U)	1A35 (∽ THA) 1A64 () TALL AA) 1A60 SAKOT 1A32 (☆ TA) 1A69 (∻ U)

^{1.} This is a very common syllable structure CVC and is probably the most obvious difference between a visual and logical encoding. The subjoined na is encoded visually in an arbitrary (by linguistic standards) but consistent order.

- 2. This word might be expected to differ between encodings but doesn't.
- 3. Only the tone moves here. This is because of the lower vowel that allows the final to be stored after it. This is because there can be a visual difference between a lower vowel before or after a final. But notice that the tone is always at the end.

- 4. What is noticeable here is that medials are now encoded after subjoined characters.
- 5. U+1A58 here is stored and renders before the U+1A63 (9 AA) vowel.
- 6. The linguistics of the proposed encoding here, are a bit of a mess, but the rendering is fine.
- 7. Here U+1A61 (☆ A) is a simple base character. And there are words like [ບໍລິງ > U+1A37 (୰ BA) U+1A55 (☆ MEDIAL RA) U+1A61 (☆ A) U+1A60 SAKOT U+1A3F (↔ YA) U+1A63 (> A) (Lao found in [1]) that also work.
- 8. An example of a polysyllabic cluster with two tone marks.
- 9. Another polysyllabic cluster: boywa. Due to the awkwardness of the combining orders, the original encoding is particularly nasty here.
- 10. There is no difference in encoding here, but the sequence is interesting.

3.2 Under Differentiation

In this section we consider the proposed encoding for issues with under differentiation.

Whereas the original encoding suffered from over differentiation in that two visually identically strings could be spelled differently, this new encoding has problems with under differentiation where necessary visual contrasts cannot be represented in the encoding. The following table shows examples of where there may be problems.

	lanna	khün	thai	original encoding	visual encoding
1.	ສູ	<u>ິ</u> ດີມ	สยะ	1A21 (ຊ) кна) 1A7B (ີ MAI SAM) 1A60 SAKOT 1A3F (ເງ YA) 1A61 (ລ A)	1А21 (ຊ кна) 1А60 закот 1АЗҒ (ເງ үа) 1А7В (ເ маі зам) 1А61 (ล а)
2.	ည့်) ที่ หย้อ ¹ 1A49 (ชา на) 1A60 รลкот 1A3F (อัางล) 1A6C (ชา на) 1A60 รลкот 1A3F (อัางล) 1A6C (อัางล) 1A60 รลкот 1A3F (อัางล) 1A76 (อัางа) 1A76 (อัางล) 1A76 (อัางа) 1A76 (อ			
3.	ର୍ଷ	S	หอย	1A49 (రా на) 1A6C (ప్రాంద) 1A60 SAKOT 1A3F (ప్రాగద)	1А49 (ഗ на) 1А6С (ప్ర оа) 1А60 sakot 1АЗF (ప్ర үа)
4.	ର୍ଣ୍ଣୁଜ	သို့ရ	ฉลอง	1A28 (ఐ cha) 1A7B (హ్ маі sam) 1A60 sakot 1A43 (్లై LA) 1A6C (్లై OA) 1A26 (్లు NGA)	1A28 (ଇ cha) 1A60 sakot 1A43 (ஜ) La) 1A6C (ஜ oa) 1A7B (ॐ mai sam) 1A26 (⇔ nga)
5.	ମ୍ବ	ျှ	คุย	1А23 (Скна) 1А69 (Эри) 1А60 sakot 1А3F (Эрүа)	1А23 (೧ кна) 1А69 (은 U) 1А60 sakot 1А3F (은 үа)
6.	ಜೆಂಗ್	သင့္လိ	อำนาต	1А4В (జృ а) 1А63 (э аа) 1А74 (ీ маі калд) 1А36 (ౖ № №) 1А63 (э аа) 1А60 закот 1А32 (్ర та)	1A4B (ఴ ౚ) 1A63 (ౢ ౚౚ) 1A74 (் Mai kang) 1A36 (ఁ NA) 1A60 sakot 1A32 (☆ ta) 1A63 (ౢ ౚౚ)

- 1. The U+1A7B (MAI SAM) here indicates a syllable break and that the following U+1A60 SAKOT U+1A3F (MAI SAM) here indicates a syllable `chaya`. The visual encoding treats it as if the U+1A60 SAKOT U+1A3F (MAI SAM) were double acting or that the cluster was the last in a reduplicated word. The question is whether there needs to be a visual rendering contrast between the two possible encoding positions of the MAI SAM. Double acting consonants are rare, which would mean that the mai sam in this position is indicating a syllable break. But this mean is really only indicated because of the following U+1A61 (A), which if missing would imply the MAI SAM is a reduplicating mark.
- 2. There is very little difference in encoding between the original and visual encoding, here. The difference here calls for some diacritic reordering during rendering. This is the only example in which the Khuen spelling and encoding is different. It is spelled U+1A49 (𝔅 HA) U+1A60 SAKOT U+1A3F (𝔅 YA) U+1A73 (𝔅 OA ABOVE) U+1A76 (𝔅 TONE-2) with no difference between the original or visual encoding.

- 3. There is no difference between the visual and the original encoding. In Khuen this word would be spelled of U+1A49 (vo HA) U+1A6D (vo OY). In general in Khuen, if there is a sequence of below diacritics and one of them has a tall stem, that is moved to the front.
- 4. This is a near identical problem to example 1.
- 5. This example is awkward. The phonetic and vertically stacked spelling calls for the U+1A69 (ເບ ບ) to come before the subjoined ya, but in Khuen the u is clearly visually located after the subjoined ya. Further, in Lanna, there are examples of the u being below the subjoined ya. The contrast is needed. Compare ບຼງ U+1A49 (ເກ HA) U+1A69 (ເບ ບ) U+1A60 SAKOT U+1A3F (ເບ YA) VS ບົງເລ U+1A49 (ເກ HA) U+1A60 SAKOT U+1A3F (ເບ YA) U+1A69 (ເບ ບ) U+1A76 (ເຕ TONE-2) U+1A26 (ເລ NGA), which do not contrast visually in Khuen with ဟု and ဟု s.
- 6. The difficulty here is that the U+1A63 (◦ aa) has no space and so we have no idea what happens underneath it. Perhaps the easiest solution is to always use U+1A36 (⊂ NA) U+1A63 (◦ AA) immediately adjacent to each other as if there was a ligature spelling for that character.

As seen in the above examples, U+1A7B (SMAI SAM) is problematic. There can be needs for contrasting rendering of the diacritic. In the case where it marks a syllable break or a double acting character, it may well be moved further left, more over the base. The final use as a reduplication mark can call for it to be pushed as far right as possible. The question is, do we need to represent this distinction and if so, how?

The need to contrast spellings where a subjoined consonant may render before or after a lower vowel (as in examples 2 and 3 above) is liable to be somewhat problematic. Because there is a contrast of order allowed in the encoding, it is not possible for the keyboard to resolve that order automatically, the way it can for most diacritics. Lack of visual differentiation is often lacking in fonts and may require users to have to think more deeply to type in the correct order, since a required difference in order results in no visual distinction.

4 Refined Visual Encoding

The proposed encoding in the previous section has some significant weaknesses:

- The encoding contrast between a subjoined character and a below vowel forces users to make a distinction that may not be visible to them or that they do not know how to resolve.
- The encoding is not refined enough to allow fonts that wish to contrast the position of tone marks or upper vowels in relation to spacing subjoined characters, like U+1A60 SAKOT U+1A3F (() YA). Which is supported by the current encoding scheme. For example compare ô

The key aspect of these issues is that neither of them actually affect spelling. The contrasts are purely visual and while there are minimal pairs by analagous environment, in the case of the below vowel, many fonts under differentiate to the point where it is sufficient to under differentiate in plain text. Therefore it is not required for a typist to provide differentiation or for computing processes to rely on contrast in these ways. By purely using order to mark visual differentiation, we are requiring users to type the contrast, which they may not see. An alternative approach to encoding contrast through ordering is to encoding contrast through extra characters. These characters are not there to provide spelling contrast but to enable layout contrasts to be encoded if so desired. As such they are format characters.

We start with a simplified version of the encoding order presented in the previous section:

- 1 $S = C CM? SUB* M? V_p? V_a* V_b? T* F?$
- 2 C = [\u1A20-\u1A54\u1A61\u1A63\u1A64\u1A80-\u1A89\u1A90-\u1A99]
- 3 CM = \u1A5B? \u1A7F*
- 4 SUB = \u1A60 [\u1A20-\u1A4A\u1A53\u1A54] | [\u1A56\u1A57\u1A5C-\u1A5E]
- 5 M = [\u1A5A\u1A55]
- $V_{p} = [\u1A6E-\u1A72]$
- 7 $V_a = [\u1A62\u1A65-\u1A68\u1A6B\u1A73\u1A74]$
- 8 $V_b = [\u1A69-\u1A6A\u1A6C\u1A6D]$
- 9 T = [\u1A58\u1A59\u1A75-\u1A79\u1A7A\u1A7C]
- 10 F = \u1A7B

The difference here is that we have removed the V_b SUB sequence. How are we, therefore, to represent the necessary contrast? We introduce two new format characters:

- U+1A8E TAI THAM SIGN INITIAL
- U+1A8F TAI THAM SIGN FINAL

Both of these characters are of category Cf and are stored as CM in the encoding order. As such they are formatting controls placed following a subjoined character. We start with a discussion of U+1A8F SIGN FINAL. This may follow a subjoined character and indicates that this and all following subjoined characters are finals in the cluster. Thus any following vowels are to be rendered 'before' the marked subjoined character. It also provides indication for those fonts that differentiate the position of vowels or tones in relation to spacing subjoined characters. This table shows two similar words. The first column shows vowel positional contrasts, while the second shows them ignored. The original encoding is given and then a proposed explicit encoding. If the two format characters are removed from the proposed strings, one would expect to see something equivalent to the simple rendering.

Contra	st Simple	Original	Proposed explicit
က်ံ့	က်ွံ့	1А49 (ల на) 1А60 закот 1АЗГ (్ర үа) 1А6С (్ర оа) 1А74 (ీ маі кало) 1А76 (్ толе-2)	1A49 (ల на) 1A60 sakot 1A3F (్రు ya) 1A8E sign initial 1A6C (్రు oa) 1A74 (ీ mai kang) 1A76 (్ tone-2)

ນ໌າ	ဟ်၊	1A49 (లా на) 1A6C (్లు оа) 1A76 (్ толе-2)	1A49 (ம на) 1A60 ѕакот 1АЗҒ (ஐ YA) 1A8F ѕідм ғіма		
~J	5	1A60 SAKOT 1A3F (:) YA)	1A6C (: OA) 1A76 (: TONE-2)		

Where there is no explicit use of format controls, the rendering of the text is up to the heuristics of the font.

The characters are format controls and are optional. For purposes of spell checking, they are ignored. Having said that, they are strongly encouraged due to the value in appropriate rendering contrasts for those fonts that show such contrasts. The rules for the automatic insertion of such marks are beyond the capabilities of a font or a shaping engine, requiring, in some cases, a dictionary. Thus they may be inserted using spell checking techniques.

The second character to consider is U+1A8E SIGN INITIAL. If one has a consonant followed by a spacing subjoined character, one does not know whether the subjoined character is definitely not a final or simply has not been marked as a final and should have been. U+1A8E SIGN INITIAL allows the former case of this character not being a final, to be clearly marked. The use of two such marks overcomes the problem that is also prevalent in Indic text that a system does not know whether the absense of a ZW(N)J indicates the desire for a conjunct or whether a more sophisticated font with more conjuncts would render the text differently from that desired.

Another use for these formatting controls is the positioning of U+1A7B ($\stackrel{*}{\hookrightarrow}$ MAI SAM). As described, there are 3 functions for the character. Some fonts show a positional contrast between the use of mai sam as a reduplication mark and for its other uses. For example:

Contrast	Original	Proposed explicit
ည်	1A48 (ఎ SA) 1A6C (పై OA) 1A60 SAKOT 1A3F (పై YA) 1A7B (పో MAI SAM)	1A48 (ఎ sa) 1A60 sakot 1A3F (్ర ya) 1A8F sign final 1A6C (్ర oa) 1A7B (్ mai sam) 1A8F sign final
သျှ်ရ	1A48 (లు sa) 1A7B (్ маі sam) 1A60 sakot 1A3F (్రు ya) 1A6C (్రు oa) 1A26 (్రు NGA)	1A48 (లు sa) 1A60 sakot 1A3F (్రు ya) 1A8E sign initial 1A6C (్రు oa) 1A7B (్రు mai sam) 1A8E sign initial 1A26 (్ర NGA)
ငည္သို	1A27 (ᢒ HIGH CA) 1A7B (் MAI SAM) 1A60 SAKOT 1A43 () LA) 1A6E (ີ E) 1A65 (் I) 1A60 SAKOT 1A3F () YA)	1A27 (☉ HIGH CA) 1A60 SAKOT 1A43 (☉ LA) 1A60 SAKOT 1A3F (☉ YA) 1A8F SIGN FINAL 1A6E (↔ E) 1A65 (☉ I) 1A7B (☉ MAI SAM) 1A8E SIGN INITIAL

In the last row, one might expect to see U+1A8E SIGN INITIAL following the U+1A60 SAKOT U+1A43 (\bigcirc LA), but this is redundant due to the occurrence of U+1A8F sign final following the next subjoined consonant.

If these characters occur in a sequence context where they are undefined, they are ignored.

There are various reasons for not wanting to use ZW(N)J for these functions:

- ZW(N)J have grown their own specific meanings and behaviours not covered by these characters. For example ZWNJ causes a cluster break in the Universal Shaping Engine.
- It is not clear which function would be covered by which of ZW(N)J. Although an arbitrary choice would be acceptable.
- ZW(N)J can only occur following a halant/sakot in the USE. This would be useable for the subjoined marking but not for the mai sam.
- ZW(N)J may be used for ligation control within Tai Tham, for example: පු U+1A2C (ව NYA) U+1A60 SAKOT U+1A2C (ව NYA) vs වූ U+1A2C (ව NYA) U+1A60 SAKOT U+200D zwJ U+1A2C (ව NYA).

The reasons in favour of using ZW(N)J are:

• No need to introduce extra Cf characters.

• These characters may be useful for resolving similar issues with regard to Old Khmer, although the analysis of that writing system and its encoding has yet to be done. As such they would cross script and using common characters would ease that crossing.

By including two such marks, it leaves the interpretation of unmarked text up to the rendering system. It can be as sophisticated as it likes and where that sophistication breaks down, the format characters can be used to resolve the problem. Or a font may be unsophisticated and require a significant level of text marking, where it wants to show a visual contrast, if at all.

The advantage this system brings is that naive text entry is regularised for purposes of text comparison with the support of whatever level of sophistication a font may bring, while also catering to the needs of those requiring precise encoding for the presentation of high quality text.

Thus the encoding order becomes, with additions in italics:

- 1 S = C CM? SUB* M? V_p ? V_a * V_b ? T* F?
- 2 C = [\u1A20-\u1A54\u1A61\u1A63\u1A64\u1A80-\u1A89\u1A90-\u1A99]
- 3 CM = \u1A5B? \u1A7F*
- 4 SUB = \u1A60 [\u1A20-\u1A4A\u1A53\u1A54] [\u1A8E\u1A8F]? | [\u1A56\u1A57\u1A5C-\u1A5E]
- 5 M = [\u1A5A\u1A55]
- 6 $V_p = [\u1A6E-\u1A72]$
- 7 V_a = [\u1A62\u1A65-\u1A68\u1A6B\u1A73\u1A74]
- 8 V_b = [\u1A69-\u1A6A\u1A6C\u1A6D]
- 9 T = [\u1A58\u1A59\u1A75-\u1A79\u1A7A\u1A7C]
- 10 $F = \1 R F = \1 R F = \1 R F$

4.1 Issues

There are outstanding issues that may or may not need to be addressed:

• Contractions where one might see something like: ♣] U+1A48 (▷ SA) U+1A69 (♠ U) U+1A60 SAKOT U+1A43 (♠ LA) U+1A66 (♠ II). Awaiting evidence of such text, although it may be argued that it is beyond scope and into the realms of calligraphy. The issue can be addressed either by allow U+1A8F SIGN FINAL following the U+1A66 vowel or by the introduction of yet another format control for the same purpose.

5 Implementation

Having proposed a visual encoding, we now examine how it might be implemented. The great value of a visual encoding is consistency, with on visual form having one encoding, which makes implementation easier. Although it can push other problems onto the font. Without it, identifying inconsistent sequences is nearly impossible. We first examine the issue of rendering and shaping and then look at the process of normalizing the order of diacritics follow base characters, both as part of text normalization but also as part of keyboarding, where we allow a relatively free typing order, but can ensure consistent output.

5.1 Shaping

A primary text rendering implementation is that of OpenType. In overly broad outline, this involves the structuring of text into clusters, based on regular expressions, much as we have been doing, and then applying lookups over the glyphs in a cluster. The primary structuring system in OpenType for Tai Tham is the Universal Shaping Engine. Can we fit the requirements for Tai Tham into the USE with minimal changes?

One advantage that rendering has over data entry is that it can be more relaxed about things occuring in the wrong order. So long as there is a visual difference between two different sequences, there is no need for the shaping engine to intervene and insert some kind of visual distinction (usually by adding a dotted circle somewhere in the sequence). Thus there is no need to stop particularly egregious sequences like U+1A61 $_$ U+1A61 $_$ A or the like.

Mapping from the visual encoding onto the USE is first a matter of mapping visual encoding structural classes onto USE classes:

```
\begin{array}{l} \mathsf{C} \rightarrow \mathsf{B} \\\\ \mathsf{CM}(\ulase\ulase) \rightarrow \mathbf{CF} \\\\ \mathsf{CM}(\ulase\ulase) \rightarrow \mathsf{CMBlw} \\\\ \ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ulase\ul
```

5.1.1 USE Structure

The proposed structure extends the current USE syllable structure as marked in bold:

```
[< R | CS >] < B [VS] | GB > (CMAbv)* (CMBlw)* (< H B [VS] | SUB > (CMAbv)* (CMBlw)* [CF])*
[MPre] [MAbv] [MBlw] [Mpst]
(VPre)* (VAbv)* (VBlw)* (VPst)*
(VMPre)* (VMAbv)* (VMBlw)* (VMPst)*
(FAbv)* (FBlw)* (FPst)* [ FM [CF] ]
```

5.1.2 OpenType Implementation

After the grapheme clusters have been identified there are still issues with ordering diacritics such that diacritic attachment tree is stored in the right order.

• U+1A74 MAI KANG may be stored following a spacing vowel, but in some cases need to be reordered onto the previous consonant, before even a tone on that consonant. But this is not true for all fonts.

- U+1A7F CRYPTOGRAMMIC DOT is stored late in the sequence and needs to be shifted to immediately following the initial consonant glyph.
- Subjoined consonants followed by U+1A8F SIGN FINAL need to be moved after the vowels. Alternatively what a vowel attaches to may be dependent on the presence of U+1A8F SIGN FINAL.
- Final superscript consonants (U+1A59, U+1A5A) will occur after the tone marks they should render below.

5.2 Normalization

How can we enforce a text to conform to the proposed encoding order described above? Unlike for rendering, the requirement here is to reorder a sequence of characters from an arbitrary keying order into a canonical spelling order, as described previously.

5.2.1 Canonical Ordering

There is a constraint on our preferred storage order and that is that marks with a Canonical Combining Class sort relative to their value after a character with CCC of 0. We must ensure that text that is ordered according to our order is not modified if normalized according to the Unicode Normalization algorithm. Thankfully, for Tai Tham, there is no difference between NFC and NFD.

It is most unfortunate that Tai Tham was landed with various diacritics being given somewhat arbitrary (in ordering terms) canonical combining classes. The overloading of classes as ordering means we have some real difficulties in the ordering of Tai Tham. And since it is impossible to fix these where they should be fix: in the standard. We need to go through some mental gymnastics to handle the resulting broken state of Tai Tham. The non-zero CCC values are:

\u1A60	9
\u1A7F	220
[\u1A75-\u1A79]	230
[\u1A7A-\u1A7C]	230

Non-zero CCC diacritics immediately preceding U+1A60 SAKOT are a problem because they are reordered immediately after the Sakot, between the Sakot and its following consonant. In the current scheme the only problem character is U+1A7F CRYPTOGRAMMIC DOT. By categorising this as CM it may occur immediately before a Sakot. But a cryptogrammic dot only modifies a base consonant. Therefore, on that understanding, it could occur anywhere in the cluster, even after a subjoined consonant, on the understanding that it is reordered during shaping to immediately after the consonant. Another constraint is that it has a CCC value less than that of VMAbv. Given we also want multiple of them, this limits us to a V of some kind or VMPre. We somewhat arbitrarily choose VMPre, since this is not a vowel. But then, neither is it a vowel modifier.

5.2.2 Sorted Ordering

Part of UTS 35, the CLDR file format description for locale information, describes keyboard layouts, and as part of that there is an algorithm given for reordering marks from their keyed order into the required underlying order. This algorithm may be simply extended to support text canonicalisation. There is no need to handle pre-vowels in such an algorithm.

One approach is to allocate every character in a string a given ordering number. Groups are those between a character with an order of 0 and the next character with order 0, while only including the first character with an order of 0. The characters are then ordered first by their order number and second by their current order (or index) in the string. The order numbers are not necessarily given by a simple lookup, but characters may given an order number based on a contextual match. All characters in the match are given the same order number (and so will keep their relative order). Character sequences are matched using a regular expression, which may include look ahead or look behind assertions. Unfortunately for most regular expression engines,

the look behind assertions are not necessarily fixed length (and so in implementation may need to be extracted and processed as another constraint on the match). All unmatched characters are given an order of 0. Negative order values indicate secondary sort relationships.

Character classes		Order
[\u1A20-\u1A54\u1A61\u1A63\u1A64]	\rightarrow	0
[\u1A8E\u1A8F]	→	-2
[\u1A5B\u1A7F]	\rightarrow	-3
\u1A60[\u1A20-\u1A4A\u1A4C]	\rightarrow	10
\u1A60	\rightarrow	10
[\u1A56\u1A57\u1A5C-\u1A5E]	\rightarrow	15
\u1A5A	\rightarrow	17
\u1A55	\rightarrow	20
[\u1A6E-\u1A72]	\rightarrow	60
[\u1A62\u1A65-\u1A68\u1A6B\u1A73\u1A74]	\rightarrow	65
[\u1A69\u1A6A\u1A6C\u1A6D]	→	70
[\u1A75-\u1A79]	→	85
[\u1A58\u1A59\u1A7A\u1A7C]	\rightarrow	105
\u1A7B	→	115

The values of the numbers are unimportant just their relative value, except the value of 0 which indicates the start of a new sequence order (and is the default for all unspecified codepoints, making the first list strictly redundant). Notice that at this point all reference to any guiding regular expressions are lost. The categories in such regular expressions are only helpful in arriving at a relative order.

The biggest problem with this implementation is that any tone [\u1A75-\u1A79] occurring after a spacing vowel is not reordered before it, as it should. This is a common problem with people often typing 'am' followed by a tone mark (which needs to be reordered across the 'aa' vowel). To support this one needs greater complexity. We, in effect, treat the spacing vowels [\u1A61\u1A63\u1A64] as diacritics onto a consonant based grapheme cluster, for the purposes of ordering.

Character classes		Order
[\u1A20-\u1A54]	\rightarrow	0
[\u1A5B\u1A7F]	\rightarrow	-2
[\u1A8E\u1A8F]	\rightarrow	-3
\u1A60[\u1A20-\u1A4A\u1A4C]	\rightarrow	10
\u1A60	\rightarrow	10
\u1A60[\u1A75-\u1A79][\u1A20-\u1A4A\u1A4C]	→	10, 85, 10 ³
[\u1A56\u1A57\u1A5C-\u1A5E]	\rightarrow	15
\u1A5A	\rightarrow	17
\u1A55	\rightarrow	20
[\u1A6E-\u1A72]	\rightarrow	60
[\u1A62\u1A65-\u1A68\u1A6B\u1A73\u1A74]	\rightarrow	65
[\u1A69\u1A6A\u1A6C\u1A6D]	\rightarrow	70
[\u1A75-\u1A79]	→	85

³ This addressed normalized text and specifies an order for each character in the matched sequence.

Character classes [\u1A58\u1A7B](?=[\u1A55-\u1A57\u1A59\u1A5E\u1A62\u1A65-\u1A68\u1A6D-\u1A7A\u1A7C-\ u1A7F]*[\u1A61\u1A63\u1A64])	→	Orde 90
[\1A61\u1A63\u1A64]	\rightarrow	100
(?<=[\u1A61\u1A63\u1A64][\u1A55-\u1A5E\u1A62\u1A65-\u1A68\u1A6D-\u1A7F]*)\u1A60[\u1A20-\ u1A49]	→	105
(?<=[\u1A61\u1A63\u1A64][\u1A55-\u1A5E\u1A62\u1A65-\u1A68\u1A6D-\u1A7F]*)\u1A60	\rightarrow	105
(?<=[\u1A61\u1A63\u1A64][\u1A55-\u1A5E\u1A62\u1A65-\u1A68\u1A6D-\u1A7F]*)[\u1A56\u1A57\ u1A5C-\u1A5E]	→	107
(?<=[\u1A61\u1A63\u1A64][\u1A55-\u1A5E\u1A62\u1A65-\u1A68\u1A6D-\u1A7F]*)[\u1A65-\u1A68\ u1A6B\u1A73\u1A74]	→	110
(?<=[\u1A61\u1A63\u1A64][\u1A55-\u1A5E\u1A62\u1A65-\u1A68\u1A6D-\u1A7F]*)[\u1A61-\u1A64\ u1A69\u1A6A\u1A6C\u1A6D]	→	112
[\u1A58\u1A59\u1A7A\u1A7C]	\rightarrow	115
\u1A7B	\rightarrow	117

The grey cells indicate characters that may be followed by a tertiary sequence. The kline star * is limited for implementation purposes at 4 occurrences. The (?<= special regular expression syntax indicates a look behind constraint. In effect it says that the characters to receive the order must be preceded by the given subexpression.

In the order 110 vowel list, we remove 1462. U+1A62 (6 A) has a secondary linguistic function as a final k. But this character does not occur aftger U+1A63 (9 AA) and so it helps if a keyboard reorders them.

There is a possible ambiguity that if someone were to type a U+1A5B (© RATHA OR PA) following a U+1A7B (© MAI SAM), with no intervening base or subjoined consonant, the U+1A5B would be incorrectly ordered. But this is highly unlikely and users can be trained to resolve this corner case.

Missing from this list are CGJ, ZWJ, ZWNJ and variation selectors which all would take a tertiary order of - 1.

5.3 Format Controls

A basic keyboarding order has already been described in the previous section. But with the move away from a logical encoding towards something much more regular and visual, more control load is placed on the keyboard. In particular the question of whether the keyboard can insert format controls automatically to reduce the interpretive load on the font which has much more limited processing capability.

It is of course possible to write a language specific keyboarding application that can even incorporate known standard spellings to help resolve format control insertion. But using more generic approaches can still yield helpful results. For example, it is possible, with care, for a keyboard to incorporate keying order to allow a user to specify whether, for example, a lower vowel should render before or after a subjoined consonant. It may not be 100% successful but can give a good default. In addition a good keyboard layout should provide the user with the capability to insert format controls if they so desire.

Where no format controls occur in a string, there are some basic principles that can be used to help provide some basic rendering. These rules should not be considered in any way normative. Consider S to be a subjoined substring: U+1A60 sakot followed by a consonant. Likewise MS is U+1A7B (:: MAI SAM). The subscript f indicates the marking of the text with a U+1A8F SIGN FINAL, and i with a U+1A8E SIGN INITIAL.

 $\begin{array}{l} U+1A49 \; S \rightarrow S_i \\ S \; S \rightarrow S_i \; S_f \\ S \; S_i \rightarrow S_i \; S_i \\ S_f \; \ldots \; MS \rightarrow S_f \; \ldots \; MS_f \\ S_i \; \ldots \; MS \rightarrow S_i \; \ldots \; MS_i \end{array}$

 $\label{eq:sigma} \begin{array}{l} S \; [\u1A65-\u1A68] \rightarrow S_f \; [\u1A65-\u1A68] \\ Fallback \; S \rightarrow S_i \end{array}$

6 Acknowledgements

Thanks go to Richard Wordingham for his careful review of this document and some examples. He also suggested a much simpler grapheme cluster structure. Andrew Glass has worked hard to help integrate Tai Tham into the USE and opened my eyes to the idea that would could change the core sequence.

Thanks go to Payap University Linguistics Institute, Chiang Mai, Thailand, under whose auspices this work is done.

7 Bibliography

Due to the multiplicity of scripts involved, entries are sorted by date rather than author.

1.	อรุณรัตน์ วิเชียรเขียว และคณะ (Arunradt Wichienkiaw et al) 1996 ธอตุตูสู่สิ่งวตุธโละกาซีเอาฏิโตโองรู เพจนานุกรมศัพท์ลานนาเฉพาะคำที่ปราฏในใบ ลาน The Northern Thai Dictionary of Palm-Leaf Manuscripts) ISBN 974-7047- 77-2			
2.	สนั่น ธรรมธิ (Sanan Tamti) 2001. ยุตตส Nuengnaiaksonlanna) ISBN 974-657	สาระ หนึ่งในอักษล้านนา (Yutotsara 7-332-2		
3.	อุดม รุ่งเรืองศรี (Udom Rungruengsri) 2004. ตอต นา-ไทย ฉบับแม่ฟ้าหลวง Phojananukr 974-386-044-4	ഭൂണ്ര്രാള് പ്രാദ്ദേദ്ദ്യ മഗ്ഗദേഷ് പ്രാസ്തര് (พจนานุกรมลาน rom Lanna-Thai: Chababmefahluang) ISBN		
4.	เกษม ศิริรัตน์พิริยะ (Kasem Siritpiriya) 2005. (Duamuang: Kanrianphasalannaphan	ตั๋วเมือง การเรยนภาษาลานนาผ่านโครงสร้างคำ khrongsangkham) ISBN 974-9942-00-0		
5.	Everson M. and Hosken M. 2006. Prop of the UCS (Unicode L2/06-258, ISO	osal for encoding the Lanna script in the BMP /IEC JTC1/SC2/WG2 N3121)		
6.	ၿಀိဘ၁ೲ၈၁ ဘိက္ခု (Patibhano Bhikku) 2008. ၀၉ဇဇ 623-711-6	က္ငံဝဲမျ-ဘ္ရက္သြိ (Tai-English Dictionary) ISBN 974-		
7.	2017. ภาคพันธสัญญาใหม่ Phrakhristthamkh	ప్లుదత్రి దర్జి స్వాద్దామిజుంగ్రం (พระคริสตธรรมคัมภีร์ namphi Phakphanotsanyamai)		
8.	ရခ္ခန်ေနေ့ စေတ္ရံဘ္သတ (Jaichechen Tengok) (Aphithansabwaihanble phaasaakhuer	ဘဘိဓါငည်ြင်ဝဟဒ္ဒဧေဂါ ဘဝသဝခ်ို n)		
9.	Unicode Technical Committee Unic Data Markup Language (LDML) <u>http</u> <u>keyboards.html#Element_reorder</u>	ode Technical Standard #35: Unicode Locale <u>o://unicode.org/reports/tr35/tr35-</u>		