

PAN Localization Project  
 Science Technology and Environment Agency of Lao PDR  
 National University Computer and Emergency Science of Pakistan  
 International Development Research Center

## Syllabification of Lao Script for Line Breaking

Phonpasit PHISSAMAY, National Project Director  
 Valaxay DALALOY, National Project Coordinator  
 Nadir DURRANI, Project consultant  
 Dr. Sarmad Hussain, Regional Project Director  
 (Ms. Oulaiphone SILIMASAK);  
 (Ms. Chitthaphone CHANHSILILATH);  
 (Mr. Khamphay INTHARA);  
 (Mr. Phonexay VILAKONE);

### I. Definition:

Lines are broken as result of one of two conditions. The first condition is the presence of an explicit line breaking character. The second condition results from a formatting algorithm having selected among available line break opportunities; ideally the chosen line break results in the optimal layout of the text.

Different formatting algorithms may use different methods to determine an optimal line break. For example, simple implementations consider a single line at a time, trying to find a *locally optimal* line break. A basic, yet widely used approach is to allow no compression or expansion of the inter-character and inter-word spaces and consider the longest line that fits. When compression or expansion is allowed, a locally optimal line break seeks to balance the relative merits of the resulting amounts of compression and expansion for different line break candidates.

When expanding or compressing inter-word space according to common typographical practice, only the spaces marked by U+0020 SPACE, U+00A0 NO-BREAK SPACE, and U+3000 IDEOGRAPHIC SPACE are subject to compression, and only spaces marked by U+0020 SPACE, U+00A0 NO-BREAK SPACE, and occasionally spaces marked by U+2009 THIN SPACE are subject to expansion. All other space characters normally have fixed width. When expanding or compressing inter-character space the presence of U+200B ZERO WIDTH SPACE or U+2060 WORD JOINER is always ignored.

Three principal styles of context analysis determine line break opportunities.

1. *Western* — spaces and hyphens are used to determine breaks: It is commonly used for scripts employing the space character. Hyphenation is often used with space-based line breaking to provide additional line break opportunities. However, it requires knowledge of the language and in addition; it may need user interaction or overrides.
2. *East Asian* — lines can break anywhere, unless prohibited: In these scripts, lines can break anywhere, except before or after certain characters. The precise set of prohibited line breaks may depend on user preference or local custom and is commonly tailorable.
3. *South East Asian* — line breaks require morphological analysis: The third style is used for scripts such as Lao, which do not use spaces, but which restrict word-breaks to syllable boundaries, the determination of which requires knowledge of the language comparable to that required by a hyphenation algorithm. Such an algorithm is beyond the scope of the Unicode Standard.

## II. Syllabification of Lao Script

### Introduction

Syllable is a unit of spoken language, which may have a common meaning or not. A unit of spoken language may have one single syllable or many syllables. Additional complexity is introduced by lack of a space character. Though humans can process multiword string while reading and extract words from it, this process is very difficult for computers. However, unless words can be separated, it is impossible to perform even the simple task of determining how to break at the end of a typed line when characters exceed the line length, without the possibility of breaking between a word and worse between syllables.

In Lao language, the breaking of the text flow after each word (like in English) is not common. One possible way line breaking can be achieved in Lao is to use a Lao lexicon which is currently incompleteness. Lao collation is complex because it does not sort on key-press order but on basis of its intricate syllable structure. Thus, before any processing is done, a Lao character string has to be syllabified. This paper develops a Lao string syllabification algorithm. This algorithm is essential for processing basic input character string to enable more advanced language processing including searching, sorting, line breaking and lexical development.

### Lao Character Set

Lao syllable structure contains characters and marks for consonants, vowels and tones. Table 1 below lists these possible characters and marks. The marks are combining characters and are shown adjacent to 'x', latter representing a consonant. Character names and their Unicode is also given (ref. to Unicode)

Table 1: Lao Characters

(a) Vowels ('x' is a placeholder for a consonant character)

Vowel	Name	Vowel	Name
ⵀ	LAO VOWEL SIGN A (0EB0)	ⵁ	LAO VOWEL SIGN UU (0EB9)
ⵂ	LAO VOWEL SIGN MAI KAN (0EB1)	ⵃ	LAO VOWEL SIGN MAI KON (0EBB)
ⵄ	LAO VOWEL SIGN AA (0EB2)	ⵅ	LAO SEMIVOWEL SIGN NYO(0EBD)
ⵆ	LAO VOWEL SIGN AM (0EB3)	ⵇ	LAO VOWEL SIGN E (0EC0)
ⵈ	LAO VOWEL SIGN I (0EB4)	ⵉ	LAO VOWEL SIGN EI (0EC1)
ⵊ	LAO VOWEL SIGN II (0EB5)	ⵋ	LAO VOWEL SIGN O (0EC2)
ⵌ	LAO VOWEL SIGN Y (0EB6)	ⵍ	LAO VOWEL SIGN AY (0EC3)
ⵎ	LAO VOWEL SIGN YY (0EB7)	ⵏ	LAO VOWEL SIGN AI (0EC4)
ⵐ	LAO VOWEL SIGN U (0EB8)	ⵑ	LAO NIGGAHITA (0ECD) (It is a Vowel Sign "OR")

## (b) Consonants ('x' marks a placeholder for a consonant character)

Cons.	Name	Cons.	Name
ກ	LAO LETTER KO (0E81)	ຝ	LAO LETTER FO TAM (0E9D)
ຂ	LAO LETTER KHO SONG (0E82)	ພ	LAO LETTER PHO TAM (0E9E)
ຄ	LAO LETTER KHO TAM (0E84)	ຟ	LAO LETTER FO SONG (0E9F)
ງ	LAO LETTER NGO (0E87)	ມ	LAO LETTER MO (0EA1)
ຈ	LAO LETTER CO (0E88)	ຢ	LAO LETTER YO (0EA2)
ຊ	LAO LETTER SO TAM (0E8A)	ຮ	LAO LETTER LO LING (0EA3)
ຢ	LAO LETTER NYO (0E8D)	ລ	LAO LETTER LO LOOT (0EA5)
ດ	LAO LETTER DO (0E94)	ວ	LAO LETTER WO (0EA7)
ຕ	LAO LETTER TO (0E95)	ສ	LAO LETTER SO SONG (0EAA)
ຖ	LAO LETTER THO SONG (0E96)	ຫ	LAO LETTER HO SONG (0EAB)
ທ	LAO LETTER THO TAM (0E97)	ອ	LAO LETTER O (0EAD)
ນ	LAO LETTER NO (0E99)	ຮ	LAO LETTER HO TAM (0EAE)
ບ	LAO LETTER BO (0E9A)	ຸ	LAO VOWEL SIGN LO (0EBC) (It functioning as consonant)
ປ	LAO LETTER PO (0E9B)	ໜ	LAO HO NO (OEDC)
ຜ	LAO LETTER PHO SONG (0E9C)	ໝ	LAO HO MO (OEDD)

## (c) Tones ('x' is a placeholder for a consonant character)

Tone	Name	Tone	Name
ˊx	LAO TONE MAI EK (0EC8)	ˋx	LAO TONE MAI TI (0ECA)
ˋx	LAO TONE MAI THO (0EC9)	ˊx	LAO TONE MAI CATAWA (0ECB)

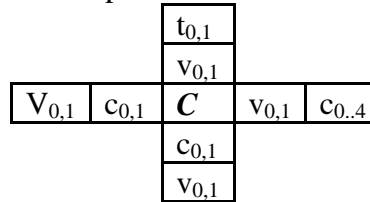
## (d) Sign ('x' is a placeholder for a alternate consonant character)

Sign	Name	Sign	Name
ʼ	LAO ELLIPSIS (OEAF)	໘	LAO KOLA (OEC6)
ˊx	LAO CANCELLATION MARK (0ECC)		

Table 1 shows there are 18 vowel marks and characters, 30 character marks and characters, 4 tonal marks and 3 special symbols. Vowels can occur before, above, below a consonantal character or on the baseline. Characters occur on the baseline. Tonal marks always occur on top of consonantal characters.

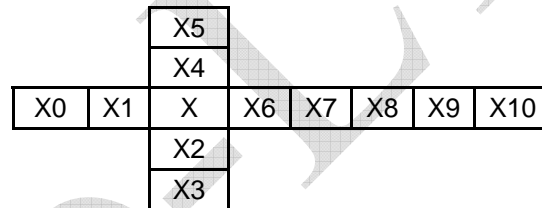
## Basic Syllable Structure

Lao language writing system is based on a central or nuclear consonantal character. This consonant may have optionally vowel character or marks around it (before, after, above or below). In addition, this nuclear consonantal character may also have optional a tonal mark above it and optionally more consonantal characters following it. This is illustrated in Figure 1 below. Capital C indicates the nuclear consonantal character. The subscripts indicate that all are optional except the nucleus C.



## Constraints on Syllable Structure

Not all characters and marks are allowed at all the placeholders shown in Figure 1. Figure 2 and Table 2 below collectively further classify these consonants, vowels and tones to indicate these limitations. The number *n* in X<sub>n</sub> represents the key-press order or typing sequence of these characters and marks relative to each other, except X (which is typed between X<sub>1</sub> and X<sub>2</sub>). Thus tone is typed after the initial vowel X<sub>0</sub>, consonant X<sub>1</sub>, nucleus X, and vowel mark X<sub>3</sub> and X<sub>4</sub> (if X<sub>0</sub>, X<sub>1</sub> and X<sub>2</sub> are not null; X can never be null).



The following table is shown the characters exist in each structure position

X0	X1	X				X2	X3	X4	X5	X6	X7	X8	X9	X10
ໄX <sub>1</sub>	ຫ	ກ <sub>01</sub>	ຂ <sub>02</sub>	ຄ <sub>03</sub>	ງ <sub>04</sub>	ຊ <sub>1</sub>	ຊ <sub>1</sub>	ໂ <sub>1</sub>	ໂ <sub>1</sub>	ວ <sub>1</sub>	ຮ <sub>1</sub>	ກ <sub>1</sub>	ຈ <sub>1</sub>	ຯ <sub>1</sub>
ໄໄX <sub>2</sub>		ຈ <sub>05</sub>	ສ <sub>06</sub>	ຊ <sub>07</sub>	ຍ <sub>08</sub>	ຮ <sub>2</sub>	ຊ <sub>2</sub>	ໂ <sub>2</sub>	ໂ <sub>2</sub>	ອ <sub>2</sub>	ຮ <sub>2</sub>	ງ <sub>2</sub>	ສ <sub>2</sub>	ງ <sub>2</sub>
ໂໂX <sub>3</sub>		ດ <sub>09</sub>	ຕ <sub>10</sub>	ຖ <sub>11</sub>	ທ <sub>12</sub>	ວ <sub>3</sub>		ໂ <sub>3</sub>	ໂ <sub>3</sub>	ງ <sub>3</sub>	ຮ <sub>3</sub>	ຍ <sub>3</sub>	ຊ <sub>3</sub>	ຯ <sub>3</sub>
ໂໂX <sub>4</sub>		ນ <sub>13</sub>	ປ <sub>14</sub>	ປ <sub>15</sub>	ຜ <sub>16</sub>	ລ <sub>4</sub>		ໂ <sub>4</sub>	ໂ <sub>4</sub>			ດ <sub>4</sub>	ພ <sub>4</sub>	
ໂໂX <sub>5</sub>		ຝ <sub>17</sub>	ພ <sub>18</sub>	ຟ <sub>19</sub>	ມ <sub>20</sub>			ໂ <sub>5</sub>				ນ <sub>5</sub>	ຟ <sub>5</sub>	
		ຢ <sub>21</sub>	ຮ <sub>22</sub>	ລ <sub>23</sub>	ວ <sub>24</sub>			ໂ <sub>6</sub>				ມ <sub>6</sub>	ລ <sub>6</sub>	
		ຫ <sub>25</sub>	ອ <sub>26</sub>	ຮ <sub>27</sub>	ໝ <sub>28</sub>			ໂ <sub>7</sub>				ປ <sub>7</sub>		
		ໝ <sub>29</sub>										ວ <sub>8</sub>		

- X0 represents a vowel which occurs before the nuclear consonant. It is can always defined the beginning of syllable.
- X1 is a combination consonant which comes before the nuclear consonant, only if nuclear consonant is one of {ງ, ຍ, ລ, ວ, ູ, ມ, ນ, ສ}.
- X is represents the nuclear consonants.
- X2 is a combination consonant which comes after the nuclear consonant, which placing under or next to the nuclear consonant.
- X3 is represents a vowel which occurs under the nuclear consonant.
- X4 is represents a vowel which occurs upper the nuclear consonant.
- X5 is represents a tone marks which occurs upper the nuclear consonant or upper vowels.
- X6 is represents consonant vowel, which occurs after nuclear consonant. It functions when the syllable doesn't have any vowels. And it always exists with X8.
- X7 is represents an after vowels. However X7<sub>1</sub> is always represents the end of syllable and it is never exist with tone mark.
- X8 is represents alternate consonants.
- X9 is represents alternate consonant to pronounce foreign language, it always exist with X10<sub>3</sub>.
- X10 represents a sign mark. X10 is always occurs at the end of syllable, but mostly people keep it separate from syllable.

The nuclear consonant always exists in syllable with some vowel or alternate consonants. The position at which vowels appear can guide to define the beginning and end of each syllable. Using following rules we can find potential syllable boundary.

### 1. For $x0_1 = \epsilon x$

- 1.1.  $\epsilon x = x0_1 (x1)X(x2)(x5)(x8)(x9:x10)$
- 1.2.  $\epsilon \tilde{x}, \epsilon \tilde{x} = x0_1 (x1)X(x2)x4_{1-2}(x5)(x8) (x9:x10)$
- 1.3.  $\epsilon \tilde{x}\Theta, \epsilon \tilde{x}\Theta = x0_1 (x1)X(x2)x4_{3-4}(x5) x6_2 (x8) (x9:x10)$
- 1.4.  $\epsilon x\epsilon, \epsilon x\epsilon = x0_1 (x1)X(x2)(x7_2)x7_1$
- 1.5.  $\epsilon \tilde{x}\eta = x0_1 (x1)X(x2) x4_6 (x5) x7_2$
- 1.6.  $\epsilon \tilde{x}(x8) = x0_1 (x1)X(x2) x4_7 (x5) x8 (x9:x10)$
- 1.7.  $\epsilon \tilde{x}(x8) = x0_1 (x1)X(x2) x4_7 (x5) x8 (x9:x10)$
- 1.8.  $\epsilon x\jmath, \epsilon \tilde{x}\jmath = x0_1 (x1)X(x2) (x4_7)(x5)x6_3$

### 2. For $x0_2 = \mu x$

- 2.1.  $\mu x = x0_2 (x1)X(x2)(x5)(x6)(x8) (x9:x10)$
- 2.2.  $\mu x\epsilon = x0_2 (x1)X(x2)x7_1$
- 2.3.  $\mu \tilde{x}(x8) = x0_2 (x1)X(x2) x4_7 (x5) x8 (x9:x10)$

### 3. For $x0_3 = \imath x$

- 3.1.  $\imath x, \imath x\omega = x0_3 (x1)X(x2)(x5)(x8) (x9:x10)$
- 3.2.  $\imath x\epsilon = x0_3 (x1)X(x2)x7_1$
- 3.3.  $\imath \tilde{x}\omega, \imath \tilde{x}\epsilon = x0_3 (x1)X(x2)x4_7(x5) x8_{3:8}$

4. **For  $x_0 = \text{ໂ}$**   $x = x_0(x_1)X(x_2)(x_5)(x_6)(x_9:x_{10})$
5. **For  $x_0 = \text{ໃ}$**   $x = x_0(x_1)X(x_2)(x_5)(x_6)$
6. **For  $x_3 = \text{ຊ}$  &  $\text{ຊ}$**   $x = (x_1)X(x_2)x_3(x_5)(x_8)(x_9:x_{10})$
7. **For  $x_{4-4} = \text{ຮ}$  &  $\text{ຮ}$  &  $\text{ຮ}$  &  $\text{ຮ}$**   $x = (x_1)X(x_2)x_{4-4}(x_5)(x_8)(x_9:x_{10})$
8. **For  $x_4 = \text{ຮ}$**   $x = (x_1)X(x_2)x_4(x_5)(x_7)(x_9:x_{10})$
9. **For  $x_4 = \text{ຮ}$** 
  - 9.1.  $\text{ຮ}(x_8) = (x_1)X(x_2)x_4(x_5)x_8(x_9:x_{10})$
  - 9.2.  $\text{ຮ} = (x_1)X(x_2)x_4(x_5)x_6x_7$
10. **For  $x_4 = \text{ຮ}$**   $x = \text{ຮ}, \text{ຮ} = (x_1)X(x_2)x_4(x_5)(x_6)x_8(x_9:x_{10})$
11. **For  $x_6 = \text{ຂ}$  &  $\text{ຂ}$  &  $\text{ຢ}$**   $x = (x_1)X(x_2)(x_5)x_6x_8(x_9:x_{10})$
12. **For  $x_7 = \text{ຂ}$**   $x = (x_1)X(x_2)(x_5)x_7$
13. **For  $x_7 = \text{ຢ}$**   $x = (x_1)X(x_2)(x_5)x_7(x_8)(x_9:x_{10})$
14. **For  $x_7 = \text{ຢ}$**   $x = (x_1)X(x_2)(x_5)x_7(x_9:x_{10})$

## Lao Syllabification Algorithm

Lao Syllabification algorithm can be defined as follows:

1. Traverse through the input array and mark syllable boundary as soon as you encounter a punctuation mark, space or a character that does not belong to Lao character set. Example:

ຄົນ | 10 | ລາວ

Syllable (boundaries due to non-Lao character)

2. Filter out Lao characters out of input array leaving behind punctuation marks, spaces, and non-Lao characters. Example:

ຄົນ 10 ລາວ → ຄົນລາວ

3. Reorder character in case of typing variations (Sometime people maybe typed X5 before X2, X3, X4 and X4 maybe typed before X2). Example:

ກຸ່ມນີ້

ກ	໌	ຸ	ມ	ນ	້	ີ
X	X5 <sub>1</sub>	X3 <sub>1</sub>	X8	X	X5 <sub>2</sub>	X4 <sub>4</sub>

Syllable In this case: X3<sub>1</sub> should typed before X5<sub>1</sub> and X4<sub>4</sub> should typed before X5<sub>2</sub>

ກ	ຸ	໌	ມ	ນ	ີ	້
X	X3 <sub>1</sub>	X5 <sub>1</sub>	X8	X	X4 <sub>4</sub>	X5 <sub>2</sub>

4. Mark each character in run with all possible  $X_n$  values it can take. Example:

ຄົນລາວ	ຄ	ົ	ນ	ລ	າ	ວ
	X	X4 <sub>6</sub>	X	X	X7 <sub>2</sub>	X
			X8	X2		X2
				X9		X6
						X8

5. Use rules discussed in previous section to find out syllable boundaries.
- In case if more then one condition suggests a syllable boundary chose the one with longest run.
  - In case none of the conditions suggest syllable boundary try including last character from previous syllable to current syllable.
  - Test all the conditions for previous syllable because removing last character might make it invalid. If it is still a valid syllable then try conditions starting from newly added character, other wise restore the previous syllable and skip first character from current syllable and try testing from next character. Keep skipping characters till find a valid syllable boundary.
  - If program finds boundary for the new syllable it should continue naturally other wise restore previously disturbed syllable by putting back the removed character. Skip the current characters and re-test the conditions, keep skipping unless you find valid syllable.

For example: In case of ເຮືອນີ້ program should identify syllable boundary like

ເຮືອນ | ື້ then it would be unable to detect any condition working for ື້ so it should shift syllable boundary one step back and try to include ‘ນ’ with ື້. But before testing this new syllable it should test the previous syllable ເຮືອ and find if it is still a valid syllable which in the current scenario is a valid syllable. Now try testing for ນີ້, this is also a valid syllable so program would continue naturally and go for testing next syllable.

Now consider this example ເຮືອນີ້ program suggests syllable boundary after ເຮື now as we try test ນີ້ we would find none of the conditions working so we try to remove ື from previous syllable and test if ເຮ is still a valid. In case it is we test ື ນີ້ we still find none of the conditions working so we restore previous syllable ເຮື skip



current character ະ and carry on testing conditions from ັ. Keep skipping characters unless you find valid syllable. In the current example you only have to skip it once.

6. Traverse till the end of array.
7. Put the Lao characters back into original array that have punctuation marks and other non-Lao characters.

### Example 1: ໂຄງການ 10N ພັດທະນາພາສາລາວ

- Step 1 & 2: Traverse through the input array mark syllable boundary where counter non-Lao character. Extract only Lao text and put it in new array.

ໂຄງການ | ພັດທະນາພາສາລາວ

- Step 3: Look for possible re-ordering of text. For example user might type ‘, X5 before ‘, X47 so we must do the re-ordering.
- Step 4: Mark each character with possible X<sub>N</sub>.

ໄ	ຄ	ງ	ກ	າ	ນ	ຮ	ພ	ັ	່	ດ	ທ	ະ	ຮ	ນ	າ	ພ	າ	ສ	າ	ລ	າ	ວ
A0	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	A11	A12	A13	A14	A15	A16	A17	A18	A19	A20	A21	A22
X0 <sub>2</sub>	X0 <sub>3</sub>	X0 <sub>4</sub>	X0 <sub>1</sub>	X7 <sub>2</sub>	X1 <sub>3</sub>		X18	X4 <sub>7</sub>	X5	X0 <sub>9</sub>	X1 <sub>2</sub>	X7 <sub>1</sub>		X1 <sub>3</sub>	X7 <sub>2</sub>	X18	X7 <sub>2</sub>	X0 <sub>6</sub>	X7 <sub>2</sub>	X2 <sub>3</sub>	X7 <sub>2</sub>	X2 <sub>4</sub>
		X8 <sub>2</sub>	X8 <sub>1</sub>		X8 <sub>3</sub>		X9 <sub>4</sub>			X8 <sub>4</sub>				X8 <sub>5</sub>		X9 <sub>4</sub>		X9 <sub>2</sub>		X9 <sub>5</sub>		X6 <sub>1</sub>
																						X8 <sub>8</sub>

- Step 5: Use rules and conditions discussed above to define syllable boundaries.

Rule 3.1			Rule 13			Rule 10			Rule 12			Rule 13			Rule 13			Rule 13			Rule 13						
ໄ	ຄ	ງ		ກ	າ	ນ		ພ	ັ	່	ດ		ທ	ະ		ນ	າ		ພ	າ		ສ	າ		ລ	າ	ວ
A0	A1	A2		A3	A4	A5		A7	A8	A9	A10		A11	A12		A14	A15		A16	A17		A18	A19		A20	A21	A22

- Step 6: Put the Lao characters back into original array.

### Example 2: ແຂວງຫລວງພະບາງ 14 ກຸມພາ 2005

- Step 1 & 2: Traverse through the input array mark syllable boundary where you counter non-Lao character. Extract only Lao text and put it in new array.

ແຂວງຫລວງພະບາງ | ກຸມພາ |

- Step 3: Look for possible re-ordering of text. For this example there is none.
- Step 4: Mark each character with possible X<sub>n</sub>.

ແ	ຂ	ວ	ງ	ຫ	ລ	ວ	ງ	ພ	ະ	B	ບ	າ	ງ	B	ກ	ູ	ມ	ພ	າ
A0	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	A11	A12	A13	A14	A15	A16	A17	A18	A19
X0 <sub>2</sub>	X	X	X	X1	X	X	X	X	X7 <sub>2</sub>		X	X7 <sub>1</sub>	X		X	X3	X	X	X7 <sub>1</sub>
		X2	X8	X	X2	X2	X8	X9			X8		X8		X8		X8	X9	
		X6 <sub>1</sub>			X9	X6 <sub>1</sub>													
		X8				X8													



- Step 5: Use rules and conditions discussed above to define syllable boundaries.

Rule 2.1				Rule 11				Rule 12				Rule 13				Rule 6				Rule 13			
ຸ	ຂ	ວ	ງ	B	ຫ	ລ	ວ	ງ	B	ພ	ະ	B	ບ	າ	ງ	B	ກ	ຸ	ມ	B	ພ	າ	
A0	A1	A2	A3	B	A4	A5	A6	A7	B	A8	A9	A10	A11	A12	A13	A14	A15	A16	A17	B	A18	A19	

- Step 6: Put the Lao characters back into original array.

### III. Lao Line Breaking Utility

#### Screen Dumps

Given below are some screen shots that might be helpful in facilitating the understanding.

ສູນເຕັກໂນໂລຊີຂໍ້ມູນຂ່າວສານ (IT Center) ຮັບຜິດຊອບ ການຄົ້ນຄ້ວາ, ການພັດທະນາ, ການບໍລິການ ການເຜີຍ ເຕັກໂນໂລຊີຂໍ້ມູນຂ່າວສານ, ເຊິ່ງນັບແຕ່ສ້າງຕັ້ງມາ ທາງສູນແມ່ນໄດ້ປະກອບ ສ່ວນໃນການພັດທະນາ ເຕັກໂນໂລຊີສື່ສານຂໍ້ມູນຂ່າວສານ ຢູ່ ສປປ ລາວ ຢ່າງແຂງແຮງ

Highlighted text represents the problem. Word is unable to break the line properly due to the continuum of words without any space or any other clue which can be used to decide where to break the line. .We have to insert these clues inform of zero width space.

For instance in first-second lines consider ນັບແຕ່ since there is no space for remaining ຕ in the currently line word shifts it to next line blindly without realizing that it was part of ແ. However, if there were a zero width space in between ນັບ and ແຕ່ then it would shift whole syllable to next line. Similar is the case with ການ in line number first where ານ being part of the syllable is dragged to next line.

User can browser for the file and let the program read the input and do the rest.



Program displays a success message and prompts the user to save the modified file with new name or replace the existing after which the output file is displayed automatically.

Output file:

ສູນເຕັກໂນໂລຊີຂໍ້ມູນຂ່າວສານ (IT Center) ຮັບຜິດຊອບ ການຄົ້ນຄ້ວາ, ການພັດທະນາ, ການບໍລິການ ການເຜີຍ ເຕັກໂນໂລຊີຂໍ້ມູນຂ່າວສານ, ເຊິ່ງນັບແຕ່ ສ້າງຕັ້ງມາ ທາງສູນແມ່ນໄດ້ປະກອບ ສ່ວນໃນການພັດທະນາ ເຕັກໂນໂລຊີ ສື່ສານຂໍ້ມູນຂ່າວສານ ຢູ່ ສປປ ລາວ ຢ່າງແຂງແຮງ

As you can see that syllable boundaries are now properly defined as word finds the zero width space after every syllable which provides it help to properly segment the words and break the line properly. If we put a space between ນັບ and ແຕ່ it would not shift only ຕ but complete syllable ແຕ່ to the next line as shown below.

ສູນເຕັກໂນໂລຊີຂໍ້ມູນຂ່າວສານ (IT Center) ຮັບຜິດຊອບ ການຄົ້ນຄ້ວາ, ການພັດທະນາ, ການບໍລິການ ການເຜີຍ ເຕັກໂນໂລຊີຂໍ້ມູນຂ່າວສານ, ເຊິ່ງນັບ ແຕ່ສ້າງຕັ້ງມາ ທາງສູນແມ່ນໄດ້ປະກອບ ສ່ວນໃນການພັດທະນາ ເຕັກໂນໂລຊີ ສື່ສານຂໍ້ມູນຂ່າວສານ ຢູ່ ສປປ ລາວ ຢ່າງແຂງແຮງ

And as you can see there is a corresponding adjustment in line three where entire syllable ການ moves onto next line.