

| | | |
|------------|---------------------|-----------|
| กึน | ข้างออก | ป่า |
| กึ่ | เขน | ป่า |
| กึ๋น | เข็น | ป่า |
| กุน | เข่น | ป่า |
| กูด | เข็ด | ปาน |
| เก็ง | แข็ง | ผิด |
| เกล้า | แข็ง | ฯพณฯ |
| เกลี้ยว | แข็ง | พาณิชย์ |
| แก้ | แข็งขวา | ย่อง |
| เกาะ | แข็งขึ้น | รอง |
| เกี้ยว | แข็งขึ้น | ฤทธิ์ |
| เกี้ยวะ | แขน | ฤษี |
| เกือก | กรรม- | ฤษี |
| แกง | กรรม | ลลิตา |
| แกะ | จุมพล | ภาซา |
| โกน | จูปล | วก |
| โกร๋น | ชาย | ศาล |
| ไกล | เฒ่า | หริภุญชัย |
| ไก่อ | เณร | หฤทัย |
| ไกล | ตลาด | หลง |
| ขึ้น | ทูลเกล้า | แห่ง |
| ขนาบ | ทูลเกล้าฯ | แห่ง |
| ข้าง | ทูลเกล้าทูลกระหม่อม | แหนม |
| ข้างๆ | น้ำ | แหนหวง |
| ข้างกระดาน | น้ำ | แหบ |
| ข้างขึ้น | นี้ | แหม |
| ข้างควาย | บุญหลง | อาน |
| ข้างๆ คูๆ | บุญ-หลง | ฮา |
| ข้างเงิน | ป่า | |

2 Algorithmic Aspect

The above principle, with appropriate character code assignment such as TIS-620 and ISO/IEC 10646, almost allows C standard library function **strcmp()** to collate Thai strings without much more complication, except:

1. Leading vowels (เ- แ- โ- ใ- ไ-), which are written before consonants, must be considered after the initial consonant. Therefore, the rearrangement is needed before comparison.
2. Diacritics and tone marks (่ ้ ๊ ๋ ๊ ๋ ๊ ๋ ๊ ๋) must be ignored in the first pass, and be considered at later pass if the first pass yields equality.

And these are the only two mandatory requirements for Thai string collation algorithms. No syllable structure nor word boundary analysis is required, as Thai lexicons are ordered alphabetically, not phonetically.

2.1 Leading Vowel Rearrangement

To fulfill this requirement, either a preprocessing or collating-element grouping is required. The preprocessor scans the string once and swaps every leading vowel with its succeeding letter. The preprocessed string is then passed to the normal weight calculation process. Another way to manage this is by means of collating-element formation. Every possible pair of leading vowel and consonant is defined as a collating-element, whose weight equals to the weight of the rearranged substring.

Note that the rearrangement of a leading vowel is performed with its immediate succeeding consonant only. No consonant cluster analysis is needed. Instead, doing so would either face ambiguity problem or yield different order from the Royal Institute Dictionary.

For example, if consonant clusters were concerned, "เพลลา" could be rearranged as either "เพลลา" (two-syllable word เพล-ลา) or "พลเลา" (one-syllable word, with consonant cluster "พล").

To illustrate the different ordering caused by consonant cluster analysis, consider this conforming order: (เพล, เพลง, เพลศ). If consonant cluster were analyzed, it would be rearranged to (พลเล, พลเง, พลศ), and would yield a different order: (เพลง, เพล, เพลศ) (if the consonant cluster "พล" were not grouped as a single collating element); or (เพล, เพลศ, เพลง) (if the consonant cluster "พล" were treated as a single collating element).

2.2 The Multiple Levels of Character Weights

The second requirement of diacritics and tone marks treatment implies multiple levels of weights. Tone marks and diacritics must be ignored in the first level, and weigh more than consonants and vowels in the second level.

There are ten Thai decimal digits (๐ ๑ ๒ ๓ ๔ ๕ ๖ ๗ ๘ ๙), each semantically equivalent to Arabic digit 0-9, respectively. Their weights are then equal to their corresponding Arabic digits in the first level, and are different in the second level, to distinguish languages.

When punctuation marks (๑ ๒ ๓ ๔ ๕ ๖ ๗ ๘ ๙) are concerned, another level of weights is required for them. This is correspondent to the fourth level in the Common Template Table. In string ordering, punctuation marks are less significant than any tone marks and diacritics, and must be ignored in all the first three levels.

For example, (ข้างจ, ข้างกบ, ข้างจ คูจ, ข้างจัน) is a valid order in the Royal Institute Dictionary. In the first level, the considered weights are ขาง, ขางกบ, ขางค, ขางจัน respectively.