

Thai spelling reform for data processing

UDOM Warotamasikkkhadit
Ramkhamhaeng University

1. Introduction

Those who know Thai know that it is impossible to tell whether a vowel *a*, *o*, *oo* or null is to be inserted after the first consonant when two consonants stand next to each other. Very often even a native Thai speaker, who knows very well the meaning of the words he reads, makes wrong syllabification cuts in words like ขนมอบกรอบ as ขน-มอ-บ-กรอบ *khǒnmwàbòkròp*, ขน-มอบ-กรอบ *khǒnmwàpkròp*, ขน-มอบ-กร-อบ *khǒnmwàpkròp*, ขน-มอ-บ-กร-อบ *khǒnmwàbòkròp*, ข-นม-อบ-กร-อบ *khànmòpkròp*, or ข-นม-อบ-กรอบ *khànmòpkròp* where only the last reading is acceptable. It, thus, is very difficult for a foreigner or a machine to syllabify the written words in Thai correctly. This paper aims to present the problems of syllabification, romanization, and data processing of Thai words and to propose the methods of correct syllabification by man or machine.

2. Syllabification problems

2.1 It is difficult to tell whether the covert vowel is an *a*, *oo*, or *o* as in สรรค์ *sǎrāsàk* where the first covert vowel is *oo* and the second covert vowel is *a*, อมร *àmmwǎn* where the first covert vowel is *a* and the second covert vowel is *oo*, ตลบ *tàlòp* where the first covert vowel is *a* and the second covert vowel is *o*.

2.2 A ´, a silence marker, usually silences a character it is placed over, but there are a number of Thai words where a silence marker silences up to three characters. It is almost impossible to tell which characters are to be silent, for example: วิทย์ *wít* where only one character, ย, is silenced, โฟล์ *phoo* where two characters -- a consonant ฟ and a vowel อล์ -- are silenced, กาญจน์ *kaan* where two consonants, ก and ญ, are silenced, ลักขณ์ *lák* where two consonants, ล and ข, are silenced, and ลักขมณ *lák* where three consonants, ล, ข, ม and ณ, are silenced.

2.3 It is impossible to tell whether ทร or ด is a syllable final or a cluster of an initial, for example: กรมฯ can be read as *konmánaa* or *krommánaa*, กรมฯ can be read as

konlámóp, *klommaphá*, *kàlámóp*, or *kàlomphá*. Those who know Thai will recognize that only *krommánaa* and *kàlomphá* are acceptable.

2.4 It is difficult to tell whether a character is a syllable final or an initial of the next syllable. การเวก can be read as *kaanwêek* or *kaaráwêek*. Those who know Thai will recognize that the second reading is acceptable and *kaanwêek* does not exist in Thai.

2.5 It is hard to tell whether ว is an initial *w*, a second member of a cluster, a vowel *ua*, or a final *w*. Similarly it is hard to tell whether อ is an initial *ʔ* or a vowel *ɔɔ*. ขวนขวาย can be read as *khàwǒnkhàwǎy*, *khwǒnkhwǎy*, or *khǔankhwǎy*. The correct pronunciation of this word according to the Dictionary of the Royal Institute is *khwǒnkhwǎy*, but there are a number of people who pronounce it as *khǔankhwǎy*. คราวอด can be read as *khraawʔòt* or *khraawɔɔt*, depending on its meaning.

2.6 An initial cluster with ฏ, ฎ, or ฏ presents problems in reading and pronunciation whether ฏ, ฎ, or ฏ is a cluster with the preceding character or is a non-cluster, for example: เปรียญ *pàrian* and เปรียบ *priap*.

2.7 A number of Thai words contain a silent ฏ or ฎ, such as กำสด *kamsòt*, เกียรติ *kiat*, จริง *cin*, ไซ้ *sáy*, ทรวง *suaŋ*, ศีรษะ *sǐisà*, พรหม *phrom*, พรหมณ์ *phraam*.

2.8 A certain number of word finals contain either a silent vowel *Ô*, or *Ø*, such as ขาติ *cháat*, ประวัติ *pràwàt*, สมมติ *sǒmmút*, ธาตุ *thâat*, เหตุ *hèet*, but the two vowels are not pronounced.

2.9 Several words in Thai contain ้ in the middle of the word, such as ชอล์ก *chǒk*, สำน *sǎan*, โหน้ม *ʔoom*.

3. The existing systems of syllabification

3.1 The use of ้ (yamakkan) over the first member of a consonant cluster and a final. This practice is widely adopted in Buddhist texts and it is still used in the *Book of Prayer* for Thai Buddhists. When a consonant-combination appears in Pali texts, a

yamakkan is placed on top of the first member of the combination¹, for example: ท้าทะสะโยทะนา dwādasayōjanā.

A yamakkan or a silent marker, ๅ, is also used in non-religious texts to signify a final. This system can be found in the manuscript of King Borommakot's version of *Chindamani* where we doubt its genuineness and suspect that the manuscript was created in the early Bangkok period by a distinguished scholar with an intention to honor King Borommakot or to use his name as a pretext, for example: คน์ *khon*, คั่น *khan*, คาน์ *khaan*, คิน์ *khin*, คิ๊น *khiin*, คึ๊น *khin*, คึ๊น *khiin*, คุน์ *khun*, คุ๊น *khuun*, เคน์ *kheen*, แคน์ *khææn*, โคน์ *khoon*, คอน์ *khoon*, ควน์ *khuan*, เคียน์ *khian*, เค็อน์ *khiian*, เคอน์ *khæon*, เค็น์ *khæon* คุวน์ *khuan*, คยน์ *khian*.

3.2 The use of ๅ a dot below, to replace ๅ, a yamakkan. It can be easily seen that the yamakkan takes up space above a character and is troublesome in writing. The use of a dot below seems to be neater in a certain aspect. This system is widely adopted in writing the Pali texts with Thai scripts up to present day, for example:

สมณ พุราหมณ เจว	ราชวสิกญาดโย
อมจเจ นาคเร จาปี	สงคณหาติ ยถารห
samaṇe brāhmaṇe ceva	rājavamsikañātayo
amacce nāgare cāpi	sangaṇhāti yathāraham

3.3 Both systems are adopted into the Thai writing system, but they do not completely solve syllabification problems of Thai for romanization and data processing. Problem 2.4 above can be easily solved, but

3.3.1 It is impossible to predict whether บว is to be pronounced *buan* or *bwaawwa* (see 2.1 and 2.5).

3.3.2 It is impossible to tell whether a character with a yamakkan on top, or a dot below is a final or the first member of a cluster (see 2.3 above). พกฺลย can be read as *phwa klaay* or *phōwklaay* depending on its meaning.

3.3.3 If we employ one of the existing systems, it is impossible to tell whether a character with a yamakkan above or a dot below in the following words is a final or

¹ It is believed that the lower part of a yamakkan covers the first member of a cluster and the top part of it covers the second member of a cluster.

the first member of a cluster, such as เกษตรศาสตร์ *kàsèttràsàat*, มาตรฐาน *mâattràthǎan*, and นิทรา *nítthraa*.

4. The proposed systems for Thai spelling reform for data processing

4.1 The use of ˆ to signify a cluster and ˙ to signify a final. If a yamakkan is placed on top of the first member of a cluster (in order to avoid interfering with a tone mark or a vowel written above the line such as ิ, ี, ื, ุ) and a dot below is placed under the final, Thai words will be more facilitated for data processing and pronunciation. This system will face difficulty in placing a dot below a consonant when a final already has a vowel ุ, ู written under the line.

4.1.1 The problems in 2.1 will be somewhat solved but not completely. Extra pronunciation rules must be devised such as:

4.1.1.1 If ก, ข, จ, ฉ, ช, ฅ, ฆ, or ฌ stands next to a single ะ which is followed by another consonant which is not a final, ะ is inserted after it; and if it stands next to a single ะ which is followed by a vowel, ะ is inserted after it, for example: กรกฎาคม *kòorákòt*, กรกฎาคมคม *kàrákkàdaakhom*, ขรณี *khòoránii*, ขรวัส *kháraawâat*, จรกา *còorákaa*, จรจร *càraacoon*, ฌณี *thòoránii*, ฌาฌ *tháraathoon*, นรเชษฐ์ *nooráchêet*, นราธิป *náraathíp*, มรดก *mòorádòk*, มรุ *máru*, วรจักร *wòorácàk*, วรุณ *wárun*, สรัสสี *sǎw rásit*, สรีระ *sàriirá*, หรดี *hòorádii*, หริน *hàrin*, อรุณ *òoránút*, อริน *òàrin*.

4.1.1.2 If ๗ precedes any consonant, except ฌ, which is followed by a vowel or a covert vowel, ะ is inserted after it, for example: บรม *bòorum*, บริบูรณ์ *bòoríbuun*, บดินทร์ *bòodin*, บริวารณ์ *bòoríbuan*².

4.1.2 Problem 2.2 and 2.9 will be solved if character deletion rules are provided (see 2.2 and 3.3.3). It must be noted that the character deletion rules must come earlier in rule ordering.

4.1.3 Problems 2.3, 2.4, 2.5, and 2.6 can be completely solved and 2.7 can be partly solved.

²It must be noted that บวร is pronounced *bòowwòon*, but บริวารณ์ is pronounced *bòoríbuan*. Even บวร appears in the latter part of บริวารณ์, they are pronounced differently.