

Fig. 2. Thai characters and their reading order

## II. THAI WRITING SYSTEM

### A. Characters and Syllable Structure

Thai Unicode characters range from  $\backslash u0E00$  to  $\backslash u0E7F$ . As displayed in Table I, there are 44 consonants in 21 phonetic groups, 18 vowel symbols (making up single and compound vowels), 4 tone marks, and 2 diacritics [2]. The characters are printed in four lines (Fig. 2). Forward characters are printed on the base line, occupying horizontal space. Dead characters are printed above or below the forward characters. According to Thai encoding standard, TIS620, characters are read from left to right. If there are multiple characters in one column, the reading order starts from forward character, dead character on the lower line, dead character on the upper line, and finally dead character on the top line.

A Thai syllable typically consists of a leading consonant, a vowel, a tone, and a final consonant. Because there is no space between syllables, determining the correct boundary of each one, i.e. syllable segmentation, is a complicated task. Indeed, it is one of challenges in linguistic and information retrieval research, such as [3]-[5].

### B. Thai Romanization

Thai proper names are usually derived from Indian origin languages, Pali or Sanskrit. Through a series of conversions from the original language, a name can be written in a few subtly different ways. Misunderstanding or typing errors are therefore very common. To write a name in Roman or English characters, Thai romanization is performed. According to the government proclamation, reported in [6], romanization rules are based on sound transcription. Thai characters are mapped to certain English characters that give the closest sounds (as in Table I). The mapping disregards tone marks, diacritics, and even meaning. Exceptions to these rules were gathered and summarized in [5].

Following the rules, “รุ่ง” and “รุ้ง”, pronounced in different tones, are converted into “Rung”. The names “พวงค์”, “พงศ์”, and “วงค์” are all converted into “Phong” as “พ” and “ก” are members of the same phonetic group while “ก”, “ศ”, and “ค” are silenced by sound killers.

However, sometimes Pali/Sanskrit romanization is applied to retain the meaning in Pali/Sanskrit, despite inexact sound mapping. For example, an international airport “สุวรรณภูมิ” is officially written “Suvarnabhumi” by Sanskrit

romanization, rather than “Suwannaphum” by Thai romanization.

Note that our framework converts Thai-version names into English, not vice versa, because Thai spelling is more diverse than English one (as seen in the aforesaid examples). When a name is romanized, it tends to be close, if not exactly equal, to the English-version name spelled by the person herself. We use Aroonmanakun’s romanization ([5], [7]) in this research. It involves three main tasks:

- 1) Syllable segmentation by predefined syllable patterns. If there are many possible outcomes, the best one is chosen based on trigram probability from a training corpus.
- 2) Syllable pronunciation. Each syllable is attached with all possible pronunciations. Again, the most probable one is chosen based on probability from another corpus.
- 3) Romanization. Each syllable’s sound is then mapped to English characters using the official standard. A hyphen is added to avoid ambiguity at the syllables’ boundaries, so “สอาด” is converted into “Sa-at”. A space is added for groups of syllables that form isolable words or subwords, so “ฉัตรสกล” is converted into “Chat Sakon”.

## III. RECORD MATCHING METHOD AND TOOL

Suppose that there are two data sources  $A$  and  $B$ . Our task is identifying record(s)  $b$  in  $B$  that belong to the same entity as record  $a$  in  $A$ . The set of matched records can be written  $A \times B$  or  $\{(a, b) \mid \forall a \in A, \forall b \in B\}$ . To determine whether  $a$  and  $b$  are matched,  $a$ ’s keys  $\{K_{a1}, K_{a2}, \dots, K_{an}\}$  are compared with the corresponding  $b$ ’s keys  $\{K_{b1}, K_{b2}, \dots, K_{bn}\}$ . This research focuses on approximate matching. The similarity between  $K_{ai}$  and  $K_{bi}$ ,  $1 \leq i \leq n$ , is measured, yielding  $\delta_i$ . Then, a decision function takes  $\{\delta_1, \delta_2, \dots, \delta_n\}$  and produces a record matching score ( $\Delta$ ) for  $a$  and  $b$ .

### A. String Comparators

To measure the similarity ( $\delta$ ) between  $K_a$  and  $K_b$ , both keys are treated as strings. The similarity score is normalized to  $[0,1]$  range. String comparators currently used in our research are as follows:

- 1) *Levenshtein*. This comparator calculates an edit distance based on the minimum number of insertions, deletions, and substitutions of characters to convert one string into another. Each operation incurs a unit cost. The similarity score is a reverse measurement of the edit distance.
- 2) *Monge-Elkan*. Like Levenshtein, it calculates a similarity score based on edit distance, but assigns decreasing costs to successive operations [8].
- 3) *Jaro-Winkler*. This comparator counts characters that appear in the corresponding and nearby positions (not farther than half of the length of the shorter string) of both strings, and the conversion from one string into another. Weights are added according to the characters’ positions because characters at the tail-end of the string are more likely to differ than those at the beginning of the string [9].
- 4) *Recursive comparator*. Strings  $K_a$  and  $K_b$  may consist of tokens or substrings delimited by punctuations, arranged in different orders. Examples are “Harry James Potter”

and “Potter, Harry J.”. The recursive method compares every token in one string with every token in the other, using a distance-based comparator (such as Levenshtein, Monge-Elkan, and Jaro-Winkler in this research), and calculates the total similarity score [10].

Our comparators compare Thai and Thai strings or English and English strings. In case that one of them is Thai and the other is English, the Thai string will be romanized prior to the comparison.

**B. Decision Function**

A decision function combines  $n$  comparison results  $\{\delta_1, \delta_2, \dots, \delta_n\}$  and gives a record matching score. This function can be simple linear regression, expectation-maximization (EM), decision tree, support vector machine, or user-defined rules. Our decision function is rule-based. A user can specify a set of matching rules, as illustrated by Fig. 3 and Fig. 4. Each rule is composed of clauses connected by logical AND operators. For example,

- Rule 1:  $JaroWinkler(K_{a1}, K_{b1}) \geq s_{11}$  AND  $JaroWinkler(K_{a2}, K_{b2}) \geq s_{12}$ .
- Rule 2:  $Levenshtein(K_{a1}, K_{b1}) \geq s_{21}$  AND  $Levenshtein(K_{a3}, K_{b3}) \geq s_{22}$  AND  $JaroWinkler(K_{a4}, K_{b4}) \geq s_{23}$ .

A similarity threshold  $s$  is set for every clause in a rule. A matching score ( $\Delta$ ) according to rule  $R$  is the average of all similarity scores ( $\delta$ 's) in that rule.

To find a record in data set  $B$  that best matches record  $a$ ,

the following is performed:

- 1)  $Candidate\_Set = \emptyset$
- 2) For each record  $b$  in data set  $B$  {
- 3) For each rule  $R$  {
- 4) If every clause in  $R$  is true {
- 5) Calculate matching score  $\Delta$
- 6) Add  $b$  to  $Candidate\_Set$
- 7) Break (i.e. skip remaining rules)
- 8) }
- 9) }
- 10) }
- 11) Choose  $b$  with the highest  $\Delta$  from  $Candidate\_Set$

If we want to find multiple matches for  $a$ , then candidates can be sorted by their matching scores and the first few can be chosen. On the other hand, if no match is found, then decision rules can be adjusted, e.g. by lowering similarity thresholds or changing string comparators.

**C. Data Integration Tool**

Our data integration with privacy protection environment, or Dipper, was developed, initially aiming to match records whose sensitive data are concealed [11] (the concealment of sensitive data is not in focus of this paper). It was written in Java, employing Cohen’s SecondString API [12] for string comparison. We have been extending it to handle bilingual matching keys, assuming that all characters in a key are either Thai or English.

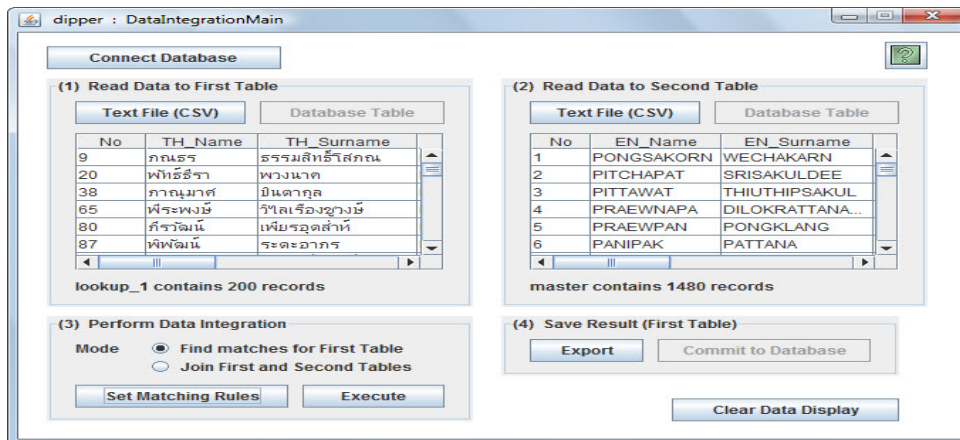


Fig. 3. Data Integrator window

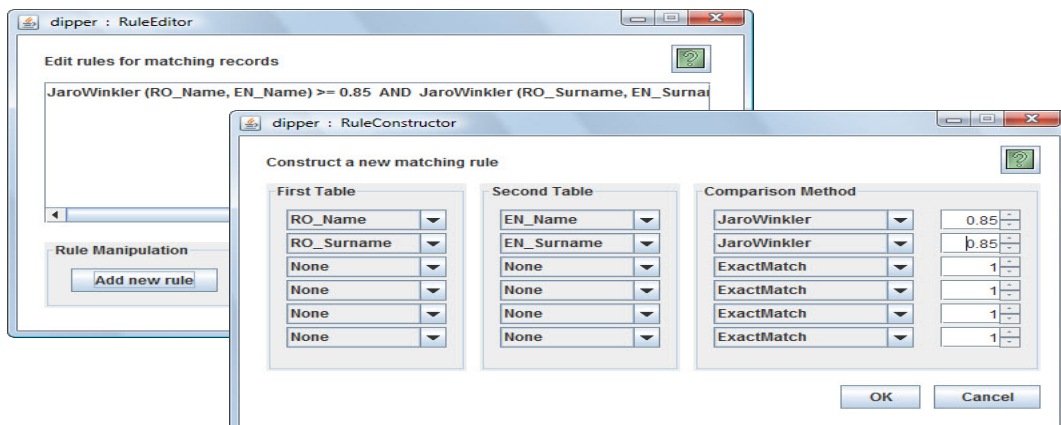


Fig. 4. Rule Editor and Rule Constructor windows

## IV. EXPERIMENTS

## A. Experimental Setup

Our Master data set contained 1,480 records of students in the Faculty of Engineering, Mahidol University. Each record had six attributes: ID, name and surname in Thai characters (TH\_Name and TH\_Surname), name and surname in English characters (EN\_Name and EN\_Surname), and major of study. We generated five Lookup sets, each containing 200 records randomized from Master. Two experiments were conducted. Their setups are summarized in Table II.

We used the same comparator and similarity threshold for comparing names and surnames. Rule setting was similar to that shown in Fig. 4. But the threshold was only 0.1 in order to allow both correct and incorrect matching, for performance evaluation purpose. Our comparators included Jaro-Winkler (JW), Levenshtein (L), and Monge-Elkan (ME).

In the second experiment, we added typographical errors to 50% of records in both Master and Lookup. The number of induced errors ranged from 1 to 4 per record. Each error was randomized from one of the following common errors:

- 1) Repetition, e.g. from “Rung” to “Runng”
- 2) Substitution with neighboring character on the keyboard, e.g. from “Rung” to “Ryng”
- 3) Substitution with look-alike character, e.g. from “Rung” to “Runq”
- 4) Substitution with shifted/unshifted character, e.g. from “รุ่ง” to “ฦๅง” (shift and “ง” yields “ฦๅ”)

The result of an experimental run was a set of 200 records belonging to the Lookup set used in that run. Each record was expanded with attributes from the matched Master’s and their matching score. We counted the number of *matches* (correct matching), *mismatches* (incorrect matching), and *unmatches* (no matching was found).

## B. Results

Fig. 5 shows the average number, across 5 Lookup sets, of correctly matched records in both experiments. Fig. 6 shows the average number of mismatches and unmatches. In the first experiment, all comparators gave more than 90% accuracy. Monge-Elkan was the best one. It was the most optimistic as its matching scores were higher than the others (Fig. 7).

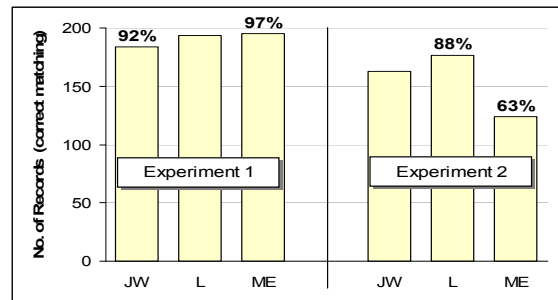
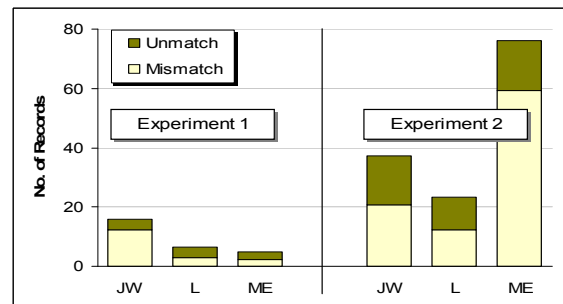
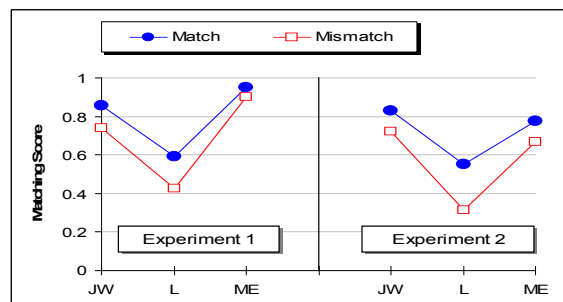
All comparators left nearly identical sets of unmatches in Lookup. It means Dipper could not find any Master’s record that passed the matching rule, despite the similarity threshold being only 0.1. We found that TH\_Name/TH\_Surname in these records had unconventional spelling that did not fit any syllable pattern in the romanization program. As a result, they were segmented and romanized incorrectly. In many cases, whole syllables were missing from RO\_Name/RO\_Surname, making them too much different from their corresponding EN\_Name/EN\_Surname.

With typographical errors in the second experiment, the number of unpronounceable strings increased, leading to even more incorrect romanization and unmatches. On the other hand, adding some errors to TH\_Name/TH\_Surname

did not affect the romanization. For example, “พงศ์” and “ภงค์” (with two errors) were both mapped to “Phong”. There was a little drop in Levenshtein’s and Jaro-Winkler’s performance. In contrast, Monge-Elkan was too optimistic and gave too many mismatches. As seen in Fig. 7, its average matching scores for matches and mismatches were both high and close to each other, compared to those of Levenshtein.

TABLE II: SUMMARY OF EXPERIMENTAL SETUP

	Matching Keys		Compa- rator	Thres- hold
	Master	Lookup		
1	EN_Name	TH_Name	JW	0.1
		romanized to RO_Name		
	EN_Surname	TH_Surname	ME	0.1
		romanized to RO_Surname		
2	EN_Name with typos	TH_Name with typos	JW	0.1
		romanized to RO_Name		
	EN_Surname with typos	TH_Surname with typos	ME	0.1
		romanized to RO_Surname		
Data set summary	Master: 1,480 records Lookup: 200 records × 5 sets			

Fig. 5. Average number of *matches* in both experimentsFig. 6. Average number of *mismatches* and *unmatches* in both experimentsFig. 7. Average matching scores for *matches* and *mismatches* in both experiment



## V. CONCLUSION

Our research aims to tackle approximate record matching, where matching keys are proper names that may be stored in different variations, either in Thai or English characters. We convert names written in Thai characters into English, and compare them with those written in English characters. The conversion is currently via transcription-based romanization. String comparators namely Levenshtein, Monge-Elkan, and Jaro-Winkler, together with user-defined rules, are employed for comparing names and matching records. Our experiments showed that these comparators were effective. But there were still problems with automatic romanization, especially when typographical errors were present and the boundary of each syllable could not be determined.

Thus, making our data integration program recognize and cleanse some errors beforehand will alleviate the problem. One way is to compare input names with those in a training corpus. To handle unconventionally spelled names, which are increasingly popular, up-to-date samples should be included for training. Besides, we will investigate other approaches to syllable segmentation and romanization. An interesting one, for example, was proposed by Chareonpornasawat and Schultz ([4]). It predicted the syllable's boundary based on entropy measures, and built the pronunciation of an unseen syllable based on the closest one in the training corpus. Another work focusing on romanization was proposed by Tangverapong *et al.* [13].

Lastly, we will incorporate other Thai-English mappings. Dictionary-based techniques can be applied for cases where Thai names are translated rather than romanized, e.g. from “ตึกช้าง” to “Elephant Building”, instead of “Tuek Chang”. Another mapping is when Thai names are neither romanized nor translated to English. For example, “กรุงเทพมหานคร” is mapped to “Bangkok”, not “Krungthep Maha Nakhon” by romanization, or “Great Angel City” by translation.

## REFERENCES

- [1] N. Kanchanawan, “Romanization, transliteration, and transcription for the globalization of the Thai language,” *The Journal of the Royal Institute of Thailand*, 31(3): 832-841, 2006.
- [2] T. Karoonboonyanan, “Standardization and implementations of Thai language,” National Electronics and Computer Technology Center (NECTEC), Thailand. Available: <http://www.nectec.or.th/it-standards/thaistd.pdf>.
- [3] C. Wutiwiwatchai and A. Thangthai, “Syllable-based Thai-English machine transliteration,” in *Proceedings of the 2010 Named Entities Workshop, ACL 2010*, Upsala, Sweden, 2010, pp. 66-70.
- [4] P. Chareonpornasawat and T. Schultz, “Example-based grapheme-to-phoneme conversion for Thai,” in *INTERSPEECH 2006 - ICSLP, Ninth International Conference on Spoken Language Processing*, Pittsburgh, PA, USA, 2006, pp. 1268-1271.
- [5] W. Aroonmanakun and W. Rivepiboon, “A unified model of Thai romanization and word segmentation,” in *Proceedings of the 18<sup>th</sup> Pacific Asia Conference on Language, Information and Computation (PACLIC)*, Japan, 2004, pp. 205-214.
- [6] The Royal Institute of Thailand, “Principles of romanization for Thai script by transcription method,” presented at the 8<sup>th</sup> United Nations Conference on the Standardization of Geographical Names, Berlin, Germany, 2002.
- [7] W. Aroonmanakun, Thai romanization [software]. Available: <http://www.arts.chula.ac.th/~ling/tts>.
- [8] A. E. Monge and C. P. Elkan, “The field matching problem: algorithms and applications,” in *Proceedings of the 2<sup>nd</sup> International Conference on Knowledge Discovery and Data Mining (KDD)*, E. Simoudis, *et al.* Eds. AAAI Press, 1996, pp. 267-270.
- [9] W. E. Winkler, “The state of record linkage and current research problems,” Technical Report RR/1999/04, US Bureau of Census, Washington DC, USA, 1999.
- [10] W. W. Cohen, P. Ravikumar, and S. E. Feinberg, “A comparison of string distance metrics for name matching tasks,” in *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, Mexico, 2003, pp. 73-78.
- [11] R. Marukatat, “Dipper: a data integration with privacy protection environment,” in *Proceedings of the International MultiConference of Engineers and Computer Scientists (IMECS)*, Hong Kong, 2009, pp. 750-754.
- [12] Second String Project Page. Available: <http://secondstring.sourceforge.net>.
- [13] A. Tangverapong, A. Suchato, and P. Punyabukkana, “Romanization of Thai proper names based on the popularity of usages,” in *PAKDD 2009, Lecture Notes in Computer Science, vol. 5476*, T. Theeramunkong, *et al.* Eds. Springer, 2009, pp. 580-587.