

# Implementations that Unify the Language Processing

Virach Sornlertlamvanich  
Thai Computational Linguistics Laboratory, NICT  
112 Phahon Yothin Road, Klong Nueng, Kong Luang,  
Pathumthani 12120, Thailand  
Email: virach@tcclab.org

## Abstract

*In this paper, we exhibit the ability of statistical-based approach that can cope with the difficulties in processing a non-segmenting input text. When there is no explicit information about the boundary of a word, parsing an input text is a formidable task. Not only the contemporary word list, but also the word usages have to be maintained to cover the use in current texts. The accuracy and efficiency in higher processing do heavily rely on the word identification process. To avoid the unreliable word segmentation output, we recognize the input text as a consecutive string. Without relying on the word segmentation performance, we directly generate a statistical model of the focused problem from the input string. The word identification problem is then defined as a part of other for performing the unified language processing in total. To exhibit the ability in conducting the unified language processing, we selectively created systems for language identification and multi-lingual search engine, and utilities for term-based ontology alignment and TCL's Computational Lexicon (TCLLEX) development.*

**Key Words:** non-segmenting language, unified language processing, statistical approach, probability, language identification, multi-lingual search engine, ontology alignment, computational lexicon

## 1. Introduction

Disambiguation is our major concern in natural language processing. Though the morphological and syntactic information play an important role in assisting the disambiguation process, degree of the allowing information can vary according to the type of the language. For instance, a space character between words reduces the task in identifying word boundary. Similarly, the grammatical markers, inflection and punctuation marks are quite meaningful information for identifying the role of

each word and structural relations between words in a sentence. Consequently, it is natural to say that many approaches in language processing have been studied in terms of word units. A word unit is definitively a unit of language that native speakers can identify. Based on the classical approaches, without knowing the entity of word, it is not so efficient to develop the language model. For the language that does not allow to use a space character (or any special markers) to separate words in natural text, so-called a non-segmenting language, i.e. Thai, Chinese, Japanese, Korean, etc., a high performance word segmentation algorithm to identify word boundary is crucial. The performance of language processing totally relies on the efficiency and the accuracy of the word segmentation algorithm.

In our recent research, we proposed a language interpretation model to deal with an input text as a byte sequence rather than a sequence of words. It is an approach to unify the language processing model to cope with the ambiguities in word determination problem. The approach takes an input text in the earliest stage of language processing when the exhaustive recognition of total word identity is not necessary.

In this paper, we present a system that can successfully identify the language of an input text based on the study of 20 different languages, a system that can index multi-lingual texts for implementing in a multi-lingual search engine. A utility for aligning two ontologies based on term vectors of conceptual nodes, and a utility for developing TCL's Computational Lexicon (TCLLEX).

## 2. Language Identification Service

Language identification is yet another challenging task when it is going to be conducted without any grammatical knowledge. Byte sequence is the only magic key in our approach to determine the language of the input text. We introduce string kernel for this language identification task. We conducted

experiments using 2 kernelized versions of centroid-based method and support vector machine (SVM). The accuracies of identification are acceptable for both methods. The accuracies reach 95 percent with only 10 training sets (2 KB per set). It is also found that the simple centroid-based classifier is comparable to the SVM classifier based on the string kernel.

String kernel is introduced to compute the common subsequences of two strings. A kernel can be simply thought of as the inner product function between two vectors,  $K(x, y) = \langle \Phi(x), \Phi(y) \rangle$ . With the recent advance in kernel methods, the kernel computation is not just limited to vectorial objects, but can be performed on sequences based on the so-called string kernels [2][12].

To describe how to compute the string kernel based on our proposed efficient method of using fast matching with suffix trees, we consider two strings of  $u = yzxxz$  and  $v = xyzxxxy$ , where the range of considering substring length ( $r$ ) is  $1 \leq r \leq 2$  and the set of characters ( $\Sigma$ ) is  $\Sigma = \{x, y, z\}$ . Figure 1 demonstrates the suffix tree for the string  $v = xyzxxxy$ .

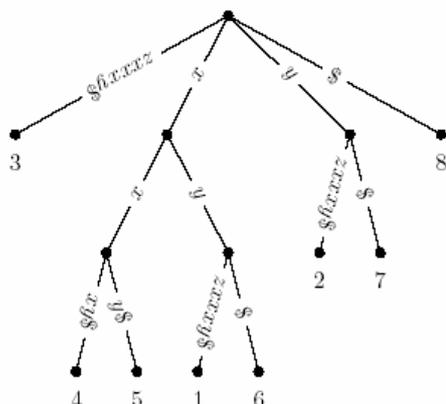


Figure 1: Suffix tree for the string  $v = xyzxxxy$

The following table lists all the matched between  $u$  and  $v$ , where the number of all occurrences of substrings in common is given in parentheses.

$u$	$v[i]$	Score
$y$	$y(2), yz(1)$	$2 \cdot \lambda^{2-1} + 1 \cdot \lambda^{2-2} = 2\lambda^2 + \lambda^4$
$z$	$z(1), zx(1)$	$1 \cdot \lambda^{2-1} + 1 \cdot \lambda^{2-2} = \lambda^2 + \lambda^4$
$x$	$x(4), xx(2)$	$4 \cdot \lambda^{2-1} + 2 \cdot \lambda^{2-2} = 4\lambda^2 + 2\lambda^4$
$x$	$x(4), xz(0)$	$4 \cdot \lambda^{2-1} = 4\lambda^2$
$z$	$z(1)$	$1 \cdot \lambda^{2-1} = \lambda^2$

As a result, we can compute the string kernels with the computational complexity of  $O(c|u|+|v|)$ .

We conducted experiments on 4 groups of language i.e. *Multi4* (Thai, Chinese, Japanese,

Korean), *Multi8-1* (English, French, Italian, Portuguese, Spanish, Swedish, German, Hungarian), *Multi8-2* (Czech, Polish, Croatian, Slovak, Slovenian, Bulgarian, Russian, Greek) and *Multi20* (all the 20 languages). The results are shown in Figure 1-5.

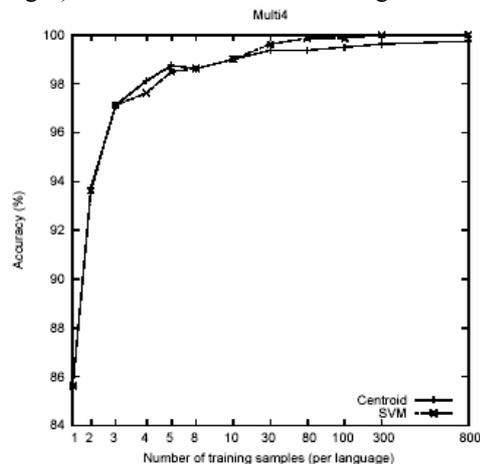


Figure 2: Classification accuracy on *Multi4*

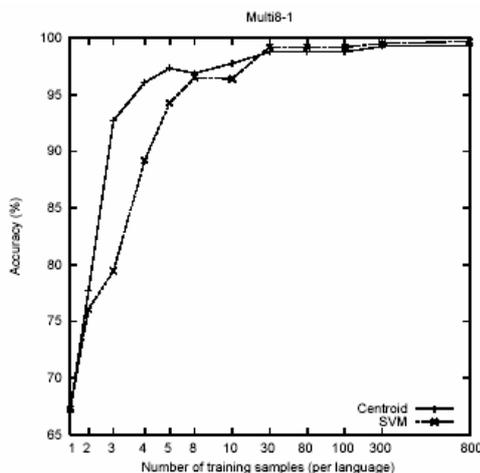


Figure 3: Classification accuracy on *Multi8-1*

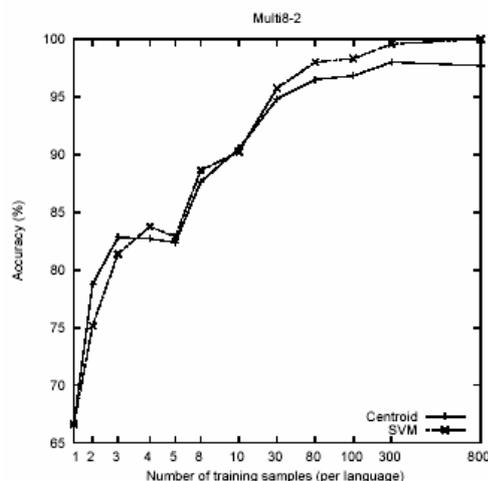


Figure 4: Classification accuracy on Multi8-2

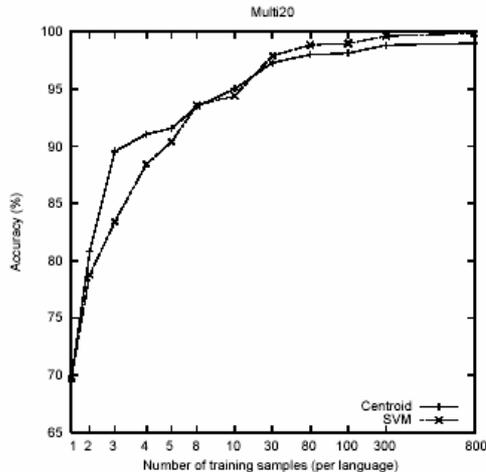


Figure 5: Classification accuracy on Multi20

In total, comparing to centroid-based method, SVM classifier yields slightly better performance when given more training samples, while the centroid-based classifier performs better for small numbers of training samples.

We hypothesize that, for the SVM classifier, only the critical support vectors are retained after the training process, whereas the centroid-based classifier can exploit the feature combination of the representative centroid that is the mean vector of all training samples. However, both had shown significant results in identifying languages of whatever groups of the languages.

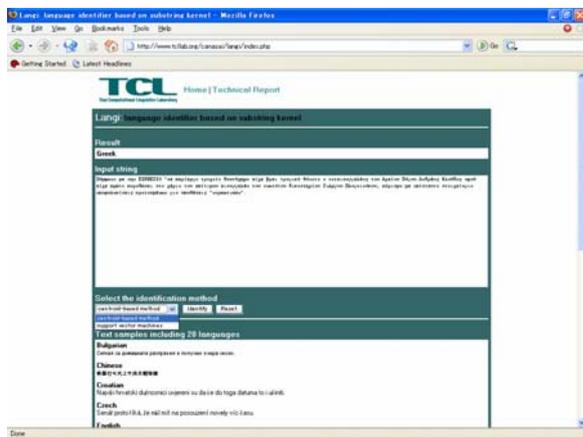


Figure 6: Language Identification Service

Figure 6 shows the interface for language identification service (<http://www.tcllab.org/canasai/langi/>). Both algorithms of the centroid-based and SVM classifiers are selectable to test on the input text. The input text can be as short as a word. However, a longer text, such as a sentence, can help raising the accuracy of the recognition. Users may

select any text from any web page that wanted to know the language, and paste it in the input string box. After selecting the identification method, the user clicks on the identify button to get the language type in the return box of result. The SVM method may consume more time than centroid-based method because of the complexity of the model.

### 3. Multi-lingual Search Engine

The performance of dictionary-based search engines is directly affected by the accuracy of word segmentation algorithms. Our previous work [11] discussed about two possible errors affected by the accuracy of dictionary-based word segmentation modules. Assuming that a dictionary contains 6 words: **a**, **b**, **c**, **ac**, **bc** and **cb**.

#### Case 1: Incorrect word segmentation

The content of the document *A* is **abc bcb**. By using a word segmentation module, the content is separated into **a|bc|bc|b**. Assuming that the correct segmentation is **a|b|cb|cb**. If the query is **cb**, it cannot be found in the document *A* or if the query is **bc**, the document *A* will be incorrectly returned.

#### Case 2: Unregistered word problem

The content of document *A* is **abcdac**. By using the word segmentation module, the content is separated into **a|bc|d|ac**. Assuming that the correct segmentation is **a|b|cd|ac** and **cd** is an unregistered word to the word segmentation. If the query is **bc**, the result from this document will be incorrect or if the query is **cd**, it cannot be found in the document *A*.

The architecture of our dictionary-less search engine is illustrated in Figure 7 comparing to the typical dictionary-based search engine in Figure 8. It is composed of 3 major modules: (1) data indexing, (2) searching and (3) document ranking.

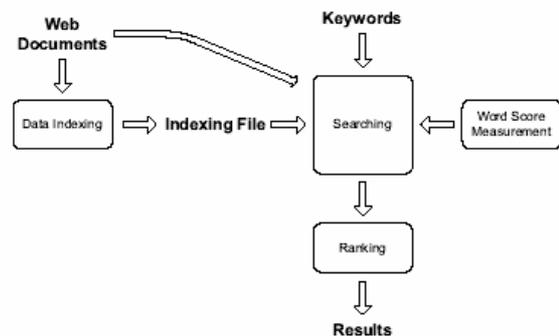


Figure 7: Dictionary-less Search Engine

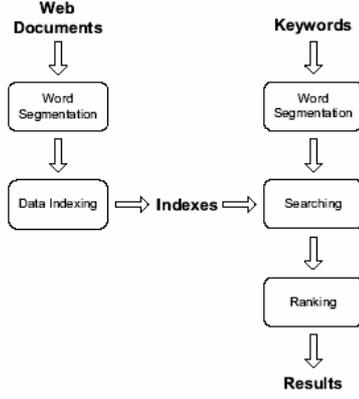


Figure 8: Dictionary-based Search Engine

### (1) Data Indexing

In typical search engines, web documents are separated into words to provide a word list for generating the indexes. In our approach, the data is considered to be the sequence of characters and indexed character by character. We adopt the enhanced suffix array [12] for indexing the data. All suffixes of the data string are indexed. Thus, the number of indexes is equal to the data size. The advantage of this indexing method is that it guarantees all search strings to be found, whereas the word indexing method depends on the word segmentation. This indexing method can also be applied to other languages since it does not require any dictionary and language-specific knowledge.

Based on the enhanced suffix array, it requires  $O(m \log N)$  to access the string in the data, where  $m$  is a length of the search string and  $N$  is the number of indexes. The use of the suffix array guarantees that all search strings will be found. However, only the meaningful strings are preferred. If the found pattern is a part of other word, that pattern is inseparable. As a result, it is not valid as a meaningful word.

For example, assuming that the search query is short and likely to be a part of other strings such as “ยา” (drug), two strings are found i.e. (1) “กินยา” (take a drug) and (2) “พืทยา” (Pattaya, name of a district in Thailand). The first string can be separated into two words: (1) “กิน” (take, eat) and (2) “ยา” (drug). Thus, the word “ยา” in the first string is a meaningful word. For the second string, the first part “พื” is a meaningless string and is strongly connected to the second part “ทยา”. Thus, the word “ยา” in the second string is meaningless since this string is inseparable.

From the example, the validity of a word can be decided from its surrounding context. If the word is strongly connected to other word and inseparable, it is likely to be a meaningless string. In contrast, the word is likely to be a meaningful string if it is loosely connected to other word and separable.

We use mutual information (MI) [1] to measure the degree of the co-occurrence of the query and its context. Let  $xy$  be a query,  $ab$  is the left context and  $cd$  is the right context of the string  $xy$ , the mutual information can be determined by the Equation 1-4.

$$MI_L(abxy) = \frac{p(abxy)}{p(ab)p(xy)} \quad (1)$$

$$MI_L(abxy) \approx \frac{Count(abxy)}{Count(ab)Count(xy)} \quad (2)$$

$$MI_R(xycd) = \frac{p(xycd)}{p(xy)p(cd)} \quad (3)$$

$$MI_R(xycd) \approx \frac{Count(xycd)}{Count(xy)Count(cd)} \quad (4)$$

If the MI value is high,  $xy$  is likely to be a part of the context. On the other hand,  $xy$  should be independent from the context if the MI value is low. We define the inverse of MI as the word score. The word score is calculated by the Equation 5-6.

$$wscore_L(xy|ab) = 1 - norm(MI_L(abxy)) \quad (5)$$

$$wscore_R(xy|cd) = 1 - norm(MI_R(xycd)) \quad (6)$$

The  $norm(\cdot)$  is the normalizing function which normalizes the argument from 0 to 1. At this point, the word score determines the probability of being a word of the string.

### (3) Document Ranking

The word score from the previous step is not only used to determine the word boundary, it is also used to rank the document. That is, the document with higher word score will attain high rank.

We conduct an experiment in order to compare the dictionary-less search engine with the dictionary-based search engine. The experiment is based on the article [3]. We assign 20 queries to 5 volunteers. For each query, the volunteer is shown the top 10 results from each system. Then, each volunteer will choose the documents that are relevant to the query. Each result is considered to be relevant if at least 3 of the 5 volunteers assigned it as relevant for the query. Finally, the satisfaction score of both systems will be calculated for each query. We define the satisfaction score as the ratio of the relevant results to all available results.

The web documents used in the experiment can be divided into two groups. The first group is obtained from websites of newspapers, consisting of 5,853 documents (approximately 65 Mb). The second group contains general articles, not related to news. The second group consists of 7,710 documents (approximately 35 Mb).

The test queries are listed in Table 4. The queries are related to news. Thus, only results from the first group are preferred. The satisfaction score of two systems for each query is presented in Figure 9. From the figure, the mean of the dictionary-less search

engine is higher than that of the dictionary-based approach. Furthermore, there are 10 queries that the dictionary-less approach is better than the dictionary-based approach, while the results of 5 queries are equal and the dictionary-based approach achieves higher results on other 5 queries.

Table 4: Queries

	Unsegmented queries	Segmented queries
1	บริจาค, สีนามิ	บริจาค, สีนามิ
2	เส้นตาย, ซัดคม	เส้นตาย, ซัดคม
3	มันส์, โชว์ตัว, จีน	มันส์, โชว์ ตัว, จีน
4	ผลกระทบ, ราคาน้ำมันแพง	ผล กระทบ, ราคา น้ำมัน แพง
5	ใช้หัวคั่น	ใช้ หัว คั่น
6	ทุจริต, การเลือกตั้ง	ทุจริต, การ เลือกตั้ง
7	จับกุม, ผู้ก่อการร้าย, ภาคใต้	จับกุม, ผู้ ก่อการร้าย, ภาคใต้
8	นโยบาย, แก้ไข, ปัญหาสุขภาพจิต	นโยบาย, แก้ไข, ปัญ หา ยา เสพติด
9	ทดลอง, ลดค่าทางด่วน	ทดลอง, ลด ค่า ทางด่วน
10	ซื้อคืน, สัมปทาน, รถไฟฟ้า	ซื้อ คืน, สัมปทาน, รถไฟฟ้า
11	ลงทุน, โนพม่า	ลงทุน, โน พม่า
12	ส่งเสริม, การท่องเที่ยว, ไทย	ส่งเสริม, การ ท่องเที่ยว, ไทย
13	เลือกตั้ง, ประธานาธิบดี, ปาเลสไตน์	เลือกตั้ง, ประ ธานาธิบดี, ปา เลสไตน์
14	เลือกตั้ง, ผู้ว่า, กทม.	เลือกตั้ง, ผู้ ว่า, กทม .
15	สินค้าไทย, การส่งออก	สินค้า ไทย, การ ส่ง ออก
16	แปรรูปรัฐวิสาหกิจ	แปรรูป รัฐวิสาหกิจ
17	อู๋ม, นายสมชาย	อู๋ม, นาย สม ชาย
18	ฉลองปีใหม่	ฉลอง ปี ใหม่
19	พรทิพย์, ลาออก	พร ทิพย์, ลา ออก
20	สวนสนุก	สวน สนุก

We observe that the dictionary-based search engine faces difficult situation when the query is

excessively segmented, and the segmented words are likely to be general terms. For example, one of the queries is “การส่งออก” (export). The word segmentation module separates this word as “การ | ส่ง | ออก”. All three terms still have some meaning in Thai, but not directly relevant to the compound word. Moreover, these terms are general words and often parts of several words. Thus, the dictionary-based approach tends to return irrelevant documents, but have several locations of these general terms.

Another observation is that the incorrect segmentation does not always affect the search performance of the dictionary-based search engine. For example, a part of one query is “ผลกระทบ” (effect). It is incorrectly segmented into “ผล | กระทบ”. All three terms are meaningless and not general terms. However, the dictionary-based search engine still effectively discovers these terms since some of this terms are quite unique. We also observe that the correctness of word segmentation is less important than the generality of segmented words. When the query is excessively segmented, the dictionary-based search engine still performs well if the segmented terms are not quite general. In contrast, the dictionary-based search engine tends to return irrelevant documents if the query is excessively segmented and the segmented words are general terms. This also explains why the dictionary-based search engine performs better on some queries. Although those queries are sometime incorrectly segmented, the dictionary-based search engine still finds related documents. The reason is that the segmented words are not general. Thus, these words are easily found.

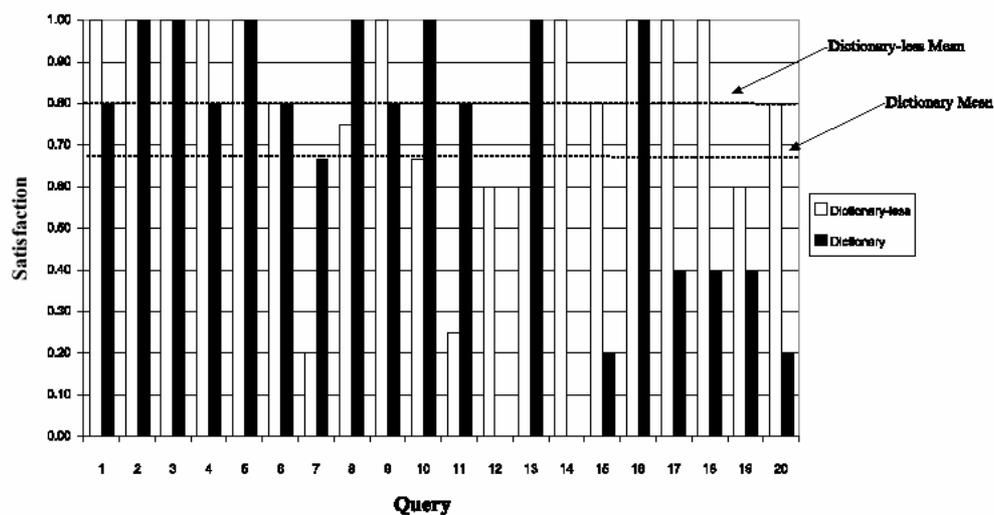


Figure 9: Satisfaction for the test queries

Based on the explained dictionary-less search engine method, we index all the texts by enhanced suffix array. To keep the texts together as a multilingual text, all are converted to UNICODE using UTF-8 character set scheme. The word score of a string in the context is computed when a key word is defined for the search. The snippets are returned with the key word highlighted, ranked according to the word score as shown in Figure 10.



Figure 10: Interface for Multi-lingual Search Engine

Users may search from all the texts or limit the search to only the web pages of a language. The web pages are classified according to the written language by using the previous explained language identification method.

#### 4. Term-Based Ontology Alignment

The shortage of language resources is potentially preventing our research in statistical based approach. A large enough corpus is necessary to capture the language model. It is obvious that the English language is majorly used in any forms, and almost all the other languages have some information related to the English language e.g. bi-lingual texts, bi-lingual dictionaries. The efforts to utilize the advantages of the English language are to increase the knowledge about the target languages. We proposed a term-based ontology alignment [10][7] to increase our language and terminology resources via the English language resources. Words in the classes are used to create a vector model for each class. By computing the extended Jaccard similarity value, we are able to align the 2 concept hierarchies of Thai and English. We are extending our research to use English as a common language for aligning the concept hierarchies that do not have the direct link.

Given two ontologies called the source ontology  $T_s$  and the target ontology  $T_t$ , our objective is to align all the concepts (or semantic classes) between these two ontologies. Each ontology consists of concepts,

denoted by  $C_1, \dots, C_k$ . In general, the concepts and their corresponding relations of each ontology can be significantly different due to the theoretical background used in the construction process. However, for the lexical ontologies such as the MMT semantic hierarchy and the EDR concept dictionary, it is possible that the concepts may contain shared members in terms of English words. Thus, we can match the concepts between two ontologies using the similarity of the shared words.

In order to compute the similarity between two concepts, we must also consider their related child concepts. Given a root concept  $C_i$ , if we flatten the hierarchy starting from  $C_i$ , we obtain a nested cluster, whose largest cluster dominates all sub-clusters. As a result, we can represent the nested cluster with a feature vector  $c_i = (w_1, \dots, w_{|V|})^T$ , where features are the set of unique English words  $V$  extracted from both ontologies, and  $w_j$  is the number of the word  $j$  occurring the nested cluster  $i$ . We note that a word can occur more than once, since it may be placed in several concepts on the lexical ontology according to its sense.

After concepts are represented with the feature vectors, the similarity between any two concepts can be easily computed. A variety of standard similarity measures exists, such as the *Dice coefficient*, the *Jaccard coefficient*, and the cosine similarity [4]. In our work, we require a similarity measure that can reflect the degree of the overlap between two concepts. Thus, the Jaccard coefficient is suitable for our task. Recently, Strehl and Ghosh [8] have proposed a version of the Jaccard coefficient called the *extended Jaccard similarity* that can work with continuous or discrete non-negative features. Let  $\|x_i\|$  be the  $L_2$  norm of a given vector  $x_i$ . The extended Jaccard similarity can be calculated as follows:

$$JaccardSim(x_i, x_j) = \frac{x_i^T x_j}{\|x_i\|^2 + \|x_j\|^2 - x_i^T x_j}$$

We now describe an iterative algorithm for term-based ontology alignment.  $A_s$  mentioned earlier, we formulate that the ontology structure is in the form of the general tree. Our algorithm aligns the concepts on the source ontology  $T_s$  to the concepts on the target ontology  $T_t$  by performing search and comparison in the top-down manner.

Given a concept  $C_i \in T_s$ , the algorithm attempts to find the most appropriate concept  $B^* \in T_t$ , which is located on an arbitrary level of the hierarchy. The algorithm starts by constructing the feature vectors for the current root concept on the level  $l$  and its child concepts on the level  $l + 1$ . It then calculates the similarity scores between a given source concept and candidate target concepts. If the similarity scores of the child concepts are not greater than the root

concept, then the algorithm terminates. Otherwise, it selects a child concept having the maximum score to be the new root concept, and iterates the same searching procedure. Algorithms 1 and 2 outline our ontology alignment process.

**Algorithm 2: FindBestMatched( $c_i$ )**  
**begin**  
 Set the starting level,  $l \leftarrow 0$ ;  
 $BestConcept \leftarrow T_l$  (root concept);  
**repeat**  
 $s_{imp} \leftarrow JaccardSim(c_i, BestConcept)$ ;  
**if**  $T_s^{(l)} \leq T_s^{(max)}$  **then**  
**return**  $BestConcept$ ;  
 Find all child concepts on this level,  
 $\{B_j\}_{j=1}^h \in T_l^{(l)}$ ;  
 Flatten  $\{B_j\}_{j=1}^h$  and build corresponding  
 feature vectors,  $\{b_j\}_{j=1}^h$ ;  
 $s_{j^*} \leftarrow argmax_j JaccardSim(c_i, \{b_j\}_{j=1}^h)$ ;  
**if**  $s_{j^*} > s_{imp}$  **then**  
 $BestConcept \leftarrow B_{j^*}$ ;  
 Set  $l \leftarrow l + 1$ ;  
**until**  $BestConcept$  does not change;  
**return**  $BestConcept$ ;  
**end**

**Algorithm 1: OntologyAlignment**  
**input** : The source ontology  $T_s$  and the target ontology  $T_t$ .  
**output** : The set of the aligned concepts  $A$ .  
**begin**  
 Set the starting level,  $l \leftarrow 0$ ;  
**while**  $T_s^{(l)} \leq T_s^{(max)}$  **do**  
 Find all child concepts on this level,  
 $\{C_i\}_{i=1}^k \in T_s^{(l)}$ ;  
 Flatten  $\{C_i\}_{i=1}^k$  and build their  
 corresponding feature vectors,  $\{c_i\}_{i=1}^k$ ;  
 For each  $c_i$ , find the best matched concepts  
 on  $T_t$ ,  
 $B \leftarrow FindBestMatched(c_i)$ ;  
 $A \leftarrow A \cup \{B, C_i\}$ ;  
 Set  $l \leftarrow l + 1$ ;  
**end**  
**end**

Figure 11 shows a simple example that describes how the algorithm works. It begins with finding the most appropriate concept on  $T_t$  for the root concept  $1 \in T_s$ . By flattening the hierarchy starting from given

concepts (‘1’ on  $T_s$ , and ‘a’, ‘a-b’, ‘a-c’ for  $T_t$ ), we can represent them with the feature vectors and measure their similarities. On the first iteration, the child concept ‘a-c’ obtains the maximum score, so it becomes the new root concept. Since the algorithm cannot find improvement on any child concepts in the second iteration, it stops the loop and the target concept ‘a-c’ is aligned with the source concept ‘1’. The algorithm proceeds with the same steps by finding the most appropriate concepts on  $T_t$  for the concepts ‘1-1’ and ‘1-2’. It finally obtains the resulting concepts ‘a-c-f’ and ‘a-c-g’, respectively.

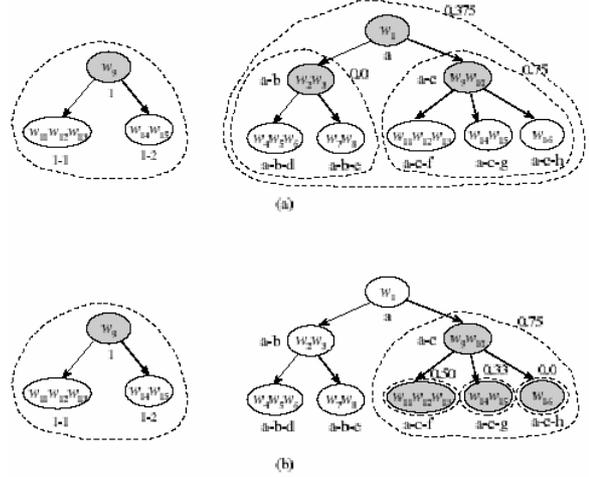


Figure 11: An example of finding the most appropriate concept on  $T_t$  for the root concept  $1 \in T_s$

In our experiments, we used a portion of the MMT semantic hierarchy and the EDR concept dictionary as the source and the target ontologies, respectively. We considered the ‘animal’ concept as the root concepts and extracted its related concepts. In the EDR concept dictionary, however, the relations among concepts are very complex and organized in the form of the semantic network. Thus, we pruned some links to transform the network to a tree structure. Starting from the ‘animal’ concept, there are more than 200 sub-concepts (containing about 7,600 words) in the EDR concept dictionary, and 14 sub-concepts (containing about 400 words) in the MMT semantic hierarchy. It is important to note that these two ontologies are considerably different in terms of the number of concepts and words.

In our experiments, we used a portion of the MMT semantic hierarchy and the EDR concept dictionary as the source and the target ontologies, respectively. We considered the ‘animal’ concept as the root concepts and extracted its related concepts. In the EDR concept dictionary, however, the relations among concepts are very complex and organized in the form of the semantic network. Thus, we pruned

some links to transform the network to a tree structure. Starting from the ‘animal’ concept, there are more than 200 sub-concepts (containing about 7,600 words) in the EDR concept dictionary, and 14 sub-concepts (containing about 400 words) in the MMT semantic hierarchy. It is important to note that these two ontologies are considerably different in terms of the number of concepts and words.

The proposed algorithm is used to find appropriate EDR concepts for each one of 14 MMT concepts. The results are shown in Table 5. In the table, there are 6 relations (marked with the symbol ‘\*’) that are manually classified as exact mapping. This classification is done by inspecting the structures of both ontologies by hand. If the definition of a given MMT concept appears in the EDR concept and the algorithm seems to correctly match the most suitable EDR concept, this mapping will be classified as exact mapping. The remaining 8 MMT concepts, e.g. ‘cold-blood’ and ‘amphibian’, are mapped to closely related EDR concepts, although they are not considered to be exact mapping. The EDR concepts found by our algorithm for these 8 MMT concepts are considered to be only the subset of the source concepts. For example, the ‘amphibian’ concept of the MMT is mapped to the ‘toad’ concept of the EDR.

Table 5: Results of aligned concepts between MMT semantic hierarchy and EDR concept dictionary

MMT concept	EDR concept dictionary
--vertebrate	vertebrate
--warm-blood	mammal
--mammal	mammal
--bird	bird
--cold-blood	reptile
--fish	fish
--amphibian	toad
--reptile	reptile
--snake	snake
--invertebrate	squid
--worm	leech
--insect	hornet
--shellfish	crab
--other sea creature	squid

By analyzing the results, we can classify the MMT words that cannot find any associated EDR words into 4 categories.

1. *Incorrect spelling or wrong grammar*: Some English words in the MMT semantic hierarchy are simply incorrect spelling, or they are written with wrong grammar. For example, one description of a tiger species is written as ‘KIND A TIGER’. Actually, this instance should be ‘KIND OF A

TIGER’. The algorithm can be used to find words that possible have such a problem. Then, the words can be corrected by lexicographers.

2. *Inconsistency*: The English translation of Thai words in the MMT semantic hierarchy was performed by several lexicographers. When dealing with Thai words that do not have exact English words, lexicographers usually enter phrases as descriptions of these words. Since there is no standard of writing the descriptions, these are incompatibility between descriptions that explain the same concept. For example, the following phrases are used to describe fishes that their English names are not known.

- Species of fish
- A kind of fish
- Species of fresh water fish

3. *Thai specific words*: The words that we used in our experiments are animals. Several animals are region specific species. Therefore, they may not have any associated English words. In this case, some words are translated by using short phrases as English descriptions of these Thai words. Another way to translate these words is to use scientific names of species.

The problems mentioned earlier make it more difficult to match concepts by the algorithm. However, we can use the algorithm to identify where the problems occur. Then, we can use these results to improve the MMT ontology.

## 5. TCL’s Computational Lexicon

The TCL’s computational lexicon (TCLLEX) [9] is developed by reusing an existing lexical database for machine translation. It contains 69,060 lexical entries of Thai words. This lexical database was originally constructed for using in the Multilingual Machine Translation (MMT) project, which is a six-year (1987-1992) cooperative project among the group of research institutes led by the National Electronics and Computer Technology Center (NECTEC) of Thailand, and the Center of the International Cooperation for Computerization (CICC) of Japan.

The structure of the lexical entry in the TCL’s computational lexicon consists of three types of information, including morphological, syntactic, and semantic information. The morphological information indicates types of word composition (TYPE). The syntactic information gives grammatical categories (CAT) and subcategories (SUBCAT), and verb patterns in sentence structures (VPPAT). The semantic information provides word concepts (AKO) and case relations (MAPS).

The types of word composition encoded in the morphological information are classified into 2 types;

single word as โทรศัพท์ ‘telephone’ and compound word as โทรศัพท์มือถือ ‘mobile phone’. A single word is a lexical unit composed of either monosyllable or polysyllable and referring to only one conceptual meaning. Contradictorily, a compound word is composed of more than one single word, e.g. the word โทรศัพท์มือถือ ‘mobile phone’ consists of 2 lexical units, โทรศัพท์ ‘telephone’, and มือถือ ‘mobile’.

The syntactic information contains all the information on the syntactic structure of each word, including grammatical categories and subcategories, and the pattern of the verb in a grammatical sentence structure. The TCL’s computational lexicon carries 11 categories with 44 subcategories for narrowing down the functions of words. Each category is divided into subcategories by their position, composition, and reference. For example, the category “determiner” is divided into 9 subcategories by its occurring position, such as after noun, after noun and classifier, between noun and classifier, etc. For the verb pattern, the lexicon classifies the pattern of the verb according to the position of the obligatory arguments and its contextual environment into 11 patterns. For example, the pattern “SUB+V+DOB” indicates that the syntactic structure of this verb consists of the other two components: a NP as a subject (SUB) and a NP as a direct object (DOB).

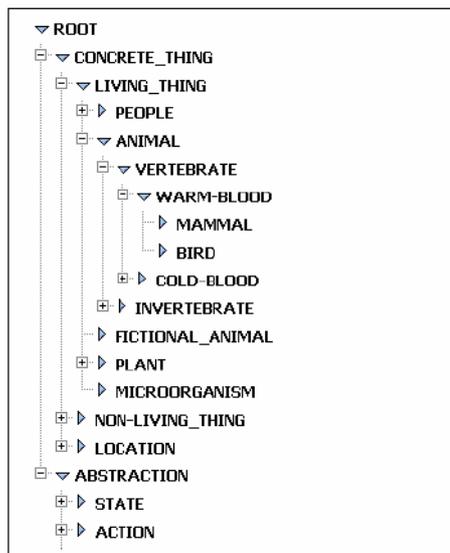


Figure 12: A partition of the semantic hierarchy of TCL’s computational lexicon

For the semantic information, the TCL’s computational lexicon gives basic information about case relations and word concepts. The case relations just bind the thematic roles with syntactic arguments, normally between the main verb and its pre- and post-components. The word concepts are organized on the semantic hierarchy. Each word concept is a group of lexical entries classified and ordered in a

hierarchical level of meanings. The semantic hierarchy is composed of 189 concept classes. Figure 12 illustrates a partition of the semantic hierarchy of the TCL’s computational lexicon.

### Forming Semantic Structure with Logical and Semantic Constraints

Our task is to redefine the structured form of the semantic information in the TCL’s computational lexicon. As addressed previously, the current structure is not expressive to all possible semantic representations. The thoroughness of the representation is required. Consequently, we propose to solve those shortcomings by forming the semantic structure with logical and semantic constraints. The main idea of our design is based on the simplicity and reusability. The logical constraints are capable of dealing with the absence of relatedness of word meanings. The semantic constraints try to discover preferences of syntactic arguments of thematic roles.

Table 6: Logical and semantic constraints

Logical Constraints	
Is-a (ISA)	a conceptual class of a given word
Equal (EQU)	a word that has the same or similar meaning of a given word
Not-equal (NEQ)	a word that has the opposite meaning of a given word
Part-of (POF)	a word that specifies a part of a given word
Whole-of (WOF)	a word that refers to the whole of which a given word is a part
Semantic Constraints	
Agent (AGT)	an entity that initiates the action
Object (OBJ)	an entity that is affected by the action
Instrument (INS)	an entity that is used in the action
Location (LOC)	a position or place where an event occurs
Time (TIM)	a point or period of time when an event occurs

Table 6 shows the entire constraints with their descriptions. One can observe that some logical constraints are similar to semantic relations used to construct WordNet. The EQU, NEQ, POF, and WOF constraints are equivalent to synonyms, antonyms, meronyms and holonyms, respectively. For the ISA constraint, it is identical to AKO of the word. The semantic constraints are the same as the thematic roles, but they are specified with selectional preferences. Here we limit the thematic roles into five important roles, including AGT, OBJ, INS, LOC, and TIM.

It is not necessary that a lexical item must be encoded with all the constraints. We can see that the logical constraints capture information of relations among nouns, while the semantic constraints connect verbs and related nouns. Furthermore, we classify the constraints into two types: obligatory and optional. While the obligatory constraints should be filled as much as possible, the optional constraints can be left empty.

In the logical constraints, there is only one obligatory constraint, ISA. It is also considered to be

the core structure of the TCL's computational lexicon. In the semantic constraints, AGT and OBJ are the obligatory constraints. The remaining constraints are categorized to be the optional constraints. Based on the idea of the frame representation, we finally consider our semantic structure as a template that provides slots of attributes and values to compose the meaning of a given word. For example, the new semantic structure of the verb 'จ่าย' is shown in Table 7.

Table 7: Semantic structure of 'จ่าย' 'pay'

จ่าย 'pay'	
Logical constraints	
ISA	GIVE
EQU	จ่ายเงิน 'spend'
NEQ	รับ 'get'
Semantic constraints	
AGT	PERSON, ORGANIZATION
OBJ	MONETARY

Figure 13 shows an image of TCL's computational lexicon interface for coding the constraints and the hyperlink of the ISA relation to the semantic hierarchy.

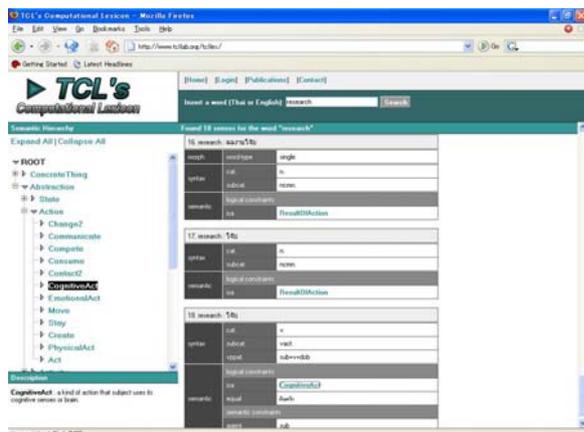


Figure 13: Interface for TCL's Computational Lexicon builder

## 6. Summary and Future Work

The tasks of language identification, word extraction, dictionary-less search engine and term-based ontology alignment had been selected to study by means of the proposed model for unified language processing. The tasks had been evaluated and resulted in a significant performance. Therefore, a non-segmenting language can be efficiently processed without relying on the word boundary information at all. It shows that the footprint of byte sequence of a language is sufficient for processing an input text. It is a direct input with most reliable

information. The advantages in these statistical based approaches are also a fundamental work for unifying multi-lingual tasks where language dependent parts can be lessened.

## 7. References

- [1] Church, K. W., Robert, L., and Mark, L., "A status report on ACL/DCL," In Proceedings of the seventh Annual Conference of the UW Centre New OED and Text Research: Using Corpora, pp. 84-91, 1991.
- [2] Haussler, D., "Convolution kernels on discrete structures," Technical Report UCSC-CRL-99-10, Department of Computer Science, University of California at Santa Cruz, 1999.
- [3] Haveliwala, T., "Topic-sensitive pagerank," In Proceedings of the Eleventh International Conference on World Wide Web, pp. 517-526, 2002.
- [4] Manning, C. D., and Schütze, H., Foundations of statistical natural language processing, MIT Press, Cambridge, MA., 1999.
- [5] Quinlan, J. R., "C4.5 Programs for Machine Learning," Morgan Publishers San Mated, California, 302p., 1993.
- [6] Shannon, C. E., "A Mathematical Theory of Communication," Bell System Technical Journal 27, pp. 379-423, 1948.
- [7] Shisanu Tongchim, Canasai Kruengkrai, Virach Sornlertlamvanich, Prapass Srichaivattana and Hitoshi Isahara, "Analysis of an Iterative Algorithm for Term-Based Ontology Alignment," to appear in Proceedings of The 2<sup>nd</sup> IJCNLP, Jeju Island, Korea, 2005.
- [8] Strehl, A., Ghosh, J., and Mooney, R. J., "Impact of similarity measures on web-page clustering," In Proceedings of AAAI Workshop on AI for Web Search pp. 58-64, 2000.
- [9] Thatsanee Charoenporn, Canasai Kruengkrai, Virach Sornlertlamvanich and Hitoshi Isahara, "Acquiring Semantic Information in the TCL's Computational Lexicon," Proceedings of the Fourth Workshop on Asia Language Resources, IJCNLP-04, Sanya City, Hainan Island, China, pp. 47-53, 25 March 2004.
- [10] Virach Sornlertlamvanich, Canasai Kruengkrai, Shisanu Tongchim, Prapass Srichaivattana and Hitoshi Isahara. "Term-Based Ontology Alignment," Proceedings of the Second International Workshop on UNL, Other Interlinguals and their Applications, Mexico City, Mexico, 2005.
- [11] Virach Sornlertlamvanich, Pongtai Tarsaku, Prapass Srichaivattana, Thatsanee Charoenporn and Hitoshi Isahara, "Dictionary-less Search Engine for the Collaborative Database," Proceedings of the Third International Symposium on Communications and Information Technologies (ISCIT-2003), Songkhla, Thailand, 2003.
- [12] Watkins, C., "Dynamic alignment kernels," Technical Report CSD-TR-98-11, Royal Holloway, University of London, 1999.
- [13] Yamamoto, M. and Church, K. W., "Using Suffix Arrays to Compare Term Frequency and Document Frequency for All Substrings in Corpus," Proceedings of the Sixth Workshop on Very Large Corpora, pp. 27-37, 1998.