

Optimized multilingual query extraction from web service

T.Mathavi Parvathi¹, Dr.Paul Rodrigues²
 Research Scholar¹, Professor,CSE²
 ManonmaniamSundaranar University¹
 King Khalid University²

Abstract -- Web service is the process of finding a suitable web service for a given user query but some times it is difficult for a user to write her request in a language which user could easily understand ,this makes cross language information retrieval and multilingual information retrieval , these methods overcome and eliminate language barriers by allowing information retrieval system to retrieve relevant documents expressed in languages other than the query language .Usually, users retrieve web data by browsing and keyword searching , but these methods have their limitations and disadvantages , today's high tech and automated business world good extraction is also necessary for survival , there have been numerous studies for extracting documents , this paper briefly describes design of a scalable and optimized document extraction for extracting documents and retrieve relevant documents expressed in languages other than the query language **Introduction .**

Index Terms— data mining , information retrieval , data extraction , parallel text collection .

I. INTRODUCTION

Cross language information retrieval (CLIR) allows an information seeker to apply a request in one language and to find information in a different language. Accordingly multilingual information research models must overcome and eliminate language barriers by allowing an Information Retrieval System(IRS) to retrieve relevant documents expressed in languages other than the query language , most of the techniques proposed to solve the problem of cross-language retrieval center around a common idea they attempt to translate the query from the user's language to the language of documents .CLIR is facing the correct terminology equivalent (translation) selection challenge from one language to another. Indeed , the terms of the various languages almost never cover the same semantic field and the sense drift is unavoidable in a translation .To minimize the sense drift and the effort on the translation we propose to avoid the query or the document translation step. CLIR systems are facing the challenge to be effective and to rank the relevant retrieved documents

CLIR is facing the correct terminology equivalent (translation) selection challenge from one language to another. Indeed, the terms of the various languages almost never cover the same semantic field and the sense drift is unavoidable in a translation. To minimize the sense drift and the effort on the translation we propose to avoid the query or the document translation step All manuscripts must be in English. These guidelines include complete descriptions of the fonts, spacing, and related information for producing your proceedings manuscripts. Please follow them.

This template provides authors with most of the formatting specifications needed for preparing electronic versions of their papers. Margins, column widths, line spacing, and type styles are built-in; examples of the type styles are provided throughout this document and are identified in italic type, within parentheses, following the example. PLEASE DO NOT RE-ADJUST THESE MARGINS. Some components, such as multi-leveled equations, graphics, and tables are not prescribed, although the various table text styles are provided. The formatter will need to create these components, incorporating the applicable criteria that follow.

The CLIR systems must use an expressive documents representation , Earlier works showed that the use of simple words as keywords is not always accurate enough to represent the documents contents.Document representation and queries representation using keyphrases , the keyphrases extracted statistically to represent document content may contain noise it may affect IRS performances.

The next approach of extracting any document from web services is translation , In this category the following approachese can be used , query translation (translated the query representation to match the document representations) , document translation (translate the document representations to match the query representation) and interlingual techniques (translate the document and the query representations into a third language or semantic space) , Interlingual techniques approach uses a language independent representations for both queries and documents of s given parallel document collection . query translation and document translation uses dictionary based translation or machine translation . CLIR is an integration of words translation into word based retrieval models .

No translation , in this category , the cognate matching is used between languages having a close linguistic relationship where unchanged words can be expected to match succееfully

we CLIR system the following challenges are faced

1. Translation ambiguity

While translating from source language to target language, more than one translation can be possible. Selecting appropriate translation is a challenge.

2. Phrase identification and translation

Identifying phrases in limited context and translating them as a whole entity rather than individual word translation is difficult.

3. Transliteration errors:

Errors while transliteration might end up fetching the wrong word in target language.

4. Dictionary coverage

For translations using bi-lingual dictionary, the exhaustiveness of the dictionary is an important criteria for performance on system.

5. Font:

Many documents on web are not in Unicode format. These documents need to be converted in Unicode format for further processing and storage.

6. Morphological analysis (different for different languages)

So we proposed CLIR model is based on good performances while choosing scalable optimize content extraction

Parallel Text Collection:

A parallel collection is a text paired with its translation into one or many languages. It is composed of a set of document pairs in different languages that are mutual translations.

Let C be a parallel collection and L the set of languages used in C where n is the number of languages used in C .

Thus, $L = \{L_1; L_2; L_3; L_4 \dots; \dots; L_n\}$ where

$C_{jj} L_{jj} = n$ with $n \leq n$. The collection C is then composed by n sub-collections:

$C = \{C_{L1}; C_{L2}; \dots; C_{Ln}\}$ where L_1, L_2, \dots, L_n . In a parallel collection, each document is translated to all languages in L .

The number of documents k in any sub-collection is the same:

a sub-collection $C_{L1}, C_{L2}, \dots, C_{Lj} = k$ where $k \leq n$.

Our model uses a dynamic structured index to represent multilingual documents and Queries. We evaluated our model on the MuchMore parallel collection. This approach is called collective because it mainly operates on two models of content extraction one based on statistical features and other on Formatting characteristics and collectively applies this model after identifying the document type to yield more accuracy in the extraction.

The methodology of content extraction from html documents [1] extracts the contents for PDA and other device. It takes each web page decompose it, determine the relationship among content and summarize it. The steps for Content Extraction from HTML documents are as follows. Content extraction by tag ratios [3] evaluates number of tags per line on HTML. Tag Ratios (TRs) are the basis by which CETR analyses a webpage in preparation for clustering. The benefits of this approach are that it is a viable and robust content extraction algorithm. It performs well even on non-news bodies and across multiple languages. It achieves better content extraction performance than existing methods works well across varying web domains, languages and styles. The drawbacks are that in some webpages wherein the HTML mark-up is written in a single line CETR would be forced to either return all text or no text. CETR does not perform well on portal home pages. In CETR the recall is high and precision is low. The webpages which do not have advertisements or menus, such as computer science professors' homepages, do not achieve high extraction accuracy. The content code blurring [9] approach aims to locate regions in a document which contain mainly content and little code.

Optimized content Extraction Using Collective Approaches

The model operates on the DOM tree representation of web page to calculate different statistical features of that web page [2], each of them namely Deviation (D), Link Density (L), Normalized Deviation (N) and Normalized Link density (NL).

The Deviation helps to identify amount of text present at each node, it is calculated as given in Eq(1).

$$D(i) = \sigma(i) - \text{Arg}\sigma(T) \quad \text{Eq(1)}$$

The deviation at each node from the arithmetic mean represents how the node contributes towards the information being rendered to the user. The higher the deviation, the more information is rendered through that node. The normalized deviation(N) close to interval [0,1] for each node is estimated as given in Eq (2).

$$N(i) = \frac{D(i) - \text{Min}(D(T))}{\text{Max}(D(T)) - \text{Min}(D(T))} \quad \text{Eq(2)}$$

The Link Density(L) helps for estimating if the node contributes towards traversing or is present there for information, it can be estimated using Eq (3), the link density can be further normalized by using Eq (4).

$$L(I) = \frac{l(i)}{\emptyset(i)} \quad \text{Eq(3)}$$

$$N(i) = \frac{L(i)}{\text{Max}(L(T))} \quad \text{Eq(4)}$$

So we extracted document in very efficient manner without noisy irrelevant documents , at the same time we get the output in expected language when users give input in any language . the query translation is done by parallel tree collection .

II. CONCLUSION

In this paper, we proposed a CLIR model avoiding translation and external resources. Our model uses a structured index to represent multilingual documents and queries. We evaluated our model on the MuchMore parallel collection as well as we focused the optimized query extraction techniques using collaborative approaches. This new collaborative approach yields more accuracy than other previous approaches. It determines the type of content the web page is representing which in turn applies the specific extraction method for each different type of web page. This approach increases accuracy along with handling the different variety of web pages. In the context of a parallel collection, using an English query to retrieve any language documents.

III. REFERENCES

- [1] Sheba Gaikwad , G. Naveen Sundar , “Optimized Content Extraction from web pages using Composite Approaches “ in International Journal of Computer Trends and Technology- volume4Issue3- 2013.
- [2] Chedi Bechikh Ali ,Hatem Haddad , “ Structured Indexing Model for Cross-Language Information Retrieval “ .
- [3] T. Weninger, W.H. Hsu, J. Han, “CETR: content extraction via tag ratios”, Proceedings of the 19th Conference on World Wide Web, WWW '10, ACM, New York, NY, USA, 2010, pp. 971–980.
- [4] Hiemstra D., de Jong F., “Disambiguation Strategies for Cross-Language Information Retrieval”, The Third European Conference on Research and Advanced Technology for Libraries, London, UK, ECCL '99, p. 274-293, 1999.
- [5] Hu R., Chen W., Bai P., Lu Y., Chen Z., Yang Q., “Web Query Translation via Web Log Mining”, The 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Singapore, p. 749-750, 2008.
- [6] S. Chakrabarti, Mining the Web : Discoverin Knowledge from Hypertext Data , Morgan Kaufmann Publishers, 2003.
- [7] C. Mantratzis, M. Orgun, S. Cassidy, “ Separating XHTML content from navigation clutter using DOM-structure block analysis”, in: Proceedings of the Sixteenth ACM Conference on Hypertext and Hypermedia, HYPERTEXT '05, ACM, New York, NY, USA, 2005, pp. 145–147
- [8] S. Gupta, G. Kaiser, D. Neistadt, P. Grimm, “DOM based content extraction of HTML documents”, Proceedings of the 12th International Conference on World Wide Web, WWW '03, ACM, New York, NY, USA, 2003, pp. 207–214.
- [9] T. Gottron, “Content code blurring: A new approach to content extraction”, Proceedings of the 2008 19th International Conference on Database and Expert Systems Application, IEEE Computer Society Press, Washington, DC, USA, 2008, pp.29–33.