

# CLOUD BASED MULTI-LANGUAGE INDEXING USING CROSS LINGUAL INFORMATION RETRIEVAL APPROACHES

Chayapathi A R<sup>1</sup>, G Sunil Kumar<sup>2</sup>, Manjunath Swamy B E<sup>3</sup>, Thriveni J<sup>4</sup>, and Venugopal K R<sup>5</sup>

<sup>1</sup>Computer Science Department, Visvesvaraya Technological University, Vijaya Vittala Institute of Technology Bengaluru, Karnataka, Indi, archayapathi@gmail.com

<sup>2</sup>Computer Science Department, Bangalore University, UVCE Bengaluru, Karnataka, India, gsuneel.k@gmail.com

<sup>3</sup> Computer Science Department, Don Bosco Institute of Technology, Bengaluru, Karnataka, manjube24@gmail.com

<sup>4</sup>Computer Science Department, Bangalore University, UVCE Bengaluru, Karnataka, India,

<sup>5</sup> Computer Science Department, Bangalore University, UVCE Bengaluru, Karnataka, India,.,

**Abstract:** The exponential growth of data sizes created by digital media (video/audio/images), physical simulations, scientific instruments and web authoring joins the new growth of interest in cloud computing. The options for distribution and parallelization of information in clouds make the retrieval and storage processes very complicated, especially when faced with real-time data management. The quantity of Web Users getting access to data over Internet is expanding step by step. An enormous measure of data on Internet is accessible in various languages which could be accessed by anyone whenever. The Information Retrieval (IR) manages finding valuable data from a huge assortment of unorganized, organized and semi-organized information. In the present situation, the variety of data and language boundaries are the difficult challenges for communication and social trade over the world. To tackle such obstructions, CLIR, the cross-language information retrieval frameworks, are these days in solid interest. The Query Expansion (Q.E.) is the way toward adding related and important terms to original inquiry to upgrade its indexing ability to improve the significance of recovered files in CLIR. In this exploration work, Q.E. has been investigated for a Hindi-English and Kannada-English CLIR in that Hindi and Kannada queries are utilized to look through English docs. After the interpretation of query, recovered outcomes are positioned making use of OkapiBM25 to organize the most important doc at the top for expanding the significance of recovered docs using QE. We proposed architecture for Hindi-English and Kannada-English CLIR making use of QE. to improve the importance of recovered reports. In the primary investigation, QE. is performed with and without OkapiBM25 ranking. The outcomes show that the pertinence of recovered archives is higher with OKapiBM25 as contrast with the one without positioning. The work docs plainly demonstrate that the presentation

of Hindi-English and Kannada-English CLIR framework can be improved altogether with query development using fitting terms located at suitable place and the recovered Snippets can incredibly fill in as the continuous test collection.

**Keywords:** Information Retrieval, MLIR, CLIR, Cloud, Query Expansion, OkapiBM25

## 1. Introduction

A high quantity of data on web is accessible in various languages which could be accessed by anyone in any place and time. The capacity to search and recover data in numerous languages is getting progressively significant and challenging at present. Thusly, Multilingual and Cross-lingual data recovery search indexes have got more research consideration and are progressively being utilized to recover data on the web. The area of data access has developed to perform many modern assignments, for example, the data retrieval, multimedia information retrieval, question answering works, summarization, and clustering and Web data recovery. Information Retrieval attempts to distinguish pertinent docs for a data need, communicated as a query. The issues that an IR framework should manage incorporate doc indexing, query investigation and query assessment. Every one of these issues has been the subject of numerous studies in IR.

In the present culture of IR, the user's or organization's process of decision making relies upon the nature of data. The ideal accessibility of pertinent and quality info causes the users to settle on making suitable decision. The importance and nature of data relies upon the user fulfilment or info need. Thus, the applicable, modern, clear, complete, and exact data are the significant boundaries to pass judgment on the pertinence of recovered docs.

The info on Web is constantly expanding because of the enormous utilization of Internet by the people. The internet

users share their data by transferring or downloading of data through Web with various other people. The General IR framework design is illustrated in Fig.1, where the user fires an inquiry via operational module to IR framework. The retrieval framework returns docs by utilizing indexing module that consists of few queries. The reason for indexing module is to introduce the positioned documents before user.

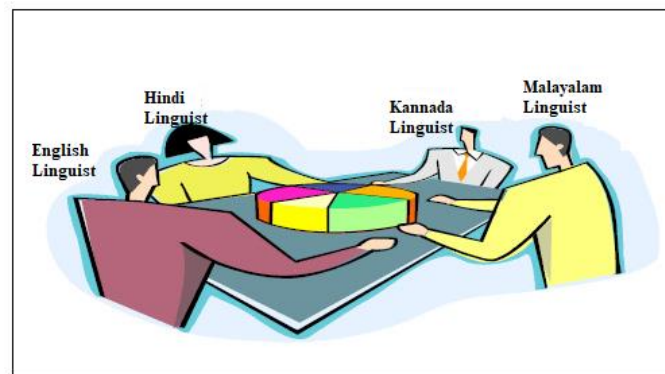


Fig. 2: Illustration of Multilingual Environment

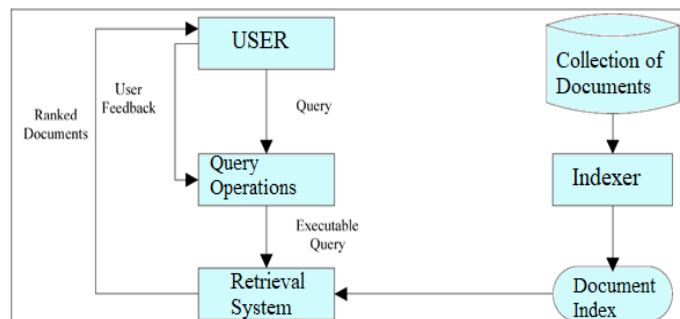


Fig.1: Architecture framework for a General IR

The IR can be broadly classified in to four groups as:

**Monolingual IR:** This is the most essential and most normal type, that like its name uncovers, utilizes just a single language for both the query and info recovery. Monolingual Information retrieval is a recovery process in that the queries presented be the user in a single language to recover info in same language.

**Bi-lingual IR:** This encourages the utilization of queries in one language to get access of the docs in different language.

In the event that two languages are included, for instance if the doc is written in Kannada and the query is done in English, data recovery is supposed to be bilingual IR.

**Multi-Lingual IR:** The MLIR encourages the utilization of queries in one language to get access to docs in various multiple languages. The MLIR gives results that are more thorough than those of mono and CLIR. The most complete way to deal with multilingual site development is to contain copies of the very content in various languages which are connected together so the framework could demonstrate the user an underlying version and afterward the user can pick an alternate interpretation whenever needed.

Multilingual accessing is a multifaceted and complicated topic, grasping specialized, useful and vital issues that have been discussed by the community of information experts. The environment of multilingual is a stage that upholds concurrence of various languages in a solitary information base framework where interaction of data isn't confined by the language. A mechanism to improve the user communication independent of the domain of the language. The Fig. 2 demonstrates environment of multilingual interactions where different individuals come from various linguistic networks and their thoughts are shared.

**Cross-lingual information retrieval:** The CLIR is rapidly turning into a developed field in the data recovery world. The objective is to permit a client to give a query in one language and have that inquiry recover docs in another language. The thought is that the user needs to give a solitary query against a collection of docs which has docs in a multitude of languages. A certain assumption which will be that the user comprehends results acquired in various languages. CLIR is a retrieving process in that the user gives inquiries in a single language to recover data in other language. The reasonable methods are needed to upgrade the performance related to IR, CLIR and MLIR. The CLIR gives a helpful way which can tackle the issues of language limits, where users may submit questions written in their self-language and docs retrieved in different language. For instance, a person is sending queries in source language and gets back the connected data of target language, as appeared in Fig.3.

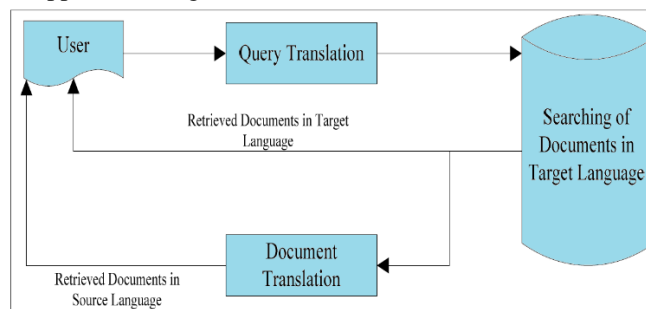


Fig.3: CLIR System Illustration

Presently, a few Indian government workplaces come up short on a robust software for searching through words from the multilingual scanned docs. Physically looking through such records is monotonous and tedious. Besides, there will be countless such docs to be looked for the ideal contents. Hence, there is a squeezing requirement for robust programmed searching app for multilingual old docs, where there is no single vigorous framework existing to perceive the complicated Indian contents. The CLIR empowers the users to recuperate the docs unique in relation to the language of the inquiry. It permits the client to enter their question in one language and recapture the set of records in different

languages. Where, query's language is unique in relation to the language of the report. CLIR framework is a framework wherein a client isn't limited to just a single language, it can plan inquiry in one language and afterward framework restores the records in another language, because in CLIR the query's language and docs both could be translated.

There are different strategies for query translation, document translation or both. There are three essential schemes of translation are parallel corpora, machine translation (MT) and dictionaries. The Query t, translation regularly utilizes either corpus-based interpretation or dictionary-based translation. The translation of document generally just utilizes MT. The Transliteration, Morphological Analyzer, and Word sense disambiguation are the significant pieces of machine translation.

MT is one of the components of language processing inside Computational Linguistic. The MT strategy interprets either the archive or inquiry by utilizing a machine interpretation framework. For deciphering text and words starting with one language then onto the next, bilingual word reference can be utilized. Bilingual word references are utilized in a word reference-based methodology. By looking into terms in a bilingual word reference, queries are interpreted. A few inquiries are likewise deciphered utilizing a few or the entirety of the interpreted terms. The Corpus based interpretation normally gives much better execution, when contrasted with word reference based. The development of parallel corpora is muddled and very costly. It tends to be tremendously unpredictable to discover parallel corpora for specific language or which are sufficiently huge to be useful. Breaking down morphology of given content is called as Morphological Analyzer, which is a part of software app. It senses or produces morphemes of an information word. With the expanding accessibility of on-line data, the people presently approach a large number of data sources and a great many reports. There has been a solid spotlight on extricating significant data from on-line gigantic data information base. To encourage the access of data according to the particular prerequisites of a user has been a most unmistakable part of the exploration in IR.

Lately, the measure of online data from the public authority, business, logical, and private areas have risen significantly. It is as of now engaged that the majority of the Web pages are written in English dialects. In any case, English isn't the local language of almost 50% of all Internet clients and the quantity of non-English talking clients is developing. In India, the greater part of individuals favour Hindi and other regional languages for interactions and different purposes, so the accessibility of non-English client is high. This spurs an interest for research that can build up the connection between English, Hindi and regional languages.

As the Internet is getting mainstream in India, the quantity of the docs written in Indian dialects is additionally expanding

step by step. Numerous news-papers and magazines of the Indian dialects are currently accessible on the web. Yet, there is no effective searching applications exists for Hindi or some other Indian dialects to look through the English reports. The upsides of CLIR are not restricted to singular clients of the Internet. Accordingly, a framework is needed that support clients to find the significant data all the more rapidly and successfully crossing language limits. The prime goal of this research is to investigate Hindi-English and Kannada-English CLIR frameworks.

**Cloud Based Search Engine Strategies:**Cloud computing includes a pay-per-use strategy for availing the services through Internet in a sizeable manner. An important criterion for the clouds is to support data intensive applications. The distributed and dynamic nature of cloud computing systems however makes the process of data management very complex, particularly in the case of real-time data database/processing upgrading. It is very difficult to know how data-intensive applications could obtain adequate performance through very high-level interfaces disclosed in clouds, like those that manipulate today's production grid infrastructures. In addition to this difficulty, there are many other fundamental challenges which needs to be adequately addressed by any data management framework implemented for clouds. An approach to cloud application virtualization that provides additional application and offers higher portability, manageability and consistency of data and applications, has yet to be thoroughly explored. In this context, the idea of universal data access for cloud computing will serve as an alternative to the incorporation of database applications, whereas using a flexible and fast data access system within the cloud can provide a standard access mechanism. The Fig.4 illustrates the diagrammatical representation of working of cloud base search engine.

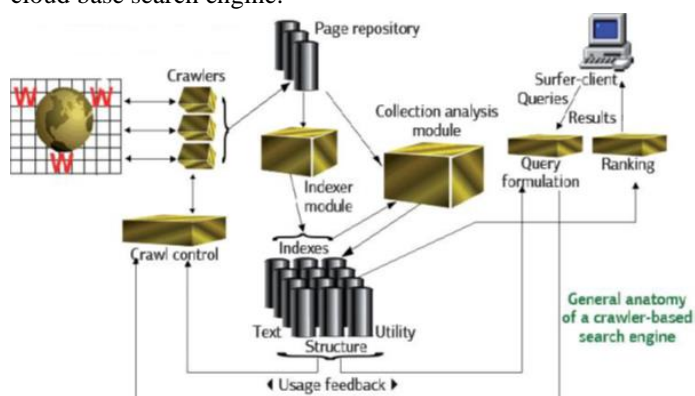


Fig.4: Working Strategy for Cloud based Search Engine

A web-oriented Cloud Service Search engine means that it enables the cloud user to locate the requested data and related information in World Wide Web data base. The primarily involved Cloud Service Searching strategy is Bing searching techniques, MSN search engine, Ask searching, Google's

search engine, Yahoo search engine and so on. The search engine for cloud service is a software program that is specially designed for searching particular documents. Every document must contain specific keyword for searching on to a large data base. The search engine can now have a list of such credentials or data that signify wherever the keyword is located.

A strategy of Cloud Service Search is a web interface which is built from huge collection of the software program. The cloud service searching engine mechanism usually makes use of software application automation. In search engine the robot otherwise spider travels through the index of the WWW data base. The agent followed the different links to the website from one page to another page. In order to scan for specific key words specified by cloud user as an input requirement, the intelligent agent often roams from one website to another website.

Once all the web sites and pages are visited the agent gathered some useful data. Here the crawler is known as an agent. Later on, this smart agent may create a deplorable catalogue page of the WWW, that may help the cloud consumer to understand and choose the appropriate cloud service provider from a huge list gathered. The provider list is prepared on the basis of ranking. The architecture given in Fig.5 illustrates the Cloud based search query analysis.

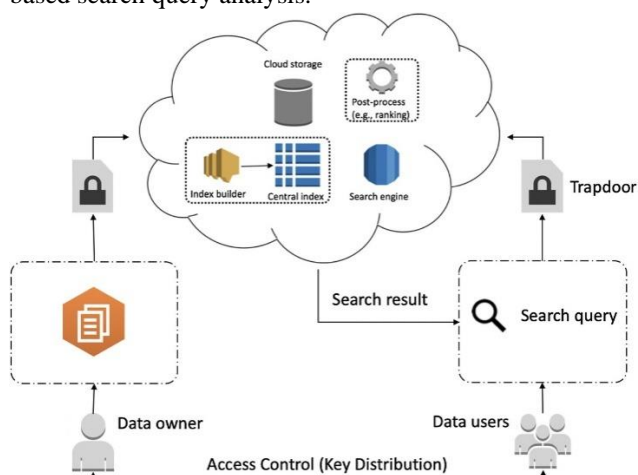


Fig.5: Cloud based search query analysis

The service discovery mechanism or simply a search engine or cloud service search engine. In the system of cloud service discovery, the user needs to show the specification and requirement of the cloud service provider. The Cloud service Search Engine has now begun its work by sending spider agent for roaming. The spider agent holds the capability of retrieve a huge cloud data and document to the extent possible. The crawler agent needs to get the specification in the form of an input keyword. The crawler agent may gather large number of information that will be kept in some data base. Then comes the next software program called as the indexer. The documents stored by crawler must be read by the indexer. The indexer is then started by creating indices on each and every

data and document based on gesture available. The various kinds of search algorithms can be used for any search mechanism. But the required result page for user queries should be shown by the algorithm.

The Intelligent Agent cloud service searching strategy system consists of the following technical module:

**Spider**– This acts as a web browser interface for downloading web Pages that provide cloud service.

**Indexer**– The Indexer makes analysis of cloud service web page downloaded through a program as well as the Cloud pages downloaded by the crawler program & spider program.

**Crawler**– This is a software program that robotically keeps a track of every web link of cloud service provider page.

**Database**– The database is utilized for storing the processed & downloaded web page that provide cloud service.

**Cloud server**– The cloud server is in charge of interacting with the clients as well as the searching technique tool module.

**Results engine**– The result engine is the one that takes the searched outcomes from stored data base.

## 2. Literature Survey

When using a weakly appeared language to search the site, the outcome of the search is very few and less important. The aim of Information retrieval system is providing the necessary information to user's questions. However, the information needed, in some situations may not be available in native language of the user. For instance, when information on "Machine Learning" is searched in Hindi language, the result tends to be very minimal and/or not significant when compared with the outcomes obtained while searching for same information in English language. The other example is, when we search information on "ruling of Hindi speaking people in India" as an English query, the results of the search might be very few when contrasted with results for the similar Hindi translated request. Hence, the prime objective of MLIR (Multilingual Information Retrieval) is to obtain correct information regardless of the language of the query. Example, for few instances which require Multilingual Information Retrieval framework, the user of the information tries to find all relevant information that is available, irrespective of the language, such as a patent lawyer searches for patent information, journalists usually look for the news and updates of other countries, business analysts want to gather information on foreign business undertaken in different countries and users check for documents that contains a combination of multiple languages. Normally, majority of the previous works of MLIR worked towards improving multilingual information retrieval system by improving translation.

The authors in paper [1] showed a study of various methods used in translating document and query and different techniques

of translation like the machine translation, dictionary, or parallel corpora to either translate the documents or the query. Researchers in [2] also concentrated on Cross lingual information retrieval issues including named entities recognizing, dictionary-based problem, Out-Of-Vocabulary, and translation disambiguation. Few others studied the influence of language on search engine by understanding user attitudes and conduct, to infer outcomes for Web practices and standards. The relation between domain/topic and availability of content across languages is studied in this work.

There's been enormous numbers of CLIR and MLIR works in which both of them used distinct methods to enhance the efficiency of MLIR/CLIR system. Few of them focus on translation improvement for enhance Multilingual Information Retrieval system. This type of attention is not shocking, as the primary factor for MLIR/CLIR improvements is the correct interpretation of the query, this directly impacts the accuracy of the results obtained; For instance: The researchers in [3][4] and [5]. Some researchers off late make use of n-gram for the process of retrieval and translation. A unified structure for MoIR (monolingual) and CLIR based on induction from comparable data of complex real-valued word vectors recognized as Word Embedding was proposed by researcher in [6]. The experts in [7] also showed that effectiveness of making use of n-grams character as indexing units, and the conversion units of Cross-Language Information Retrieval systems. In the query translation process and the indexing–retrieval process the benefits of n-grams was used.

Experiments are carried out in CLEF collections for 7 European languages by considering English as target language. The back-translation technique was used for calculating the precision of query translation for Hindi-English CLIR method in work described in [8]. They found out that, the back translation significantly increases the precision of query interpretation of three translators used for the research process (i.e. Babylon Microsoft and Google). However, they stated that Google was the best option for their requirements. The research work in [9] used LSA (Latent Semantic Analysis) method for constructing semantic space by utilizing the parallel corpus of two related languages. The outcomes of Latent Semantic Analysis by utilizing semantic space and without merging documents together gain excellent outcomes than the traditional Latent Semantic Analysis approaches. A new probabilistic method that makes use of a probability-to-possibility conversion for dictionary-based query translation was proposed by the researchers in [10]. This method overcomes the probabilistic success and for few state-of-the-art CLIR instruments.

A method for calculating the amount of online content of a collection of languages at domain level is implemented in works [11]. This calculation is used for constructing a MLIR framework that identifies which languages are highly expressed on the internet about a particular query

subject. The system design includes two modules; the offline module and the online module. The off-line module constructs an index of linguistic diversity for languages at the topic level and, on-line module, on the other hand identifies where the suitable language for search is based the index for extracting the related documents to user query in that language.

Multilingual integrations open up the chance of transmitting information through languages and constructing complicated systems for even languages with less quantity of supervised resources [12][13]. Training a multilingual word embedding model is by far the most popular approach for learning multilingual embeddings, which is then used to obtain representations by structure for files and sentences. Typically, these models are trained solely on sentence or word aligned corpora and composition models are often simple predetermined features such as averages over word embeddings [14] or models of parametric composition learned with word embeddings [15].

Cross-lingual word representations allows us to reason on meaning of words in multilingual contexts and are key facilitators of cross-lingual transition while designing models with natural language processing for low resource languages. A detailed typology of cross-lingual word embedding models was presented in this research. Compare their objective functions and data criteria. The recurrent concept of the research survey is that majority of the models presented in literature strive for similar goals, and that apparently different models are equal most of the time, hyper-parameters, modulo optimization strategies, and so on. The numerous ways in which cross-lingual word embeddings is assessed are also discussed, as well as research horizons and future challenges [16][17].

In [18], The authors worked on Cross-lingual classification of sentiments to instantaneously forecast the sentiment polarity of information in a label-scarce target language by utilizing labelled data through a label-rich language. The researchers have made use of the given multilingual sentiment classification data sets, containing Amazon product feedbacks in 4 languages German, French, English and Japanese of the three categories Music, DVD, Books. The research work in [19], explains Edinburgh's neural machine translation systems for WMT16 shared news translation challenge. The experimental research was performed for four pairs of language in both directions: English-Russian, English-Romanian, English-German and English-Czech.

The authors in [20], using Q.E. technique worked on the personalized Chinese–English CLIR. In this research work, initial query was extended by making additions of relevant words that were extracted from the user's past information of one language. Article [21] gives importance for improving the searching and ranking elements of CLIR framework by making use of PSO (Particle Swarm Optimization) algorithm and Naïve Bayes. For translating English queries into Hindi

language bilingual translator was used, further Hindi queries were expanded with the synonyms of the same.

### 3. CLIR Approach

IR (Information Retrieval) is a much-needed activity in a cross-lingual community as it creates a medium for addressing the problem of language barriers. In Cross-lingual information retrieval, three forms of translation may be used for retrieving the information: document translation, query translation and both document - query translation.

**Query Translation Approach:** This is method of translating each word present in user query of one language into another language. The usefulness of the translation of queries relies on the translational tool that can communicate the needs of the consumer. Translation of queries may be accomplished by machine translation, dictionary and corpus. Query words will be analysed linguistically in dictionary translation and only the keywords are translated by utilizing machine readable dictionaries. Query words are translated on basis of multilingual terms got from comparable documents or parallel collection in the case of corpus translation. In the case of parallel corpus, group of texts are translated into one or many languages. In machine translation, query terms is being translated automatically from one language to another by making use of a context.

**Document Translation Approach:** The Document Translation Approach could be a most ideal scenario in Cross-Lingual IR, if the aim is allowing the users to look for the documents which varies from their language & get the outcomes in the language of the user. In this way, it's really a better choice because it does not require the user to have a passive knowledge of foreign languages. In this approach, all the target languages are being translated to source language. Twofold is the function of this translation. First, post translation or 'on-the-fly translation' or 'as-and-when-needed', where records of any language being searched by the user is being translated into user language at the time of query.

**Dual Translation (Both Document & Query) Approach:** In this dual translation approach – both documents as well as queries are converted into a common representation. This dual approach needs additional storage space for translated files but offers usability when several languages need the same set of documents. Controlled vocabulary systems are one of the examples for such type of approach. Using a pre-defined list of language-independent concepts, these systems reflect all documents and implement queries in same definition space. This approach defines the precision or granularity of possible searching. The key problem of regulated vocabulary frameworks is that, generally non expert users require some amount of training and often need vocabulary interfaces in order to generate successive queries.

Hybrid translation approach also called as dual translation approach may also be performed by pivot language. Due

to limitations of translation facilities, direct translation between two languages may not always be feasible. A third language or a resource named pivot language is needed to perform such type of translation between these languages. Two types of approaches are possible in this process: either document or the query is transferred to pivot language first, and later transferred to the target language; transfer both query and document into pivot language as shown below in Fig.6.

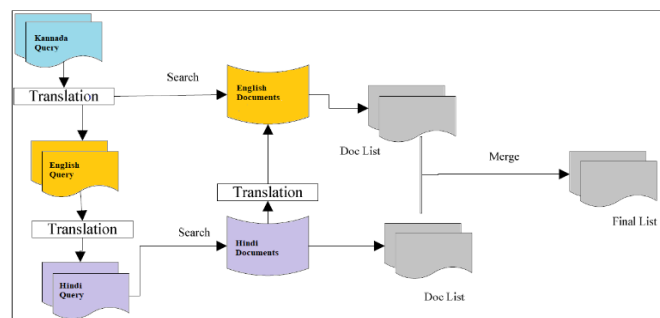


Fig.6: Block Diagram showing Dual Translation using the Pivot Language

In the above block diagram, query of Kannada (i.e. source language) is being fired by the user for searching documents from target language (i.e. Hindi). Surprisingly, due to lack of translator availability, translation process cannot be created between these two pairs of language.

Hence, firstly the query is converted into intermediate or pivot language (i.e. English) which may establish relationship between document and query. Now, pivot language query is again converted into target language for searching the documents. The documents that are retrieved (That is Hindi) are being translated to English language. Further no translation process of the thesis to source language is necessary if user is able to read documents in English. The Table-1 differentiates the three of CLIR.

Table-1: Comparisons of Translation Approaches

Arguments	Query Translation	Document Translation	Dual Translation (Query & Document)
Translation Time	Low	Higher than Query	Higher than both Query & Document
Information Retrieval	Bilingual	Bilingual	Multilingual
Ambiguous nature	Higher	Lower	Higher than both Query & Document
Flexibility	Higher	Lower	Lower
Excess Storage Space	Not Required	Required	Not Required
Nature of Working	At a time provides interface between two languages	At a time provides interface between two languages	At a time provides interface between more than two languages

Over the past couple of years, Cross Lingual IR research has improved and too many frameworks have been developed. Some of the important CLIR systems are explained in Table-2. The importance of CLIR is continuously increasing in each and every field because knowledgeable people are available in almost every languages of World and wants to enhance as well as share their knowledge.

Table-2: Some of the well-known CLIR Tools

Name & Year	Languages	Researcher & Center of development
KANSHIN, 2005	English, Japanese, Korean & Chinese	University of Tokyo Japan by Tomohiro Fukuhara et.al.
UCLIR, 2004	English, Arabic, Korean & Japanese	By Ahmed Abdelali et.al at Computing Research Laboratory at New Mexico State University
MIRACLE, 2003	English,Hindi,German, French, Spanish & Cebuano	By Julio, Sara et.al at Spanish University
MULINEX, 2000	English, German & French	By Capstick et.al at German research Centre for artificial Intelligence
KEIZAI, 1999	English, Korean & Japanese	By William Ogden & James Cowie et.al at Computer Research Lab New Mexico State University, Las Cruces USA
SAPHIRE, 1998	English, Russian, Spanish, German, Portuguese & French	By William R. Hersh, M.D., Laurence C. Donohoe, M.L.I.S at School of Medicine Oregon Health Sciences University Portland, OR, USA

#### 4. Proposed Framework for Kannada-English, Hindi-English CLIR

The proposed framework of our CLIR system as depicted in Fig. 7, can be bifurcated into various phases: Test Collections, ranking of documents, candidate term filtering, query translation, searching of documents, user’s query and Web searching with elongated query. Forevaluating the length of each retrieved document and frequency of each query terms UAM corpus tool is used.

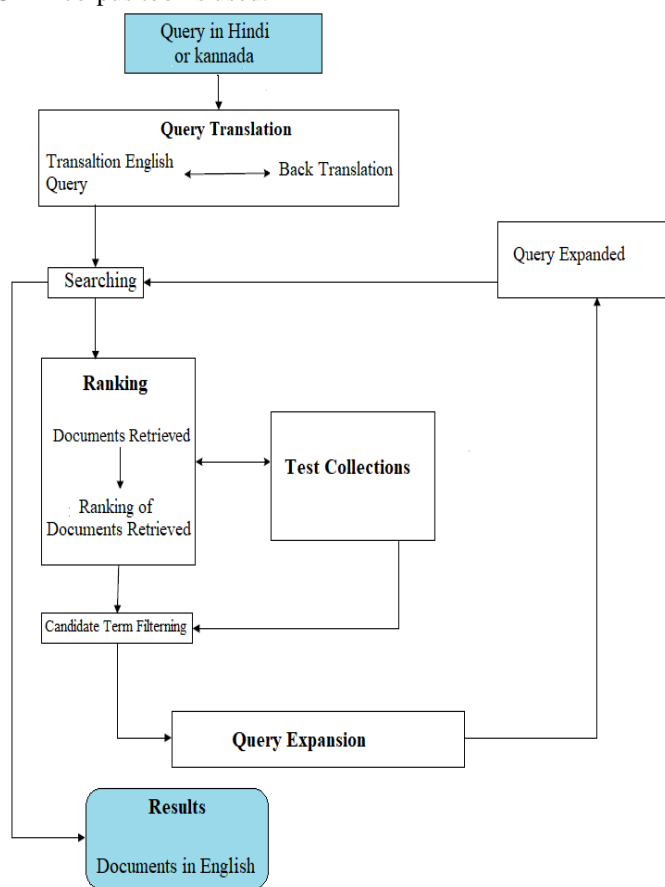


Fig.7: Proposed CLIR Architecture system

The user enter query with the help of the search engine indicating the type of data he wants to obtain. For testing the system, 20 Hindi queries are received from FIRE (Forum for Information Retrieval Evaluation) test collection to retrieve English documents against each question.

The process of query translation includes translation of queries with the help of Google translator (MT tool). Incorrect conversion may lead to poor performance due to lack of efficient communication between documents and query in cross-linguistic IR. Further, for improving the precision of automatic query translation, we used the approach of back translation. Back-translation [22] keeps a check on the accuracy of query translation by converting it to user’s language. It is a very effective technique for improving accuracy of automatic query translation. It also decreases the requirement of manual linguistic assistance for some of the queries where human intervention may be required. After the query translation, Google has been used for fast processing reliability and capability to receive documents against each query.

The retrieved documents that were received from Google search against each query is being re-ranked. The top 10 retrieved documents of each query are being ranked using OkapiBM25 [23]. OkapiBM25 is a mixture of BM11 and BM15, here BM stands for “Best Match”. Ranking done using OkapiBM25 arranged necessary document at the top by making adjustments in the order of retrieved documents which is highly relevant and mines information from top ranked documents for making changes in the user query. For extracting the candidate terms filtering is performed to create various test collections.

For the most part, query of user consists not many catchphrases to communicate required data. In the event that these catchphrases are discovered equivocal, web crawlers may neglect to recover records according to the user requirements. To address the issues of uncertainty in queries, the procedures termed Query Expansion (Q.E) is utilized. The Q.E cycle includes the two significant stages. The main stage includes the looking of most fitting competitor term to extend the first query. The subsequent stage, annexes the proper candidate terms in query at suitable place. In the research works, the extension terms are gotten from different sources by using the documents which are retrieved.

#### 5. Experimentation and Results Analysis

The queries been interpreted utilizing Google interpreter and the precision of interpretation has been checked through the back-translation method. Be that as it may, some queries have been inaccurately deciphered as resulting in back-translation works. For these inquiries just human mediation was required. The records got after Google looking against each inquiry are positioned utilizing OkapiBM25 [24]. Numerous specialists have demonstrated that OkapiBM25 is a powerful technique

and assumes a significant part in Q.E. The reports ranking is done to recognize the most important docs against each inquiry which could be utilized as corpus for Q.E. The OkapiBM25 values are calculated by using the condition:

$$bm25(q,d) = \sum_{t \in q} \log \left( \frac{N - f_t + 0.5}{f_t + 0.5} \right) \times \frac{k_1 + 1}{k + f_{d,t}} f_{d,t}$$

Where: ‘q’ is a terms ‘t’ contained in query; ‘d’ is document; ‘N’ is count of docs in the collection; ‘f<sub>t</sub>’ is quantity of docsconsisting terms ‘t’and ‘f<sub>d,t</sub>’ is the quantity of occurrences of ‘t’ in d; and value of k is got by using the formula:

$$k = k_1 \left( (1 - b) + b \times L_d / A_L \right)$$

In this, constants k<sub>1</sub> and b are set respectively to 1.2 and 0.75; ‘L<sub>d</sub>’ and ‘A<sub>L</sub>’ are doc length and average doc length respectively. The length of everydoc as well as frequency of every query terms are calculated by UAM corpus [25] tool for completely20 queries. The screenshots are displayed in figures Fig.8.a and Fig.8.b are the snapshots of the experimentation done using UAM corpus tool. The length of the docs ‘L<sub>d</sub>’ and ‘A<sub>L</sub>’ are calculatedmaking use of this tool.

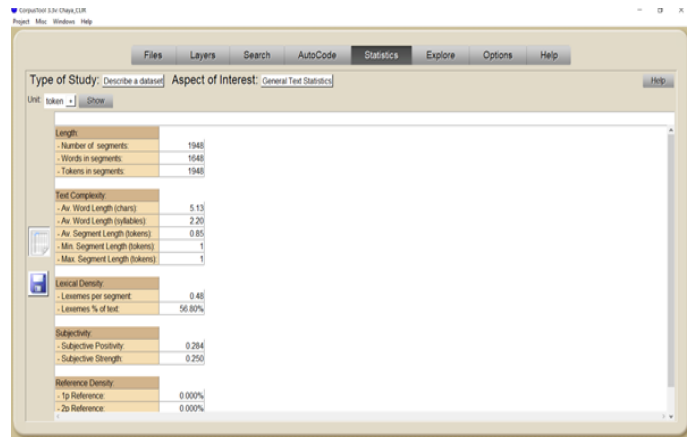


Figure 8a: Snapshots of UAM corpus tool experiment

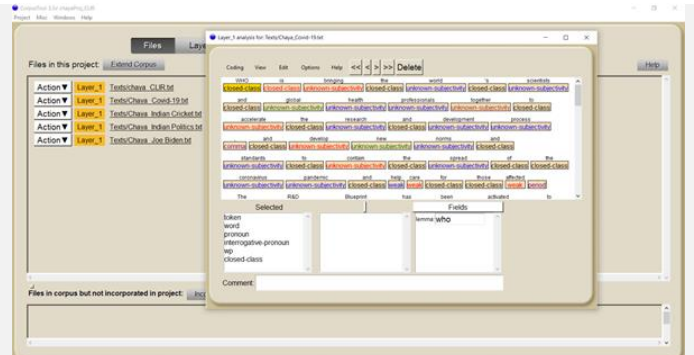
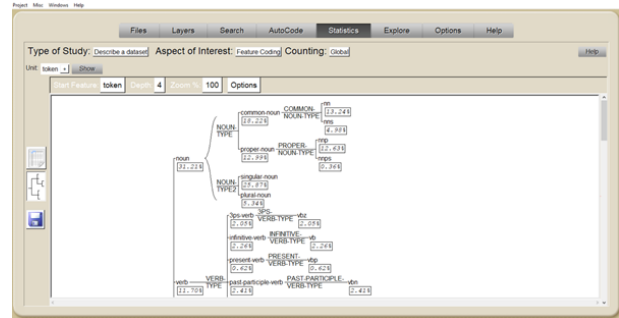
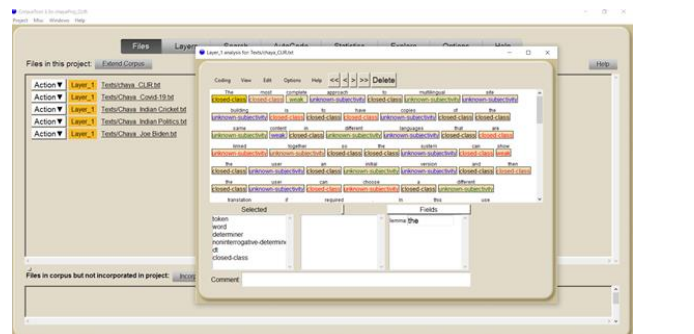


Figure 8b: Snapshots of experimentation with tools of UAM corpus

For the assessment of relevancy, three measures are viewed as, for example, Precision, Average Precision and Mean Average Precision. For computing these, many of important reports are chosen physically from the top @10 recovered docs. Exactness or Precision is the extent of number of significant



recovered docs to the quantity of recovered docs. Precision for each inquiry is figured utilizing the formula:

$$\text{Precision} = \frac{\text{Relevant Retrieved Documents}}{\text{Retrieved Documents}}$$

The Average Precision is the normal of the accuracy value got for the arrangement of top K docs existing after every applicable doc is recovered, and such value is then found the average value of over data needs. The Map: Mean Average Precision for a bunch of queries is the mean of the normal accuracy scores for each inquiry and is processed utilizing condition:

$$\text{MAP} = \frac{\sum_{q=1}^Q \text{AP}(q)}{Q}$$

Where, Q (20 in our case) is the total count of queries.

In the experimentation done by us, the candidate\_terms are acquired from FIRE test collection for every inquiry. The Q.E. is performed with and without ranking of recovered docs If there should be an occurrence of ranking of reports, the OkapiBM25 ranking is thought of. The expansion of query is acted in both the cases and the extended inquiries are recorded. The Table-3 portrays the Precision, AP and MAP values for all cases prior and later the Q.E. The correlation of Precision, AP and MAP estimations of the relative multitude of cases are demonstrated in Fig. 9, Fig.10 and Fig.11 separately.

Table-3: Precision, AP and MAP before and after Q.E.

Queries	Precision value @10 Retrieved Docs			Average Precision		
	Before QE	After QE		Before QE	After QE	
		With OkapiBM25 Ranking	Without OkapiBM25 Ranking		With OkapiBM25 Ranking	Without OkapiBM25 Ranking
Q1	0.4	0.8	0.6	0.355	0.7067	0.4546
Q2	0.5	0.8	0.7	0.413	0.8663	0.5907
Q3	0.3	0.5	0.5	0.1766	0.398	0.398
Q4	0.5	0.6	0.6	0.4264	0.5133	0.5133
Q5	0.6	0.8	0.6	0.4073	0.6817	0.4142
Q6	0.5	0.7	0.6	0.343	0.5546	0.4507
Q7	0.4	0.6	0.7	0.3666	0.5049	0.5541
Q8	0.7	0.8	0.8	0.5541	0.762	0.762
Q9	0.5	0.8	0.7	0.3707	0.7888	0.6241
Q10	0.5	0.7	0.6	0.3216	0.6467	0.5464
Q11	0.4	0.6	0.5	0.3082	0.5049	0.438
Q12	0.4	0.6	0.6	0.2987	0.524	0.524
Q13	0.6	0.7	0.7	0.5133	0.616	0.616
Q14	0.4	0.6	0.6	0.355	0.524	0.524
Q15	0.6	0.7	0.7	0.569	0.6365	0.6365
Q16	0.5	0.6	0.8	0.413	0.3557	0.696
Q17	0.5	0.5	0.5	0.3796	0.3796	0.3796
Q18	0.4	0.6	0.4	0.2932	0.513	0.267
Q19	0.4	0.6	0.6	0.2582	0.4546	0.4546
Q20	0.5	0.5	0.5	0.3057	0.3057	0.3057
<b>Total</b>	<b>9.6</b>	<b>13.1</b>	<b>12.3</b>	<b>7.4282</b>	<b>11.237</b>	<b>10.1495</b>
	<b>Mean Average Precision</b>			<b>0.37141</b>	<b>0.56185</b>	<b>0.507475</b>

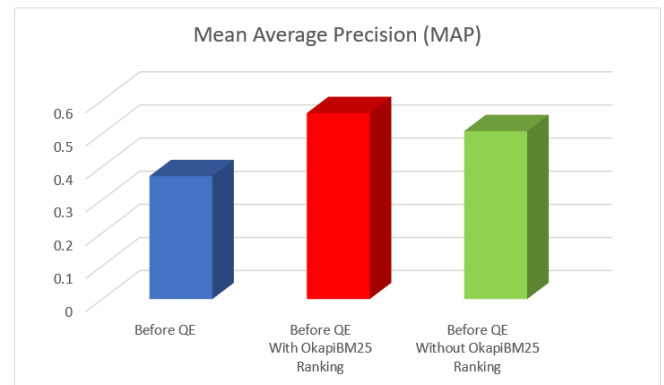


Fig.11: MAP before and after QE (with & without OkapiBM25 ranking)

While making the comparison of the results before and after Q.E., we observed a relevancy improvement without OkapiBM25 and with OkapiBM25. This observation depicts that QE provide better to web searchers. Further, we also compared relevancy (in term of MAP) of Q.E. without & with OkapiBM25. The outcomes of these two observations shows that relevancy of Q.E. with OkapiBM25 is higher than without OkapiBM25. Hence, QE with OkapiBM25 is considered to be one of the best tools for improvement of relevant documents for web searchers in Cross Linguistic IR.

**6. Conclusion**

The complete analysis of the research expresses that expansion of query is a compelling procedure that expands the importance of our Hindi-English and Kannada-English CLIR framework. The ranking assumes a significant function in creation inquiry extension more successful for recovering the

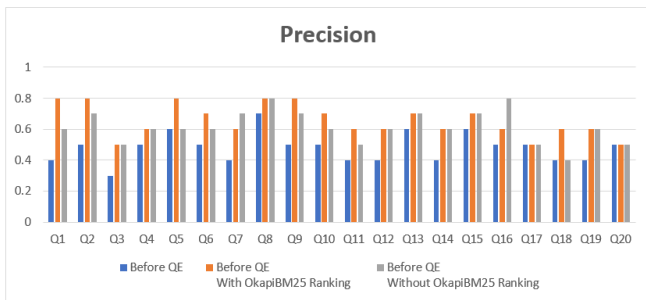


Fig.9: Precision before and after QE (with & without ranking using OkapiBM25)

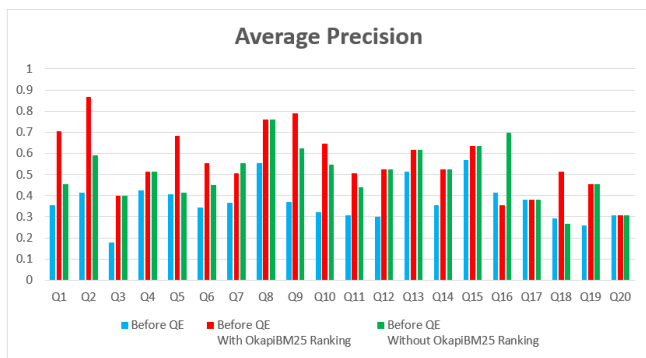


Fig.10: Average Precision (AP) before and after QE (with & without ranking with OkapiBM25)

most applicable docs. At the time of query expansion, the adding of candidate terms in a query at suitable place further expands the recovery pertinence, for which we presented location-dependent calculation. Out of all the test collections, Snippets are discovered to be much esteemed test collection as they give the ongoing assortment of candidate terms for every inquiry to perform powerful query extension. The result of this exploration work opens a stage for creating indexing of search for more languages belong to various regions.

## References

- [1]. Dwivedi SK, Chandra G. A Survey on Cross Language Information Retrieval. *IntJ Cybernet Inform (IJCI)* 2016;5(1):127–42.
- [2]. Sharma M, Morwal S. A survey on cross-language information retrieval. *Int J Adv Res Comput* 2015;4(2):384–7.
- [3]. Kornai A. A new method of language vitality assessment. *Linguistic and Cultural Diversity in Cyberspace* 2015;132.
- [4]. Giang LT, Hung VT, Phap HC. Building structured query in target language for vietnamese-english cross-language IR systems. *Int J Eng Res Technol* 2015;4(04):146–51.
- [5]. Kumar MA, Rajendran S, Soman KP. Cross-lingual preposition disambiguation for machine translation. *Procedia Comput Sci* 2015; 54:291–300.
- [6]. Vulic´ I, Moens M. Monolingual and cross-lingual information retrieval models based on (bilingual) word embeddings. In: *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*. p. 363–72. ACM.
- [7]. Vilares J, Vilares M, Alonso M, Oakes M. On the feasibility of character n-grams pseudo-translation for cross-language information retrieval tasks. *Comput Speech Lang* 2016:136–64.
- [8]. Chandra G, Dwivedi SK. Assessing query translation quality using backtranslation in hindi-english CLIR. *Int J Intelligent Syst Appl* 2017;9(3):51–9.
- [9]. Layfield C, Ivanovic D, Azzopardi J. Multi-Lingual LSA with Serbian and Croatian. *An Investigative Case Study*; 2017.
- [10]. Elayeb B, Romdhane WB, Saoud NBB. Towards a new possibilistic query translation tool for cross-language information retrieval. *Multimedia Tools Appl*. 2017:1–43.
- [11]. Mohamed, Ebtsam & Elmougy, Samir & Aref, Mostafa. (2019). Toward multi-lingual information retrieval system based on internet linguistic diversity measurement. *Ain Shams Engineering Journal*. 10. 10.1016/j.asej.2018.11.009.
- [12]. Waleed Ammar, George Mulcaire, Miguel Ballesteros, Chris Dyer, and Noah A Smith. 2016. Many languages, one parser. *arXiv preprint arXiv:1602.01595*.
- [13]. Melvin Johnson, Mike Schuster, Quoc V Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Vi´egas, Martin Wattenberg, Greg Corrado, et al. 2016. Google’s multilingual neural machine translation system: Enabling zero-shot translation.
- [14]. Aditya Mogadala and Achim Rettinger. 2016. Bilingual word embeddings from parallel and nonparallel corpora for cross-language text classification. In *Proceedings of NAACL-HLT*, pages 692–702.
- [15]. Holger Schwenk, Ke Tran, Orhan Firat, and Matthijs Douze. 2017. Learning joint multilingual sentence representations with neural machine translation.
- [16]. Sebastian Ruder, A Survey of Cross-lingual Word Embedding Models, *Journal of Artificial Intelligence Research* 65 (2019) 569-631
- [17]. P. Paranjape, N. Funde, M. Thakur, M. Dhabu and P. Deshpande, "A Robust and Automated Approach for Multilingual Indian Document Indexing," 2019 5th International Conference on Advanced Computing & Communication Systems (ICACCS), Coimbatore, India, 2019, pp. 457-462.
- [18]. Zhou, Guangyou, et al. "Transfer learning for cross-lingual sentiment classification with weakly shared deep neural network," *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*. ACM, 2016.
- [19]. Sennrich, Rico, Barry Haddow, and Alexandra Birch. "Edinburgh neural machine translation systems for wmt 16", *arXiv preprint arXiv:1606.02891* (2016).
- [20]. Zhou, Dong, et al. "Query expansion for personalized cross-language information retrieval", *Semantic and Social Media Adaptation and Personalization (SMAP)*, 2015 10th International Workshop on. IEEE, 2015.
- [21]. Katta, Eva, and Anuja Arora. "An Improved Approach to English-Hindi Based Cross Language Information Retrieval system", *Contemporary Computing (IC3)*, 2015 Eighth International Conference on. IEEE, 2015.
- [22]. Imran, H., & Sharan, A. "Thesaurus and Query Expansion", *International journal of computer*

- science & information Technology (IJCSIT), 1(2), 89-97, 2009.
- [23]. Sari, Syandra, and MimaAdriani. "Learning to rank for determining relevant document inIndonesian-English cross language information retrieval using BM25", AdvancedComputer Science and Information Systems (ICAC SIS), 2014 International Conferenceon. IEEE, 2014.
- [24]. Ermakova, Liana, and Josiane Mothe. "Document re-ranking based on topic-commentstructure", Research Challenges in Information Science (RCIS), 2016 IEEE TenthInternational Conference on. IEEE, 2016.
- [25]. <http://www.corpustool.com>