# Cross Lingual Information Retrieval

**Special Emphasis for Sibling Language Hindi-Gujarati**

*Author*
**Dr. Kalyani A. Patel**

# PREFACE

India is one of the richest countries, with eighteen constitutional languages which are written in ten different scripts in the field of linguistics. Along with this vernacular linguistic richness, about a billion people use these languages as their first language. A huge amount of regional news and cultural information is usually found on the web in these languages and inaccessible to people of other regions within the country, essentially because of a language barrier. This in turn has sharpened the focus on the demand to cross language barrier at a rapid pace, with no end in sight for this growth in appetite.

Indian languages are highly inflectional with a rich morphology, relatively free word order, and default sentence structure as SOV (Subject Object Verb). Many of them are structurally similar which are better known as sibling languages. The study on cross lingual information extraction and retrieval for similar structured languages with special emphasis to Hindi and Gujarati has been undertaken by us keeping in view the following points.

- Gujarati is still lesser touched language. With a Gujarat Government's emphasis on the basics of e-governance and the development of the Information Technology sector, Gujarati language is slated to attain the significance it aspires towards Digital Era. Therefore, it is required to overcome language barrier from Gujarati to any other Indian language or English language and vice versa.

- The structurally similar language Hindi is considered as an Interlingua to cross language barrier for Gujarati because

    o Hindi is a National language and the transformation applications form Hindi language to Indian languages, Hindi language to English language and vice versa has been developed and tested.

o    Machine translation between closely related languages is easy and efficient.

The principal aim of this research study is to boost up Gujarati on Internet via Hindi as an Interlingua, considering fascinating factor of closely related languages. It is an encouragement of cross lingual information retrieval systems in true sense. The performance of cross lingual information extraction and retrieval system heavily depends on how the query is processed. The study focuses on disambiguated query expansion and translation tasks of query processing. The major contribution of the book can be outlined as follows:

- Query Expansion using Hindi WordNet

- Word Sense Disambiguation using Lesk like algorithm [Lesk,1996].

- Word-for-Word Translation

    o    Development of a generalized model for storage of typologically different words in Hindi and Gujarati.

    o    Design and development of rule-base for

        §    String-pattern-based transformations of morphologically rich and lexically similar words in Hindi and Gujarati.

        §    Mapping of lexical, grammatical and structural divergent words between Hindi and Gujarati languages.

- Development of a system called GH-MAP for rule based token mapping for translation from Hindi to Gujarati and vice-versa.

- Development of search engine prototype to retrieve documents in Hindi-Gujarati.

- Development of evaluation software to calculate PER, METEOR and BLEU score considering into-Hindi / into-Gujarati characteristics.

- A set of machine judgment like F-score, PER, BLEU and METEOR score, has been calculated and compared to establish the relevance of model.

The book is organized in eleven chapters.

Chapter 1 provides research prologue which includes motivation, aim, objectives, and benefits of research. It also includes description of the work already done. Chapter 2 provides overview of information processing, to appreciate the resemblance among information extraction, information retrieval, cross lingual information retrieval, query processing and multilingual information retrieval. Chapter 3 describes possible approaches of machine translation along with the related work done by various researchers. Chapter 4 includes comparative study of Hindi and Gujarati language to understand the similarities and differences between the both. It provides an inspiration to create effective mapping rules amongst languages within the same group with fewer efforts to overcome the language barrier. Chapter 5 provides overview of Hindi and Gujarati Cross Lingual Information Extraction and Retrieval System along with the related work done by various researchers. Chapter 6 describes Query Expansion

Engine along with algorithm and output. Chapter 7 describes Query Translation system along with related work. Provide overview of GH-MAP system. Chapter 8 describes rule base and data structure used for GH-MAP. Chapter 9 describes translation mechanism used by GH-MAP. Also presents algorithm and output of GH-MAP system. Chapter 10 explains evaluation and discusses result. Chapter 11 concludes the contribution of research work in this book. Future work possibilities are discussed in this chapter. Appendix A contains WX alphabet notation used for transliteration Hindi and Gujarati characters, for readers' convenience. Appendix B contains GH-MAP: User Interface screen shots for addition of new rules. Appendix C contains Example of a block of 50 Sentences. Appendix D contains screen shot for BLEU, METEOR and PER scores; exemplary sentences are given in appendix C.

**— *Kalyani Patel***

# ACKNOWLEDGEMENT

# CONTENTS