# Data Extraction and Alignment using Natural Language Processing

M.S Vinu[1],L.Dhanam[2], Dr.Sudha Mohanram[3]

[1]*Professor, Department of CSE,* [2] *M.E Computer Science and Engineering,* [3]*Principal,*
*Sri Eshwar College of Engineering, Coimbatore, India*

*Abstract*— **Web databases generate question result pages based on a user's question. Mechanically extracting the info from these question result pages is extremely necessary for several applications, like information integration, which require to join forces with multiple internet databases. a unique information extraction and alignment methodology known as Combining Tag and value Similarity (CTVS) that mixes each tag and worth similarity. CTVS mechanically extracts information from question result pages by initial distinctive and segmenting the question Result Records (QRR) within the query result pages then aligning the metameric QRRs into a table, within which the info values from constant attribute area unit place into constant column. a replacement techniques to handle the case once the QRRs aren\'t contiguous, which can be due to the presence of auxiliary data, like a comment, recommendation or promotional material, and for handling any nested-structure that may exist within the QRRs. a replacement record alignment algorithmic program has been designed that aligns the attributes during a record, initial pair-wise then holistically, by combining the tag and information worth similarity data.**

*Keywords*— **Database, Information. Data Transfer,internet.**

## INTRODUCTION

Online databases known as internet databases comprise the deep web. Compared with sites within the surface internet, which might be accessed by a novel uniform resource locator pages within the deep internet area unit dynamically generated in response to a user question submitted through the question interface of internet information. Upon receiving the user's question, an internet information returns the relevant knowledge, either structured or semi structured, capsulate in a very markup language pages. Many internet applications, like meta-querying, knowledge integration and comparison looking, want the information from multiple internet databases. For these applications to more utilize the information embedded in markup language pages, automatic knowledge extraction is critical. This paper focuses on the matter of mechanically extracting knowledge records that area unit encoded within the question result pages generated by internet databases. In general, a question result page contains not solely the particular knowledge, however additionally alternative info, like steering panels, advertisements, comments, info concerning hosting sites then on. The goal of internet information knowledge extraction is to get rid of any impertinent info from the question result page, extract the question result records from the page and align the extracted QRRs into a table such the information values

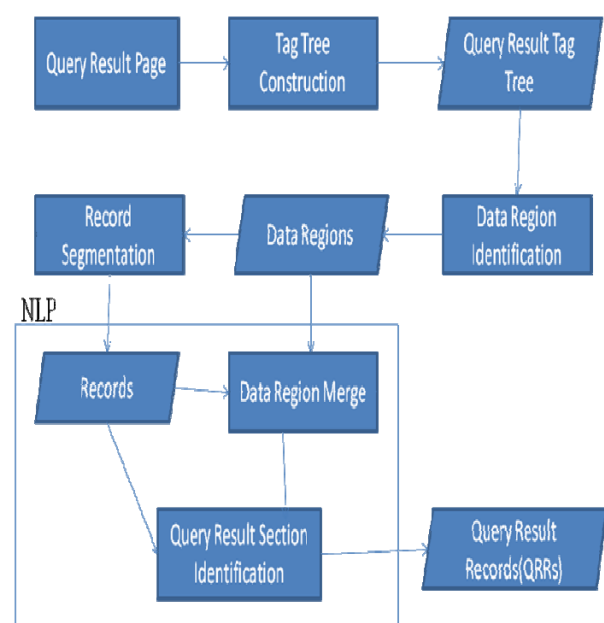happiness to constant attribute area unit placed into constant table column.



**FIG 1.1 ARCHITECTURE DIAGRAM**

## I. EXISTING SYSTEM

The problem of web data extraction has received a lot of attention in recent years. It focuses on the problem of extracting data records that are encoded in the query result pages generated by the web databases. In general query result page contains not only the actual data, but also other information, such as navigation panels, advertisements, comments, information about hosting sites and so on.The system removes any irrelevant information from the query result page, extract the records and then align the records into table such that data values of the same attribute are placed into the same table column.

The system proposes a brand new measure, revision, to evaluate the performance of internet data extraction tools. it is the percentage of the web informationbases whose data records or data things cannot be absolutely extracted. the net databases, manual revision of the extraction rules are required to realize excellent extraction. a completely unique technique is planned to perform information extraction from deep sites and realize similarity tag values exploitation information science.

## II. PROPOSED SYSTEM

### 2.1 Content Extraction

In this module, the content of hypertext mark-up language File is found. By victimisation classification technique, the complexness of finding clone is reduced simply. classification is principally accustomed count the tags of same sort. It makes issues easier. This module compares the index of tags between 2 internet applications. If 2 internet applications contain similar tags in similar variety, it is over that clones square measure existing. when classification of tags, indices of tags square measure computed. Then, analyzing of tags and tag classification provides the trail of finding clone share of tags existing between 2 internet applications. If clone share is additional, then plagiarism of code is conformed else plagiarism of code is extremely less.

### 2.2 Query Result Section Identification

In this module the info records square measure the contents focussed on deep websites. website designers perpetually have the region containing the info records centrally and prominently placed on pages to capture the user's attention. By work an oversized variety of deep websites, 2 fascinating facts square measure found. First, knowledge regions square measure perpetually situated within the middle section horizontally on deep websites. Second, the scale of {a knowledge|a knowledge|an information} region is sometimes giant once there square measure enough knowledge records within the data region. the particular size of a knowledge region could modification greatly as a result of it's not solely influenced by the quantity of information records retrieved, however conjointly by what info is enclosed in every knowledge record. Therefore, the projected approach uses the magnitude relation of the scale of the info region to the scale of whole deep website rather than the particular size.

### 2.3 Structure Retrieval

This module method 2 layouts. In model 1, the information records area unit organized in a very single column equally, tho' they'll differ wide and height. low frequency implies that the information records have an equivalent distance to the left boundary of the information region In model 2, knowledge records area unit organized in multiple columns, and therefore the knowledge records within the same column have an equivalent distance to the left boundary of the data region.

As a result of most deep websites follow the primary model, the main focus is especially on the primary model during this paper, and therefore the second model are often addressed with minor implementation growth to our current approach. additionally, knowledge records don't overlap, which implies that the regions of various knowledge records are often separated.

### 2.4 Data Value Similarity Calculation

The Web pages are preprocessed ahead, and their signatures are saved within the indexed system.The user will will specify a collection of sensitive keywords (e.g., "eBay" for World Wide Web.ebay.com) throughout searching method to indexed Server.

If a webpage contains the sensitive words, the URLs within the indexed files are parsed out. URLs within the protected list are deleted and also the URLs found within the link list ar reported like a shot. the remainder of the URLs are suspected ones, that are compared with the protected website related to the sensitive word(s) exploitation the planned approach. If the visual similarity is larger than the edge related to the protected website, the suspected website is classifed as grouping structure and an alert is reported . initial decide their information varieties then work them as deeply as doable into the nodes n1 and n2 of the information sort tree In this tree structure webpage is the root therein it extract every and each image, video and content of the page.And the datatype tree is made per the datatype worth that ar utilized in the webpage.

## III. CONCLUSION

A novel knowledge extraction methodology, CTVS has been bestowed to mechanically extract QRRs from a question result page.CTVS employs 2 steps for this task. the primary step identifies and segments the QRRs. It improves on existing techniques by permitting the QRRs during a knowledge region to be non-contiguous. The second step aligns the info values among the QRRs. a completely unique alignment methodology is planned within which the alignment is performed in 3 consecutive steps: pair-wise alignment, holistic alignment and nested structure process.

## FUTURE ENHANCEMENTS

The long run development of the project is that it reduces time quality and will increase house quality. It is utilized in Google for simple extraction and alignment of the merchandise for giant quantity of information.

## REFERENCES

[1] A. Arasu and H. Garcia-Molina,(2003) "Extracting Structured Data from Web Pages", Proc. ACM SIGMOD, pp. 337-348.
[2] R. Baeza-Yates, (1989) "Algorithms for String Matching: A Survey", ACM SIGIR Forum, vol. 23, no. 3-4, pp. 34-58, 1989.
[3] R. Baumgartner, S. Flesca, and G. Gottlob, (2001) "Visual Web Information Extraction with Lixto", Proc. 27th Int'l Conf. Very Large Data Bases, pp. 119-128.
[4] W. Cohen, M. Hurst, and L. Jensen, (2002) "A Flexible Learning System for Wrapping Tables and Lists in HTML Documents", Proc. 11th World Wide Web Conf., pp. 232-241.
[5] P. Bonizzoni and G.D. Vedova,(2001) "The Complexity of Multiple Sequence Alignment with SP-score That Is A Metric", Theoretical Computer Science, vol. 259, no. 1-2, pp. 63-79.
[6] D. Buttler, L. Liu, and C. Pu,(2001) "A Fully Automated Object Extraction System for the World Wide Web", Proc. 21st Int'l Conf. Distributed Computing Systems, pp. 361-370.
[7] K. C.-C. Chang, B. He, C. Li, M. Patel, and Z. Zhang,(2004) "Structured Databases on the Web: Observations and Implications", SIGMOD Record, vol. 33, no. 3, pp. 61-70.
[8] C.H. Chang and S.C. Lui,(2001) "IEPAD: Information Extraction Based on Pattern Discovery", Proc. 10th World Wide Web Conf.,pp. 681-688.
[9] L. Chen, H.M. Jamil, and N. Wang, "Automatic Composite Wrapper Generation for Semi-Structured Biological Data Based on Table Structure Identification," SIGMOD Record, vol. 33, no.2, pp. 58-64, 2004.

[10] W. Cohen, M. Hurst, and L. Jensen, "A Flexible Learning System for Wrapping Tables and Lists in HTML Documents,"Proc. 11th World Wide Web Conf., pp. 232-241, 2002.

## AUTHORS

**M.S.VINU** has obtained her Post Graduate degree, M.E.(Computer Science and Engineering) in Nandha College of Engineering, Erode and obtained her Graduate degree B.E.,(Computer Science and Engineering) from VSB College of Engineering, Karur. She is currently serving as Assistant Professor of Department of Computer Science and Engineering at Sri Eshwar College of Engineering, Coimbatore, Tamil Nadu with a teaching experience of 2 years. She is specializing in the area of Network Security. India). Her area is Network Security and Wireless Sensor Network.

**L.DHANAM** received her B.E(CSE) Degree from Hindusthan college of Engineering and Technology,Coimbatore, Tamilnadu, India and pursuing M.E (CSE) Degree from Sri Eshwar College of Engineering, Coimbatore, India. Her field of Interest is Network security, Operating system and Theory of Computation.

**Dr.Sudha Mohanram,** She graduated her Bachelor of engineering from Government College of Engineering, Salem and pursued her Masters in Engineering at Coimbatore Institute of Technology. She did her PhD at Government College of Technology, Coimbatore and was awarded Doctoral degree in Electrical Engineering by Anna University Chennai in 2010. She started her teaching profession as Lecturer in Government College of Technology, Coimbatore. She possesses 18 years of teaching experience. When she was about 13 years into teaching profession, her family founded Sri Eshwar College of Engineering in 2008. She has been playing the role of Secretary till 2011 and became the Principal of the institution in 2011. She has steered the institution to be one of the most sought after institutions in Coimbatore, within a short span of time through its laudable achievement in Academic excellence and Placement. She has also been interviewed by media for her adeptness and positive leadership style. She has published many papers in leading journals.