

## BOT Virtual Guide

Ms. Nitisha R. Tungar<sup>1</sup>, Ms. Nutan V. Avhad<sup>2</sup>, Ms. Pranoti P. Gayakhe<sup>3</sup>, Ms. Rutuja V. Musmade<sup>4</sup>,  
<sup>5</sup>Mr. Uttam R. Patole

<sup>1</sup>Computer Engineering, SVIT, Nashik,

<sup>2</sup>Computer Engineering, SVIT, Nashik,

<sup>3</sup>Computer Engineering, SVIT, Nashik,

<sup>4</sup>Computer Engineering, SVIT, Nashik,

<sup>5</sup>Assistant Professor, Dept of Computer Engineering, SVIT, Nashik, Maharashtra, India

\*\*\*

**Abstract** - This paper proposes a general solution for the School timetabling problem. As all staff is busy and the end time lecture conduct is severe problem for college. So to automatic Virtual Guide is been implemented which will extract web content based on recent topic been taught. An enormous amount of learning material is needed for the e-learning content management system to be effective. This has led to the difficulty of locating suitable learning materials for a particular learning topic, creating the need for automatic exploration of good content within the learning context. We aim to tackle this need by proposing a novel approach to find out good materials from www for eLearning content management system. This work presents domain ontology concepts based query method for searching documents from web and proposes concept and term based ranking system for obtaining the ranked seed documents which is then used by a concept-focused crawling system. The set of crawled documents so obtained would be obtained an appropriate set of content material for building an e-learning content management system. The filtered data crawled will be provided with speech output.

**Key Words:** DOM Parser, Web Crawler, text to speech, speech to text.

### 1.INTRODUCTION

This work proposes that Information Retrieval (IR) techniques and technologies could be specifically designed to traverse the WWW and centrally collect educational resources, categorized by topic area. IR systems are generally concerned with receiving a users information need in textual form and finding relevant documents which satisfy that need from a specific collection of documents [3]. Most existing content retrieval techniques rely on indexing keywords. Unfortunately, keywords or index terms alone cannot adequately capture the document contents, resulting in poor retrieval performance [7]. Typically, the information need is expressed as a combination of keywords and a set of constraints. However, here we use learning terms associated with topic under consideration extracted from the domain ontology. These topics and learning terms are used in the concept based query method. In addition, this work proposes a concept and term based ranking system for ordering the documents from search engine to obtain a ranked list of seed documents. With the appearance of sophisticated search engines, finding materials for e-learning is not a problem. However, the resources that one discovers might have varying styles and may be targeted at different type

of audiences. The resources may not have a complete coverage of topics which the instructor actually requires for content authoring. Moreover, a number of resources which are retrieved are highly redundant [4]. Hence, appropriate ranking of documents using concept and topic learning terms possibly will help in retrieving topic related documents and reducing redundancy from retrieved content. In this work, the ranking system exploits the concept-document similarity of the document collection. These ranked documents could then be used as seed documents for our proposed crawling system.

Similar to the work described earlier, the proposed system also used the concepts of the ontology to query the web to obtain seed documents. The ontology used by us is however specially designed a compute science ontology based on the ACM classification hierarchy. The association of terms to concepts for specific purposes has been used by Info Web a filtering system using user profiles in a digital library scenario. Here the semantic network used to represent the user profile has nodes representing concepts and as more information is gathered about the user the profile is enhanced by associating additional weighted keywords with these concept nodes. This idea has been used in the work described in this paper where in the ontology, each node in addition to having concepts from ACM classification, has an associated set of topic learning terms typically used when teaching this topic. At present this set of associated topic learning terms is manually obtained from typical texts covering the topic. As a future enhancement we propose to enhance this ontology through machine learning techniques. The search using concepts and topic learning terms from the ontology retrieves a set of seed documents.

The pace of growth of the world-wide body of available information in digital format (text and audiovisual) constitute a permanent challenge for content retrieval technologies [1]. The popularity of exchange and dissemination of content through the web has created a huge amount of educational resources and the challenge of locating suitable learning references specific to a learning topic has become a big challenge [2]. As the web grows it will become increasingly difficult for educators to discover and aggregate collections of relevant and useful educational content. There is, as yet, no centralized method of discovering, aggregating and utilizing educational content [3]. This work proposes that Information Retrieval (IR) techniques and technologies could be specifically designed to traverse the WWW and centrally collect educational resources, categorized by topic area. IR systems are generally concerned with receiving a users information need in textual form and finding relevant documents which satisfy that need from a specific collection of documents [3]. Most existing content retrieval techniques rely on indexing keywords. Unfortunately, keywords or index terms alone cannot adequately capture the

document contents, resulting in poor retrieval performance [7]. Typically, the information need is expressed as a combination of keywords and a set of constraints. However, here we use learning terms associated with topic under consideration extracted from the domain ontology. These topics and learning terms are used in the concept based query method. In addition, this work proposes a concept and term based ranking system for ordering the documents from search engine to obtain a ranked list of seed documents. With the appearance of sophisticated search engines, finding materials for e- learning is not a problem. However, the resources that one discovers might have varying styles and may be targeted at different type of audiences. The resources may not have a complete coverage of topics which the instructor actually requires for content authoring. Moreover, a number of resources which are retrieved are highly redundant [4]. Hence, appropriate ranking of documents using concept and topic learning terms possibly will help in retrieving topic related documents and reducing redundancy from retrieved content. In this work, the ranking system exploits the concept-document similarity of the document collection. These ranked documents could then be used as seed documents for our proposed crawling system.

## 2. Problem Statement

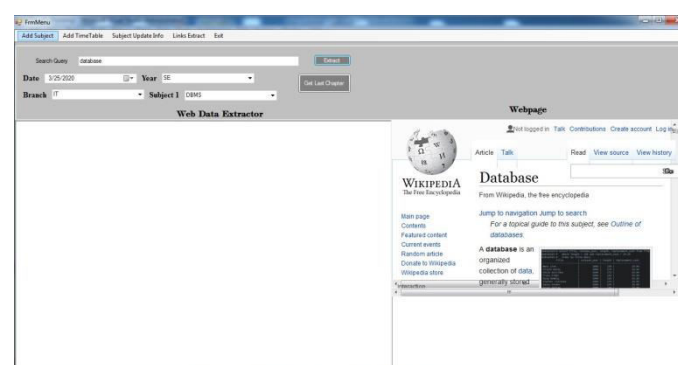
Timetabling is known to be a non-polynomial complete problem i.e. there is no known efficient way to locate a solution. To provide virtual guide using web crawler BOT. As the web grows it will become increasingly difficult for educators to discover and aggregate collections of relevant and useful educational content. There is, as yet, no centralized method of discovering, aggregating and utilizing educational content.

## 3. Proposed Methodology

### Project Modules:

#### A.Web Parsing:

The DOM Parser interface provides the ability to parse XML or HTML source code from a string into a DOM Document. DOM parser is intended for working with XML as an object graph (a tree like structure) in memory – so called “Document Object Model (DOM)“. In first, the parser traverses the input XML file and creates DOM objects corresponding to the nodes in XML file. These DOM objects are linked together in a tree like structure. Once the parser is done with parsing process, we get this tree-like DOM object structure back from it. Now we can traverse the DOM structure back and forth as we want – to get/update/delete data from it.



#### B.Text To Speech:

Text-to-Speech (TTS) encoder decoder architectures. These auto encoders learn features from speech only and text only datasets by switching the encoders and decoders used in the ASR and TTS models.

#### C.Pattern Mining:

This pattern step is designed to handle set-typed data, where multiple values occur; thus, a naive approach is to discover repetitive patterns in the input. However, there can be many repetitive patterns discovered and a pattern can be embedded in another pattern, which makes the deduction of the template difficult. The good news is that we can neglect the effect of missing attributes (optional data) since they are handled in the previous step. Thus, we should focus on how repetitive patterns are merged to deduce the data structure. In this section, we detect every consecutive repetitive pattern (tandem repeat) and merge them (by deleting all occurrences except for the first one) from small length to large length world situations. It requires a small amount of training data to estimate the parameters.

#### D.Speech Recognition:

Speech recognition is an interdisciplinary subfield of computational linguistics that develops methodologies and technologies that enables the recognition and translation of spoken language into text by computers. It is also known as automatic speech recognition (ASR), computer speech recognition or speech to text (STT). It incorporates knowledge and research in the linguistics, computer science, and electrical engineering fields. Some speech recognition systems require "training" (also called "enrollment") where an individual speaker reads text or isolated vocabulary into the system. The system analyzes the person's specific voice and uses it to fine-tune the recognition of that person's speech, resulting in increased accuracy. Systems that do not use training are called "speaker independent" systems. Systems that use training are called "speaker dependent".

To provide virtual guide using web crawler BOT. As the web grows it will become increasingly difficult for educators to discover and aggregate collections of relevant and useful educational content. There is, as yet, no centralized method of discovering, aggregating and utilizing educational content.

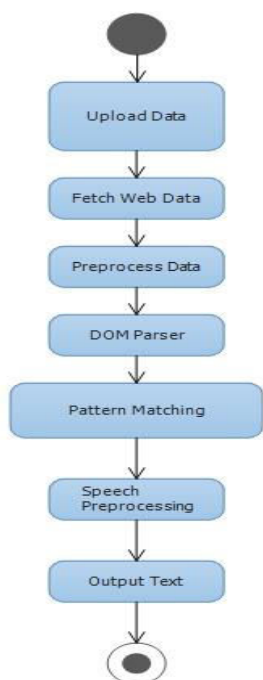


Fig 3: Activity diagram of proposed system

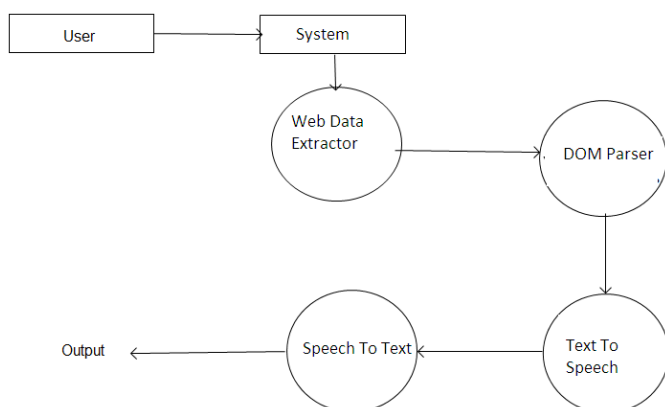


Fig 2: Data Flow diagram of proposed system

#### E.DOM Parser:

According to our page generation model, data instances of the same type have the same path from the root in the DOM trees of the input pages. Thus, our algorithm does not need to merge similar subtrees from different levels and the task to merge multiple trees can be broken down from a tree level to a string level. Starting from root nodes `<html>` of all input DOM trees, which belong to some type constructor we want to discover, our algorithm applies a new multiple string alignment algorithm to their first-level child nodes. There are at least two advantages in this design. First, as the number of child nodes under a parent node is much smaller than the number of nodes in the whole DOM tree or the number of HTML tags in a Webpage, thus, the effort for multiple string alignment here is less than that of two complete page alignments in RoadRunner. Second, nodes with the same tag name (but with different functions) can be better differentiated by the subtrees they represent, which is an important feature. Instead, our algorithm will recognize such nodes as peer nodes and denote the same symbol for those child nodes to facilitate the following string alignment.

After the string alignment step, we conduct pattern mining on the aligned string  $S$  to discover all possible repeats (set type

data) from length 1 to length  $|S|$ . After removing extra occurrences of the discovered pattern, we can then decide whether data are an option or not based on their occurrence vector. The four steps, peer node recognition, string alignment, pattern mining, and optional node detection, involve typical ideas that are used in current research on Web data extraction. However, they are redesigned or applied in a different sequence and scenario to solve key issues in page-level data extraction.

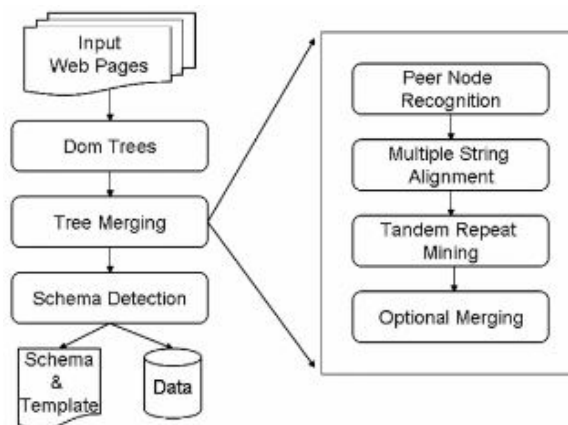


Fig 3: DOM Parser Working

### 3. Software Testing

#### .Net Automated Testing Tool:

HP Win Runner software was an automated functional GUI testing tool that allowed a user to record and play back user interface (UI) interactions as test scripts. As a functional test suite, it worked with HP Quick Test Professional and supported enterprise quality assurance. It captured, verified and replayed user interactions automatically, in order to identify defects and determine whether business processes worked as designed. The software implemented a proprietary Test Script Language (TSL) that allowed customization and parameterization of user input.

Two type of recording in Win Runner.:

1. Context Sensitive recording records the operations you perform on your application by identifying Graphical User Interface (GUI) objects. Winrunner identifies all the objects in your window you click like menus, windows, lists, buttons and the type of operation you perform such as enable, move, select etc.
2. Analog recording records keyboard input, mouse clicks, and the precise x- and y-coordinates traveled by the mouse pointer across the screen i.e Winrunner records exact co-ordinates traveled by mouse.

What is the purpose of loading WinRunner Add-Ins?

Add-Ins are used in WinRunner to load functions specific to the particular add-in to the memory. While creating a script only those functions in the add-in selected will be listed in the function generator and while executing the script only those functions in the loaded add-in will be executed else WinRunner will give an error message saying it does not recognize the function.

What are the reasons that WinRunner fails to identify GUI

object?

WinRunner fails to identify an object in a GUI due to various reasons.

1. The object is not a standard windows object.
2. If the browser used is not compatible with the WinRunner version, GUI Map Editor will not be able to learn any of the objects displayed in the browser window.

#### Unit Testing:

Unit testing, also known as component testing refers to tests that verify the functionality of a specific section of code, usually at the function level. In an object-oriented environment, this is usually at the class level, and the minimal unit tests include the constructors and destructors. These types of tests are usually written by developers as they work on code (white-box style), to ensure that the specific function is working as expected. One function might have multiple tests, to catch corner cases or other branches in the code. Unit testing alone cannot verify the functionality of a piece of software, but rather is used to assure that the building blocks the software uses work independently of each other.

#### Integration Testing:

Integration testing is any type of software testing that seeks to verify the interfaces between components against a software design. Software components may be integrated in an iterative way or all together. Normally the former is considered a better practice since it allows interface issues to be localised more quickly and fixed. Integration testing works to expose defects in the interfaces and interaction between integrated components (modules). Progressively larger groups of tested software components corresponding to elements of the architectural design are integrated and tested until the software works as a system.

#### Validation Testing:

The process of evaluating software during the development process or at the end of the development process to determine whether it satisfies specified business requirements. Validation Testing ensures that the product actually meets the client's needs. It can also be defined as to demonstrate that the product fulfills its intended use when deployed on appropriate environment.

#### GUI Testing:

GUI testing is a process to test application's user interface and to detect if application is functionally correct. GUI testing involves carrying set of tasks and comparing the result of same with the expected output and ability to repeat same set of tasks multiple times with different data input and same level of accuracy. GUI Testing includes how the application handles keyboard and mouse events, how different GUI components like menu bars, toolbars, dialogs, buttons, edit fields, list controls, images etc. reacts to user input and whether or not it performs in the desired manner. Implementing GUI testing for your application early in the software development cycle speeds up development improves quality and reduces risks towards the end of the cycle. GUI Testing can be performed both manually with a human tester or could be performed automatically with use of a software program. TO test whether .net and java GUI is properly managed as per flow in use case diagram. To test all controls

of In GUI testing check weather .Net module GUI is been Working properly.

## 4.Test Cases and Test Results

**Table -1:** Test cases result

Sr. No	Description	Expected Result	Actual Result	Test Result
1	GUI Working	All the menus and buttons of the data should work properly.	System GUI options are working properly.	Pass
2	Database Connectivity	System should do proper connectivity with database. To retrieve Lecture Information	System is doing proper connectivity with database. To retrieve Lecture Information	Pass
3	Speech to Text	System should convert speech to text properly	System is converting speech to text properly	Pass
4	Web Extraction	System should extract Web information as required using DOM parser	System is extracting Web information as required using DOM parser	Pass
5	Text To Speech	System should perform text to speech properly	System is performing speech to text properly.	Pass

## 5. CONCLUSIONS

The motto of this BOT Virtual Guide is obtaining seed documents from search engine and presented a concept-focused crawling system for the discovery of educational content from the web. In Future we will provide an Android application for the same working.

## REFERENCES

1. Nitisha R. Tungar, Nutan V. Avhad, Pranoti P. Gayakhe, Rutuja V. Musmade, Mr. Uttam Patole, "BOT Virtual Guide", IRJETPublication, Volume: 04, Issue: 06 June 2018
2. Khairil Imran Bin Ghauth, Nor Aniza Abdullah, Building an E-Learning Recommender System using Vector Space Model and Good Learners Average Rating, Advanced Learning Technologies, 2009. ICALT 2009 Ninth IEEE International Conference (15-17 July 2009), pp 194 - 196
3. Lawless, S."Leveraging Content from Open Corpus Sources for Technology Enhanced Learning", Ph.D Thesis, Submitted to the University of Dublin, Trinity College, 2009.



4. Brusilovsky, P. Henze, N. Open Corpus Adaptive Educational Hypermedia. In The Adaptive Web: Methods and Strategies of Web Personalisation, Lecture Notes in Computer Science, vol. 4321, Berlin: Springer Verlag, pp. 671-696. 2007.
5. Chakrabarti, S., Punera, K., Subramanyam, M. Accelerated Focused Crawling through Online Relevance Feedback. In proceedings of the Eleventh International World Wide Web Conference, WWW2002, Honolulu, Hawaii, USA. May 7-11, 2002.
6. Sparck-Jones, K (1972).A statistical interpretation of term specificity and its application in retrieval, Journal of Documentation 2004, Volume 60 Number 5 pp. 493-502
7. DIK L. LEE, Document Ranking and the Vector-Space Model, Software, IEEE (Mar/Apr 1997) Volume: 14 Issue: 2 pp 67 75.
8. Mehrnoush Shamsfard, Azadeh Nematzadeh, and Sarah Motiee, ORank: An Ontology Based System for Ranking Documents, International Journal of Computer Science, Vol. 1, No. 3. (2006), pp. 225-231.
9. Udit Sajjanhar, Focused Web Crawling for E- Learning Content, M.Tech Thesis to be submitted Indian Institute of Technology Kharagpur , April 2008
10. Jun Li Kazutaka Furuse Kazunori Yamaguchi, Focused Crawling by Exploiting Anchor Text Using Decision Tree, Proceeding WWW '05 Special interest tracks and posters of the 14th international conference on World Wide Web ACM , pp.1190-1191
11. Hiep Phuc Luong Susan Gauch Qiang Wang, Ontology-based Focused crawling, International Conference on Information, Process, and Knowledge Management, Cancun, Mexico, Feb. 1-7, 2009, pp123-128.
12. Marc Ehrig Alexander Maedche, Ontology- Focused Crawling of Web Documents, Proceeding SAC '03 Proceedings of the 2003 ACM symposium on applied computing, ACM
13. Gentili, G., Micarelli, A., Sciarone, F.: Infoweb: An Adaptive Information Filtering System for the Cultural Heritage Domain. Applied Artificial Intelligence 17(8-9) (2003) 715-744
14. Chakrabarti, S., van den Berg, M., Dom, B. Focused Crawling: A New Approach to Topic- Specific Web Resource Discovery. In The International Journal of Computer and Telecommunications Networking, Vol.31(11-16), Elsevier North-Holland, Inc. New York, USA. pp. 1623-1640, May 1999.



**Name:** Rutuja V. Musmade  
**Educational Details:**  
BE Computer (Pursuing)



**Name:** Uttam R. Patole  
**Educational Details:**  
M.Tech(CSE)

## BIOGRAPHIES (Optional not mandatory )



**Name:** Nitisha R. Tungar  
**Educational Details:**  
BE Computer (Pursuing)



**Name:** Nutan V. Avhad  
**Educational Details:**  
BE Computer (Pursuing)



**Name:** Pranoti P. Gayakhe  
**Educational Details:**  
BE Computer (Pursuing)