# Effective Method for Creation of Domain Module Information from Electronic Textbook

**Satyasheela T.Kadam**
*Research Scholar, Department of Computer Engineering,*
*Bharati Vidyapeeth Deemed University College of Engineering, Pune, India.*
*Orcid Id: 0000-0001-9035-1397*

**Prof. Dr. Devendrasingh Thakore**
*Professor and Head, Department of Computer Engineering,*
*Bharati Vidyapeeth Deemed University College of Engineering, Pune, India.*
*Orcid Id: 0000-0003-4381-3645*

## Abstract

In the majority education situation technology supported systems have turned out to be very supportive. This type of systems has need of the idyllic re-presentation of the information to become learned, called domain module. The domain module authoring is charge and effort demanding, however its progress charge may be decreased simply by earning through semi-automatic Site Module authoring methods along with encouraging again and again use of information. Natural language processing techniques, ontologies, heuristic reasoning is used typically for DOM-Sortze system to the semi-automatic building of the domain module from electronic-textbooks. It has been examined with an electronic-textbook to locate how it assist in the domain module authoring methods and the collected information has been compared  through the domain module , that instructional designers build up by hand. How the accuracy of identification of pedagogical relationships is increased is analyzed. This system shows DOM-Sortze and explains the research takeaway and used for present result of different type of the electronic textbooks and result generated will be the domain module information which stores more appropriate pedagogical relation.

**Keywords:** Domain module, Ontology, Learning Domain Object, Learning Object, Didactic Resources.

## INTRODUCTION

The education has affected by information and communication technology (ICTs), providing techniques to improve together teaching in addition to learning methods. At the present time, technology- supported mastering systems (TSLSs), similar to intellectual teaching techniques (ITSs), adaptive hypermedia techniques (AHSs), as well as, particularly, learning supervision systems (LMSs) like Moodle1 or Blackboard, are extensively utilized in lots of educational organizations in addition to becoming necessary for learning. Additionally, an optimistic connection between using web-based mastering technology in addition to student suggestion and desirable learning results has been observed.

To work, TSLSs need the correct Domain Component, i.e., the educational representation from the domain to get learned. Website module is the center regarding any TSLSs for the reason that it corresponds to the records concerning a topic material to get corresponds for the beginner. Website module permits also these learners to master through their own, during the condition of investigative mastering structures, or direct learners above the education practice in instructive TSLSs. A partial domain Component might glow a structure that is simply capable in the direction of supply piece of the teaching necessary within domain, the computerized or semiautomatic production from the Learning Domain Object (LDO).

Domain component intended for TSLSs has been hardly ever tackled. Chen et al. offered a head item meant for regularly construction ITSs throughout instrument understandable representations regarding textbooks, and wished for a place in the direction of construct ITSs throughout worksheets. The primary necessitates these instructional producers in the direction of print out the workbook to your official re-presentation which might be developed, as the concluding is restricted to this mathematics site. However, there have been lots of creates and try to semi- routinely collect site ontology's during various resources. TEXCOMON trace site ontology from a small number of manuscript-supported LOs through the intention of enhancing them by means of an enlarge of information. Onto-Learn have been use to generate ontologies with regard to tourism as well as economy. The thought utilizes on the net glossaries, corpora, and manuscripts because resource for that ontology mastering method. This category of development presents DOM-Sortze, a construction for that semiautomatic creation from the Domain Component from digital textbooks. DOM-Sortze plans to obtain domain independent, so completely no domain-detailed information is in employment apart from this practiced digital textbook and also the information collected from that. DOM-Sortze is just

not intended in construction exhaustive site ontology, but in given that supports to enlarge ontology with regard to educational reasons. While many ontology mastering approaches merge a lot of positive features or are bounded to definite detailed domains, DOM-Sortze is domain-independent, and is dependent completely about the electronic book presented .The structure intends to provides on several manuscript, for the reason that it does not think about the domain relates to it. The components of the system do not rely on inbuilt domain-specific information the connected information believe as the domain matters and the connections among them defined on the learning domain ontology, it presents the input information for the Learning Objects withdrawal procedure collected with the manuscript to be inspected. Artificial intelligence methods supply the means for the semi-automatic building of the domain modules from e- workbooks which may considerably supply in the direction of decrease the charge of the domain modules.

## LITERATURE REVIEW

The moral fiber of information re-presentation intended for semantic web is ontologies. Within the Domain of computer supported teaching, its concluded with the purpose of ontologies be able to take part in a most important responsibility within the prospect of intellectual teaching methods and also electronic-education information bases. The instructive Semantic- Web is an proposal intended to develop education surroundings among Semantic-Web languages and re-presentations. here, Ontologies can perform as a general and again usable information support so as to instruction methods can be used again for learning reasons, supplied that this type of systems stick on to the Domain information observation conveyed in the ontology[1].

It is truth that, information is not at all a confirmed thing: it develops with new innovation and usages. With the purpose of maintain ontologies improved through such progress, techniques in the direction of construct them in addition to expand otherwise revise them should be present set up. efficient character of information entails with the intention of physical process utilized to construct Domain Ontologies are unscalable, they are moment as well as power-consuming and re-present information since a situate construction set up at the instant the Ontology was visualized as well as construct. With the purpose of minor these negative aspects and keep away from the incredible try of constantly opening in excess of once more, mechanical techniques used intended for Domain Ontology construction should be approved. a variety of domain manuscripts can be used as a resource of information. To take into custody the outlook of a assured community using domain manuscripts can assist conserve a agreement between community members. The ontology "emerges" as of wordings; it's feasible to give details the existence of an exacting idea, belongings, example, or characteristic.

Ontology learning from manuscripts be able to assist hold to some extent of a semantic authority via given that signifies to submit to the new manuscripts. Sections and sentences are acquired from every manuscript throughout IBM-depend Java annotators. Basic verdicts are taken out through running a basic verdict miner that pull together verdicts that consist of assured keywords. Keywords are extracted throughout a keyword finding algorithm. Basic verdicts finding assists to decrease the dimension of the mass to be investigated via a linguistic parser. It also assists to spotlight on statistically appropriate statements as well as their interaction through additional statements or ideas. Every verdict is after that parsed during the Stanford Parser, that produces a formed dependant net-work, is called as a grammatical conception chart. Linguistic investigation can also depend on constituency or addiction grammar. While dependency connections are instinctively appropriate for semantic explanations, a dependency re-presentation is chosen. in addition, dependency pathways have been utilize in a number of forms to take out information for example query and responding, rephrasing, etc., and have shown their legality as information removal patterns. While the TEXCOMON. Intention is to stay Domain Independent, a grammar directed technique is used to model lexico-syntactic outlines addicted to sorted dependencies sub trees. Every outline is categorized as a hierarchy approximately a root expression t that re-presents a variable that enter and crop grammatical connections. Every joint in the outline re-presents a variable. Throughout the investigation method, every time an incidence is originated these outlines are sought after in the manuscript and instantiated with information [2].

The generated LOs are stored in a ZIP file that contains the XML-based representation of the LO based on the ALOCOM formalism, as well as the referenced images or resources. The XML-based ALOCOM document format allows content reusing. Harmonization or duplication of content is an additional important use case. Such as, in the MELT project, construction of a doorway that utilizes federated investigate methods among a variety of repositories has been completed. In this network, each doubt is distributed to each and every one collaborator in the network. A investigate similar to „natural science" proceeds further more15,000 outcomes, which sets in excess of progression weight on the doorway as data about data for all outcomes are conveyed above the network. Harmonization of data about data with a vital data about data store avoids from having to concern disseminated investigates. These have need of an API that permits for harmonization of data and/or content [3] .While the initial large case learning of the latest AMG structure, we have been raising a structure that indexes the entire matter that is formed in the circumstance of ProLearn Network of superiority on specialized education. To handle the entire ProLearn manuscripts a shared workplace structure is used, which is the Agora Groupware Web Server. In a primary footstep we

spotlight on the deliverables that are manufactured inside ProLearn. After that, the further matter like the documents is processed. even though the creation of data about data for a learning entity is the key focal point of our employment so far, we also would like to create use of that data about data afterwards, for instance for investigating. To generate this investigate functionality we wrote an addition of AMG with the grouping of a Lucene index, and the Simple Query Interface (SQI). Lucene is build up by the Apache crowd, and supply a high-performance, filled-featured wording investigate engine library, which we used for store up the created data about data. SQI is a meaning of web services that allow querying Learning Object Repositories in a ordinary method. In this method, it describes the query interface that we will apply for investigating the produced data about data [4].

Dublin Core metadata element set (DCMES) is utilized via a number of e-learning initiatives for the account of learning things. An ISO standard for data about data is DCMES, proposed for annoyed-domain source explanation. The data about data standard consist of two ranks: easy and competent. Center of easy Dubin includes fifteen factors: heading, inventor, topic, explanation, publisher, supplier, blind date, category, layout, identifier, resource, speech, relative, reporting and human rights. competent Dublin Core consist of three supplementary factors (viewers, attribution and privileges holder), and a collection of component modifications that build the sense of an component easy or more specific. The teaching working cluster of the Dublin Core data about data proposal is budding teaching exact components, component qualifiers and controlled vocabularies to be used with DCMES for telling instructive matters. Among others, the DC-Education application suggests the use of three components from the LOM data about data standard: Interactivity category, representative education moment and Interactivity level. If this is requirement of learning management system then user will provide great number of learning object [4]. So, there was need for automatic generation of generating domain ontology and then use this ontology in computer based education as a domain module [5].

To recognize such sections of manuscript, the method, initial an inner hierarchy-like re-presentation is build from the e-document. To obtain the part-of-speech information Linguistic investigation is carry out on the hierarchy-like manuscript re-presentation. This information might be not so essential for various languages such as English, wherever usual expressions can be used to recognize illustrations and meanings, it is necessary for further languages for example Basque. In Basque for the construction of statements the vocabulary entry gets each one of the components required for the dissimilar purposes (syntactic case included).For this cause Basque is called as agglutinative language, that is., More purposely, the affixes keep in touch to the concluder,

numeral and declension case are in use, in this type sort separately of each other. Since prepositional purposes are recognized by case suffix within word form, Basque shows a comparative more control to create inflected word forms that builds morpho-syntactic investigation extremely vital for being capable to take out information from manuscript pieces. Improving the part-of-words information by tagging the commonness of the domain subject matters recorded in the domain ontology is completed in the afterward footstep. Such ontology identifies the subject matters and the connections amongst them. Formerly the indications to domain subject matters have been noticeable, DRs, such as meanings, illustrations, problem statements, etc., are recognized. The recognition of DRs is carry out by identical prototypes normally used to describe contents, current examples, etc. The amount produced of this footstep is a collection of minute DRs that are improved subsequently in two techniques: it has been observed that dissimilar DRs are joined by human teachers to obtain additional correct matter consequently, if they are alike or secure sufficient, according to the content and the type of DR successive DRs are come together. For sake of brevity, a detail on the composition of DRs is not supplied [6], [7].

To elicit keywords from text documents, there are some studies on this kind of text mining technique widely known as keyword extraction [8],[9],[10].To ensure that the developed curriculum was comprised of eligible contents, many studies have proposed curriculum evaluation methods to accomplish this task[11], [12]. Learning objects (LOs) and their reusability are one of the most important current research topics in the learning technology community. Ontology provides a mechanism to capture information about the ideas, concepts, and the relationships between them in some domain [13,15]. Preprocessing [14] is important to work on pdf.

## PROPOSED SYSTEM

The key plan of system is to permit learner to study through their own, in the case of shortage of additional education system or supposed to be supply them management regarding teaching scheme in significant method. So building of Domain Module consist of
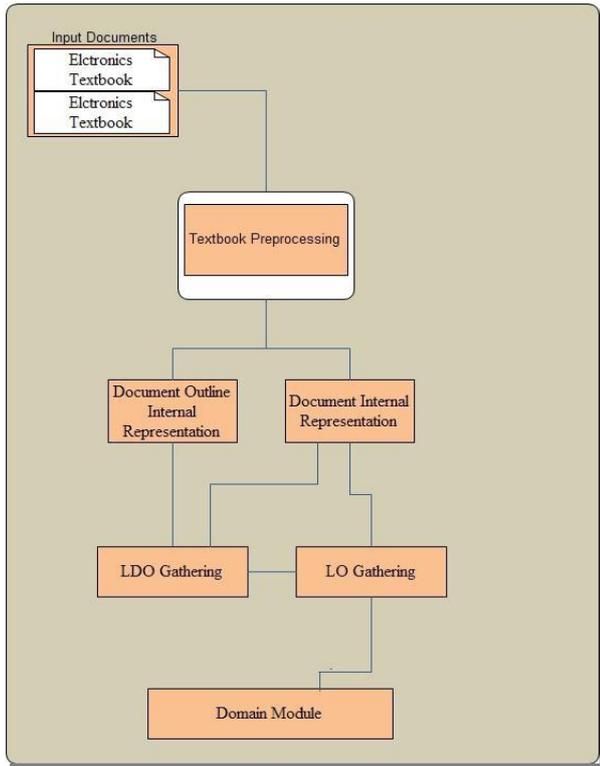
**Figure 1:** System architecture

## PREPROCESSING OF TEXTBOOK

For information gaining procedure we arrange the manuscript primary and then result are utilized to gather information instructed at two stages of information into domain module.

### *Textbook preprocessing*

With the purpose of run information acquirement procedure we have to build e-document and collects identical re-presentation of it but as these manuscripts are accessible in a lot of dissimilar design like text file, pdf, etc. this type of pre-processing is essential. There is utilization of i.e tree-like construction to arrange content of the manuscript. It will consist of hierarchy i.e tree like structure as manuscripts which consist of chapters which supplementary consist of different section. after converting file Part of speech analysis by using NLP is done.

### *Keyword Extraction*

TF-IDF:

The Term frequency–Inverse document frequency (tf-idf) which is invented at 1988 by salton & buckley technique merges the manifestations of word into a manuscript among frequencies of the manuscripts wherein the word is discover in a indication mass to conclude its termhood. At other side, word frequency calculates the consequence of the word. Extra

repeatedly a word come into views in a manuscript, the extra related it. Alternatively, the IDF determines the elasticity of the word. The term is less specific when it appears in lots of manuscripts.

Method for term withdrawal permits the recognition of term since a document. In addition TF, this type of method as well think about the consequence of the words within the manuscript. Term frequency-inverse document frequency necessitates a lots of quantity info to differentiate ordinary words from the related ones, individuals which come out in the investigated manuscript however are not regularly utilized in any supplementary document. As this intended that this work is domain independent. Natural language processing techniques uses tf-idf.

### *Term Frequency:*

In this case tf(t,d),Make use of raw count of the words in the manuscript is straightforward option ,that is the amount of moments that word $t$ take places in manuscript $d$. If we represent the unprocessed tally through $f_{t,d}$, after that the straightforward tf method is

$$\mathrm{Tf}(t,d) = f_{t,d}. \qquad (1)$$

Further promises consist of

- Boolean'frequencies': tf($t,d$)= 1 if $t$ arises in $d$ and 0 or else;

- For document length term frequency is corrected, length: $f_{t,d} \div$ (number of words in d)

- logarithmically scaled frequency: tf($t,d$) = 1 + log $f_{t,d}$, or zero if $f_{t,d}$ is zero;

### *Inverse document frequency*

It is a used to determine of how much information the word presents, that is, whether the term is ordinary or uncommon crossways all manuscripts. It is the logarithmically scaled inverse fraction of the manuscripts that include the word, gained by dividing the whole number of manuscripts by the number of manuscripts holding the term, and after that taking the logarithm of that quotient.

$$\mathrm{idf}(t, D) = \log \frac{N}{|\{d \in D : t \in d\}|} \qquad (2)$$

After that term frequency-inverse term frequency is computed as: A more load in tf-idf is attained through a more tf in the specified manuscript and a less manuscript frequency of the word in the total gathering of manuscripts; the loads therefore are inclined to sort out ordinary words. Because the proportion within the inverse document frequencies log purpose is forever larger than or equivalent to 1, the

assessment of inverse document frequency and tf-idf  is larger than or equivalent to 0. As a word come into views in much more manuscripts, the proportion within the logarithm advances 1, taking the idf and tf-idf nearer to 0.

$$tfidf(t,d,D) = tf(t,d) * idf(t,D) \qquad (3)$$

## Gathering of LDO

In this step educational association among domain subject matters that is mastered, supposed to be recognized and represented into the appearance of LDO that assists technology supported learning structure to permit learner to direct their own throughout education procedure.

### Analysis of outline

This step consists of different most important stages:

### Basic investigation

For outline internal representation major subject matter of domain and connection between these subject matters are extracted from the outline internal re-presentation. So guide point is believes as key subject matter and subordinate thing which explains part of it so structural association is explained among point and sub point.

### Heuristic investigation

It extracts latest connection stand on the beforehand experienced collection of the heuristic for structural relationship:

These permits to recognizing type of connection among thing of outline and its part of things, it labors on investigation that just single type of relative survives among a thing and part of thing and highest moment it is occurred, that connection which is to be establish is a category of part of relation.

At this type of situation we observed that:

1.  Here survive is A connections among the outline thing and its entire part of thing if assembly heuristics triggers.

2.  Or else for each part of thing an each heuristic that competitions should be applied.

## Collection of  LOs From Manuscript

In this step the meaning of LOs, illustrations, implements, etc are utilized the duration of the education procedure that are recognized as well as created every footstep explained in point.

### Collection of LOs

This footstep is educational association among the Domain subject matter is stock up by Learning Domain Object therefore educational connection consist of structural relations that is., is A, piece of, precondition, after that in that P is A Q connection point out that P is a kind of Q.P part Of Q point out that P is a part of Q, P requisite Q point out that P have to be master to tutor Q.

The creation of learning objects intended for  domain subject matters be realized by recognizing and assembly Didactic Resources, that is, reliable portions of the manuscript associated to more than one subjects, through a exacting instructive reason recognition and removal of this parts is take out in an ontology-determined procedure, which also utilizes Natural Language process method. As the LO creating method described in this research means to be not depend on domain , simply domain-detailed information utilized is LDO  which is collected from the e- workbook in the preceding stage. From this time forth, a DRs will pass on to a part of the manuscript destined to be utilized in the education meetings (for example meaning, implement…, ) at the same time as a learning objects make reference to a again used Didactic Resources developed by data about data. The LOs creation method explained is happened by erauzOnt that is element of the DOM–Sortze construction. DR grammars, LDOs, communication indicators as well as a educational ontology are used to collect DR from the inside re-presentation of the e-workbook with the part of speech information. To construct LOs from the collected DRs the LDOs as well as ALCOM ontology and in conclusion, the LOs are stored in the LOR to make easy their use and again use.

## RESULT

The mainly significant preparation of this system is to evaluate DOM-Sortze module wherever it tends the teacher's to construct the module through processing the information in the LDO and LOs collected from the textbook resources. The system process is on LO where the e-workbook in those workbooks the images was processed. To calculate approximately the method of production of the Module utilizing DOM -Sortze module by combining an indicated DR and LDO results. The instruction values generated related to domain areas and the educational connections to describe the Learning Domain Object structure. The invention of LDO is estimated and processed supported for automatically gathered knowledge and is gathered among domain topic and pedagogical relationships. The evaluation of the LO generation is now being considered as the identified LOs, to which end the automatically gathered LOs were compared to the identified DRs.

In this research, figure 2.shows electronic-textbook is an input given for purpose of experimentation, the electronic-text book may be any text file, pdf etc. you can upload any electronic-textbook for creating the information domain module.
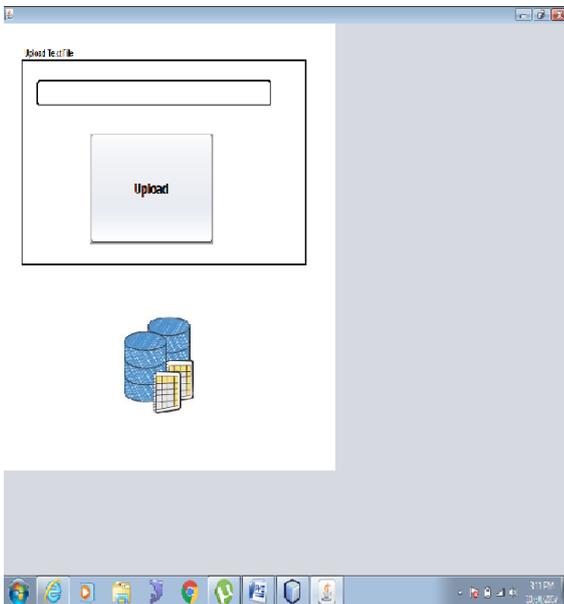


**Figure 2**: Result of Uploading File

The following result shows that, for generating the domain module of related topic you can give that particular topic as a keyword for example, "program" is one of the keyword which is taken as input and result will be generated. It will take very less time for generating output.
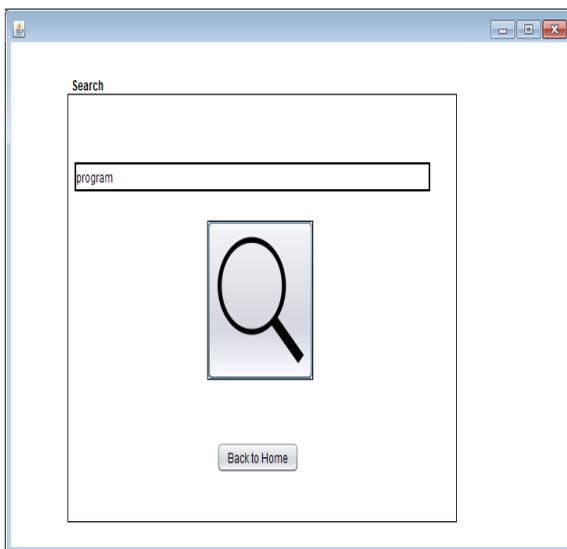


**Figure 3**: Result of Keyword

Figure 4.Shows that, there will be first conversion happens and after that LOS and LDO gathering happens. In this project

we observed that the application developed is to process e-books in the form of PDF file to generate LO's.  As all topics are not studied at a time, the proposed application gives an option to generate LO's from user's choice of pages in a file. The application will analysis the PDF file and displays information like, total size of file in and total pages in a file. Then application will ask for from and to page numbers for further processing. Application will extract the contents from selected pages and present to user. Then   user either goes for generation of learning modules or selects another from-to page numbers. After pre-processing of e-book, the LOs are generated from user's choice of from -to pages.
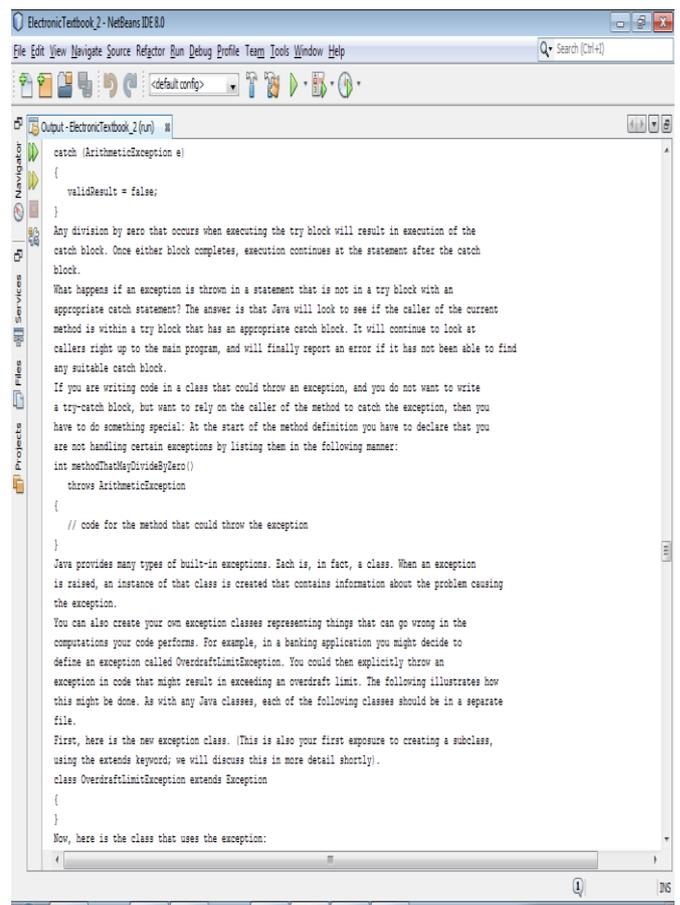


**Figure 4**: Result of Conversion

Figure 5.shows that Accuracy is increased and we got more appropriate results which shows pedagogical relation. After generation, generated LOs are presented to user in a form of title list. On selecting the title, the detailed LOs will be displayed. All the accepted learning modules are formed together in the sequence of selection for generating final learning module to be saved or take confirmation by user.
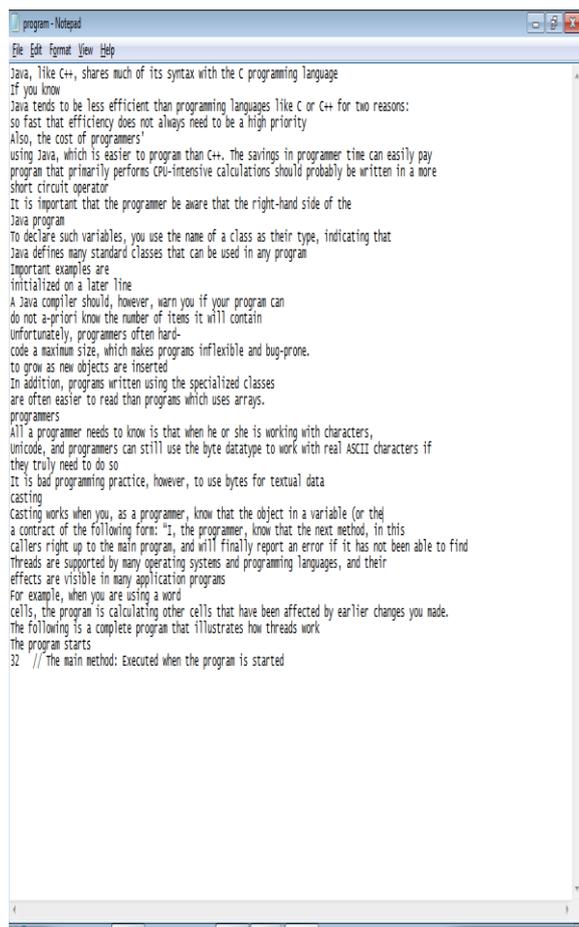
**Figure 5**: Information Domain Module

**Table 1**: Execution time

|  | File 1 | File 2 | File 3 | File 4 |
|---|---|---|---|---|
| Size of input file(kb) | 110 | 80 | 64.8 | 95 |
| Time to execute the file in msec | 2378 | 1298 | 1507 | 2095 |

Research shows the file size and time required. For keyword extraction, time will be in (millisecond) and the size of the input file in (kb).for the purpose of experimentation we take different files of the different size and work will be done on it, so we got different time for the purpose of the execution.
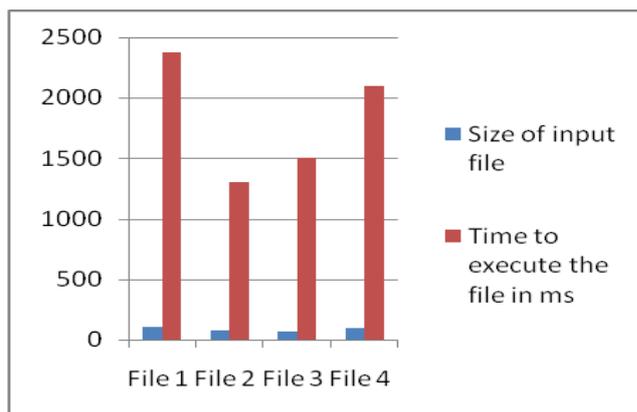


**Figure 6**: Comparison of Size of Input File and Time of Execution

In this system, we plotted the graph like input files in (kb) with respect to time to execute the file in (ms) for comparisons of various files. In this way, we have selected some Input files of different sizes for better output and also executing the time in (ms) which is varying continuously that shown in above graph. In this Graph we observed that file1 have required more time to execute file in ms, also file 2 have required less time to execute file.

**Table 2**: Precision and Recall

| File No. | Precision | Recall |
|---|---|---|
| File1 | 84.5 | 77.27 |
| File2 | 89.36 | 73.68 |

DOM-Sortze a Domain module shows precision which is increased and also improvement in the recall also. In this research there are two files which are taken in the kb and the result will be calculated.
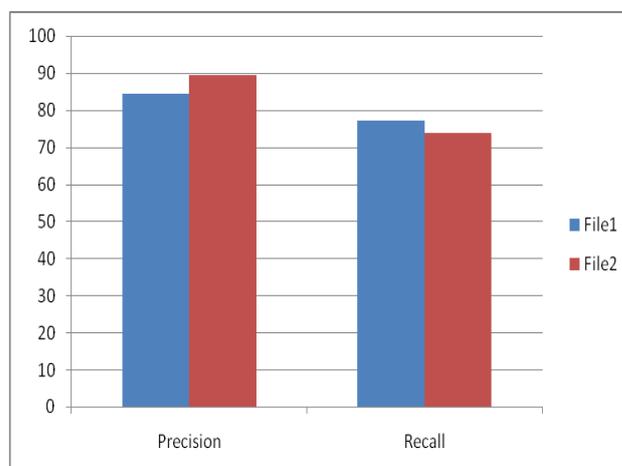


**Figure 7**: Precision and Recall for files

**Advantages**

1. If the input manuscripts are given or presented appropriately then existing system work   more effectively but proposed system works effectively on different types of manuscripts .

2. Additional pedagogical connections are to be recognized. Accuracy is increased.

3. The accuracy of Identification of pedagogical relationships is increased.

4. System is very simple to read the content of Electronic textbooks and gives accurate result.

**CONCLUSION**

These kinds of technique have introduced DOM-Sortze, method for real semi-automatic production of the Domain Module from e-textbooks. The method makes use of technique like heuristic thinking, ontologies, and NLP technique to the information acquire methods. DOM-Sortze may be tested by using an e-textbooks and contrasting on auto-pilot generated elements, while using the Domain Module not by automatically put together via educational markers. Every effort was to estimate how DOM-Sortze plays a role in Domain designing. The e-manuscript employed for try was a list of the volumes. Because experiment geared to determine the data acquirement by wording, an edition not including diagrams of the manuscript seemed to be utilized because resource is connected with information. At this time, DOM-Sortze is being developed to help newly arrived languages for example English.  This system is used for present result of different type of the electronic textbooks and result generated will be the domain module information. which stores more appropriate pedagogical relation and time required for execution is less.

**REFERENCES**

[1]   Mikel Larranaga and Jon A. Elorriga" Automatic Generation of the Domain Module from    Electronic Textbboks: Method and Validation", IEEE Trans. Knowledge and Data Eng., vol. 26, no. 1, Jan .2014.

[2]   A. Zouaq and R. Nkambou, "Evaluating the Generation of Domain Ontologies in the Knowledge Puzzle Project," IEEE Trans. Knowledge and Data Eng., vol. 21, no. 11, pp. 1559- 1572, Nov. 2009.

[3]   K. Verbert, X. Ochoa, and E. Duval, "The ALOCOM Framework Towards Scalable Content Reuse," J. Digital Information, vol. 9,no. 1, 2008.

[4]   M. Meire, X. Ochoa, and E. Duval, "SAmgI: Automatic Metadata Generation v2.0," Proc. World Conf. Educational Multimedia, Hypermedia, and Telecomm. (ED-MEDIA ΄07), pp. 1195-1204, June 2007.

[5]   A. Zouaq and R. Nkambou, "Evaluating the Generation of Domain Ontologies in the Knowledge Puzzle Project," IEEE Trans. Knowledge and Data Eng., vol. 21, no. 11, pp. 1559-1572, Nov. 2009.

[6]   A. Maedche and S. Staab, "Ontology Learning for the Semantic Web," IEEE Intelligent Systems, vol. 16, no. 2, pp. 72-79, Mar. 2001.

[7]   M.T. Pazienza and A. Stellato, "Semi-Automatic Ontology Development: Processes and Resources" eds., IGI Global, 2012.

[8]   Lott, B., "Survey of Keyword Extraction Techniques." UNM Education, (2012).

[9]   Hasan, K. S. and V. Ng.,  "Automatic Keyphrase Extraction: ASurvey of the State of the Art". ACL (1), (2014).

[10]   Salton, G. and C. Buckley , "Term-weighting approaches in automatic text retrieval." Information processing & management 24(5):513-523, (1988).

[11]   Welch, W. W. and H. J. Walberg , "A national experiment in curriculum evaluation." American Educational Research Journal: 373- 383.

[12]   Dai, Y., et al., "Course Content Analysis: An Initiative Step toward Learning Object Recommendation Systems for MOOC Learners."

[13]   M.-j. YUN Hong-yan, XU Jian-liang and X. Jing,"Development of domain ontology for e- learning course," ITIME-09IEEE international symposium, 2009.

[14]   Siddarth Banga et.al, "Regression and Augumentation Analytics on Earth Surface Temperature", IJCST 2017.

[15]   Duval, E., "Standard for learning Object Metadata". editor. 1484.12.1 IEEE June 2002.