

# Survey: Analysis and Design of Plagiarism Software for Regional Language

Prashanth Kumar HM<sup>1</sup>, Dr. Subramanya Bhat S<sup>2</sup>

<sup>1</sup>Student, College of Computer Science, Srinivas University, Mangalore, India

<sup>2</sup>Professor, College of Computer Science, Srinivas University, Mangalore, India  
E-mail: [prashanth.hm02\[at\]gmail.com](mailto:prashanth.hm02[at]gmail.com)

**Abstract:** Copying and pasting ideas or data across languages has created many challenges in copyrights. Though many methods have been developed to detect plagiarized content in English and English related language, however, detecting all language it is like left to right (Ex: English or other) and right to left (Ex: Urdu or other) language plagiarism is still a challenge. This study uses advanced Encrypted Text Matching (ETM) and Text Encryption Indexing (TEI) to build all language semantic space, from which it checks the contextual similarity in indexed data and internet data. Still there is no plagiarism detection software for language to providing single platform checker, as we studied in a major plagiarism detection software's in a world like 'Turnitin' providing 30+ language, 'Urkund' Providing 20+ languages and 'DrillBit' Supporting 20+ Languages majorly English related. Here we are working for a single platform to show all language text matching process by using UTF-8, 16 and other Unicode text matching technology.

**Keywords:** Text Matching, Indexing, Internet Data, Encryption, Decryption

## 1. Introduction

Plagiarism is the demonstrations of taking somebody else written job, conversation, tune, or even supposed and transient it off as yours own. This integrates information from website books, pages, tunes, network shows, interviews, email messages, published articles, work of art or whatever other standard. At whatnot point you reword, abbreviate, or expressions, take words or sentences from someone else work or job, it is significant to establish the well of the data inner your paper developing an interior citation. It is insufficient to merely list the sources in a list of source near the finish of your article. Neglecting to legitimately refer to, quote, or know another person's data or thought with an inward reference is counterfeiting, as we know that we have many English or English related open or paid plagiarism detection tool in worldwide, those are checking similarity between authors work and indexed data from company database or internet world data. But challenges are not in English literature, it's a main challenge in some other Unicode or UTF-16 or other special character similarity checking in same indexed data from company database or internet world data. Here we are going to start research and implementation to the same idea will establish into market.

## 2. Objectives

DrillBit SoftTech India Pvt. Ltd 'DrillBit Plagiarism' India By now we comprehend what DrillBit is and what it does. Presently how about we perceive how this remarkable Software functions. The working of it begins when it is nourished with an information i.e. say any record. DrillBit exclusively takes the lines from the document and checks for the comparability over web and indexed data.

**String Matching:** The problem of find presence(s) of a pattern string to another string or body of text data. There are so many various algorithms for good searching. Also called as exact text matching, text searching, string

searching. The name "exact text matching" contrasts with string matching with errors. Text matching concept or algorithms gives a key role in different real-world problem or application in plagiarism world. A one or other of its imperative application are Spell Checker, Spam Filter, Intrusion Detection, Search Engine, Similarity Detection, Bioinformatic, Digital Forensic with Information's Retrieval System etc.

**Index:** Indexing is the practice of compiling and finding data into a single string or comparing data to such a group of string. Indexing is also used to refer to passively finding attached in meaning full information indexes to replicate broad data returns rather than actively selecting individual inputs. One common type of indexing in plagiarism is called "search engine indexing." Here, plagiarism tools aggregate and interpret search engine data, again, to streamline data retrieval. This type of indexing is also sometimes called data indexing. Plagiarism developer explain that indexing helps to make searches less labour intensive without an index, the search engine would have to search every document at its disposal equally, whereas with an index, much of this work is eliminated. One example of indexing is the legacy Microsoft Indexing Service, which maintained an index of files on a computer or in an operating system environment. Another example is database indexing, which involves creating an index for a database structure to help expedite retrieval of data.

**Data Computing:** In computing, data is information that has been translated into a form that is efficient for movement or processing. Relative to today's computers and transmission media, data is information converted into binary digital form. It is acceptable for data to be used as a singular subject or a plural subject. Computing is any activity that uses computers to manage, process, and communicate information between plagiarism operation. It includes development of both matching algorithm and indexed data. Major plagiarism data computing disciplines include in a base of artificial intelligence and machine

learning with area referred in computer engineering, software engineering, computer science, information systems, and information technology.

**Encoding Basic: UTF-8:** For the standard ASCII (0-127) characters, the UTF-8 codes are identical. This makes UTF-8 ideal if backwards compatibility is required with existing ASCII text, in plagiarism majority of English language and English related languages support only here in UTF-8. Other characters require anywhere from 2-4 bytes. This is done by reserving some bits in each of these bytes to indicate that it is part of a multi-byte character. In particular, the first bit of each byte is 1 to avoid clashing with the ASCII characters.

**UTF-16:** For valid BMP (Basic Multilingual Plane) characters, the UTF-16 representation is simply its code point. However, for non-BMP characters UTF-16 introduces surrogate pairs. In this case a combination of 2 two-byte portions map to a non-BMP character. These two-byte portions come from the BMP numeric range but are guaranteed by the Unicode standard to be invalid as BMP characters. In addition, since UTF-16 has two bytes as its basic unit, it is affected by endianness. To compensate, a reserved byte order mark can be placed at the beginning of a data stream which indicates endianness. Thus, if you are reading UTF-16 input, and no endianness is specified, you must check for this.

### 3. Methodology: Data Collection & Analysis

**Input Level:** The most common document file extensions we are use in plagiarism input file formats. Those are listed called doc, docx, html, htm, odt, pdf, xls, xlsx, ppt, pptx or txt. During implementation, there are many occasions to use to extract data from certain filetype by using proper extraction algorithm. Our ddocument Conversion API offers a way to convert between document files programmatically with maintain exact user data to further process.

**Character and String Data Types:** We are encoded data during the plagiarism processing in different level. If they are raw bytes or symbols in report, the time we try to give output non-ASCII characters is very economic, here we may run into a few problems. Also, even if the character type is based on a UTF, that doesn't mean the strings are proper UTF. They may allow byte sequences that are illegal. Generally, we will have to use a library that supports UTF/Unicode, such as support for C, C++, Python, dot net and Java. In any case, if we want to input/output something other than the default encoding, we will have to convert it first.

**Comparing for Equality:** Identifying different language characters is not a big task in programming, but it's a big task when comparing different language string to other equivalent string in multilingual comparing concept. Example A, A, and A look the same, but they're Latin, Cyrillic, and Greek respectively. We also have cases like C and C, one is a letter, the other a Roman numeral. In addition, we have the combining characters to consider as well.

**Multilingual Use Unicode, UTF-8, UTF-16:** Multilingual is most used to describe someone who can speak or understand multiple languages, especially someone who can write report in several languages with some level of fluency. The ability to write multiple languages or the use of multiple languages.

Unicode: A set of characters used around the world.

UTF-8: A character encoding capable of encoding all possible characters (called code points) in Unicode.

Code unit is 8-bits: use one to four code units to encode Unicode

00100100 for "\$" (one 8-bits);11000010 10100010 for "¢" (two 8-bits);11100010 10000010 10101100 for "€" (three 8-bits)

UTF-16: Another character encoding.

code unit is 16-bits. use one to two code units to encode Unicode

00000000 00100100 for "\$" (one 16-bits);11011000 01010010 11011111 01100010 for "𐀀" (two 16-bits)

#### Buffer Compatability: Writing to Buffer

If we write to a 4-byte buffer, symbol あ with UTF8 encoding, your binary will look like this:

00000000 11100011 10000001 10000010

If we write to a 4-byte buffer, symbol あ with UTF16 encoding, your binary will look like this:

00000000 00000000 00110000 01000010

As we can see, depending on what language we would use in our content this will affect our memory accordingly during plagiarism checking process in a software.

E.g., For this symbol あ UTF16 encoding is more efficient since we have 2 spare bytes to use for the next symbol. But it doesn't mean that we must use UTF16 for Japan alphabet.

#### Reading from Buffer

Now if we want to read the bytes, we must know in what encoding it was written to and decode it back correctly.

E.g., If we decode this: 00000000 11100011 10000001 10000010 into UTF16 encoding, we will end up with 𐀀 not あ.

30 42 (hex) - > UTF8 encoding - > E3 81 82 (hex), which is above result in binary.

30 42 (hex) - > UTF16 encoding - > 30 42 (hex), which is above result in binary.

### 4. Survey

**Survey 1: Turnitin Text Matching** “In existing framework we have taken US Established Plagiarism programming TurnItIn, this creation is world best running more than 15,000 foundations including India and they are servicing 30+ languages majority in English literature. Turnitin is an Internet-based unoriginality avoidance benefit made by iParadigms, LLC, first propelled in 1997. Ordinarily, colleges and secondary schools buying licenses to submit

articles to the Turnitin site, which authorizations the records for probable substance.”[1]

“The outcomes can be employed to recognize likenesses to existing sources or can be utilized as a part of growing evaluation to enable understudies to figure out how to stay away from written falsification and enhance their composition. Understudies might be required by schools to submit expositions to Turnitin, as an obstacle to copyright infringement. This has been a wellspring of comment, with a few understudies diminishing to do as such in the conviction that requiring it constitutes an assumption of blame. Furthermore, commentators have claimed that utilization of the creation damages instructive security and licensed innovation laws.”[1]

**Survey 2: Ouriginal Text Matching:** “According his company information, Ouriginal supports and enables academic institutions, secondary schools, and corporates to improve efficiency and deliver quality content by providing an automated system for assessing the authenticity and originality of any text. Ouriginal was established in 2020 when two of the industry’s leading giants, Urkund and PlagScan, joined forces to improve academic integrity and promote original thinking, and serving including India and they are servicing 30+ languages majority in English literature also. With more than three decades of combined knowledge and expertise, Ouriginal delivers cutting-edge technology that helps enhance the potential of students to think originally, saves time for teachers when evaluating assignments and assists corporates preserve their reputation.”[2]

#### A. New Related Issues

In the above two major text matching company survey, I found that, why they are supporting text matching for only limited languages. So, I came some survey conclusion in non-English languages text matching:

- 1) Text matching does not support an Indian regional language.
- 2) Text matching does not support Arabic, Aramaic, Azeri, Dhivehi/Maldivian, Hebrew, Kurdish (Sorani), Persian/Farsi, Urdu.
- 3) Knowledge of extraction data for store and indexing in storage pools.
- 4) Language wise data collection form worldwide websites.
- 5) Difficult to maintain in originality report.
- 6) Algorithms compatibility during text comparison.
- 7) Problems in extracted data to be process in different levels.
- 8) Comparison result or grade.

#### B. Proposed System

It is vital to be selective. To pick out the right alternative it frequently takes a mellow recognition on whatever it is. Written falsification as we probably are aware is such a significant issue and it takes part of legal coming through and furthermore numerous ethical that tally under it. To convey the best the association is endeavoring hard for the best outfit. We convey the best of the administrations for the understudies. Understudies here allude to everybody. In the instructive spree for the most part Drill Bit targets and

handles each test. Give it a chance to be for the instructors or for the understudies.

It is to get profited those matters for us. An adaptable utilization system without many inconveniences, thought in a laymen kind. The method for introduction has a considerable measure effect which champion, we are facing it. Innovation utilized will dependably have a refreshed diagram, adapting to the changing necessities of the clients.

- **Detection of Erroneous words:** Erroneous words allude to inaccurate words. In any case, with regards to Plagiarism, wrong words have a gigantic impact in coordinating and skipping counterfeited work. It is very testing undertaking to identify these parts in the sentences, since a few characters in the words go shrouded thus making it seem lawful word. Yet, those words can't be overlooked. This issue has been focused and explained here.
- **Identifies Case Ignored Words:** As like incorrect words, case disregarded words too hugely affect last consequence of copyright infringement.
- It is when words are inaccurately composed as far as cases keeping in mind the end goal to be dealt with as new word. Each time coordinating experiences these words they may not coordinate with the current substance but rather pass on a similar importance. To defeat this issue extraordinary care is taken in DrillBit.
- **Detection of Similarity in Data is done Line by Line:** DrillBit checks for comparative information on hold premises over the web and puts away documents in the database. It ought to be realized that line by line implies sentences. On the off chance that comparability i.e. 'five back to back words' match then the checker will stamp the information as appropriated.
- Provides 100% Direct Link To site :Here we acquainted an immediate connection idea with the counterfeited information, it is made easy to understand in order to distinguish all coordinated web copyrighted information.
- Client can confirm identified information by a simply click. Those connected information may show up in pdf, html, htm, rtf and other referenced record designs. Here the fundamental objective is clarified i.e. to make the client recognize the wellspring of literary theft effectively rather checking superfluously.
- Implemented Plagiarism lead (Five Consecutive word).
- Implemented Erroneous word's location.
- Ignored Case Sensitive received.
- Text situated outcome bolster.

#### C. System Design

##### Detection of Erroneous words:

Erroneous words refer to incorrect words. But when it comes to Plagiarism, erroneous words play a huge part in matching and skipping plagiarized work. It is quite a challenging task to detect these parts in the sentences because some characters in the words go hidden hence making it appear legal word. But those words cannot be ignored. This issue has been keenly concentrated and solved here.

##### Identifies Case Ignored Words:

Like erroneous words, case ignored words too have a huge impact on final result of plagiarism. It is when words are incorrectly written in terms of cases to be treated as new word. Every time matching encounters these words they may not match with the existing content, but they convey the same meaning. To overcome this issue special care is taken in DrillBit.

#### Detection of Similarity in Data is done Line By Line:

DrillBit checks for similar data on the line basis over the web and stored files in the database. It should be known that line by line means sentences. If similarity i.e. 'five consecutive words' matches then the checker will mark the data as plagiarized.

Provides 100% Direct Link To website:

Here we introduced a direct link concept to the plagiarized data, it is made user friendly so as to identify all matched web copyrighted data. User can verify detected data by a just click. Those linked data may appear in pdf, html, htm, rtf and other referenced file formats. Here the main goal is made clear i.e. to make the user identify the source of plagiarism easily rather than checking unnecessarily.

- Implemented Plagiarism rule (Consecutive word).
- Implemented Erroneous words detection.
- Ignored Case Sensitive adopted.
- Text oriented result support

#### D. Advantages of Proposed System

- 1) User easily buys the plagiarism using internet.
- 2) Quickly generate report.
- 3) Accuracy in detection and matching.
- 4) Remote access
- 5) Data Integrity: You will always have control on whom to share your document with. Only after your approval it is published for further works. So nowhere your data will be shared.
- 6) A A note on Deletion: when you upload your document for plagiarism check and done with the check. Later when you wish to be deleted from our servers, the deletion will be performed on your request. So always your sensitive data is kept with keen security.

#### E. Disadvantages

Without internet user drill bit server cannot be accessed.

#### F. System Architecture

Framework conformation is the way toward characterizing the design, modules, interfaces, and info for framework to fulfill indicated prerequisites. Frameworks configuration could be watched as the utilization of frameworks hypothesis to item advancement. Because of the organization secure approach, we ought not include the finish handle plan. The accompanying figure demonstrates fundamental outline module of unoriginality process.

DrillBit solely takes the lines from the document and instructions for comparability over web. It works in a precise route by picking the source where the likeness matches most extreme. All the organized lines are as needs be filed and highlighted alongside its coordinating connections. At last DrillBit produces a report with finished investigation of the record and even creates the rate of Originality work and Similarity work.

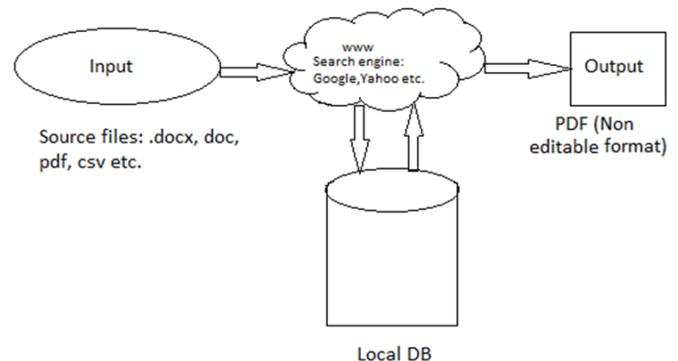


Figure 1: Architecture of the literary theft checker

Figure [1] determines the source files as doc, pdf and so on as an input to access the data to search through internet data sources to plagiarise the information to identify the duplication of data. Once it identifies the data it stores as pdf file and results into local database. Usually, the system specifies an illustration of any data to be identified in google search engine.

More often than not for programming reasons we have favored java. As it has gigantic points of interest and that is the purpose behind picking it. The following are a portion of the motivations to pick it: Portable-Code written in Java can be taken from one PC to the next without worrying about framework design points of interest. Strong: Java bolsters dependable special case taking care of that can withstand all the real sorts of incorrect and exemption conditions without breaking the framework. Secure: Upon gathering, source code written in Java gets assembled into byte code, which is later deciphered by the Java Virtual Machine.

Byte code is impervious to altering by outer operators. Stage Independent: Most of the frameworks have worked in java runtime condition, the main essential for running an application that has been outlined in Java. Accordingly, no setups or conditions must be infused into a framework before executing a Java application. Self-Memory Managed: The coder does not need to be worried about the memory co-ordination, designation and de-allotment of items. JVM takes care of it.

Elite: Both regarding memory and proficiency, Java has ended up being flawless. Prior in its history, byte code understanding was thought to be an extra duty of the compiler, which requested escalated handling and memory utilization.

Multithreading: Synchronization and multitasking come as a complimentary blessing because of Java's multithreading highlights

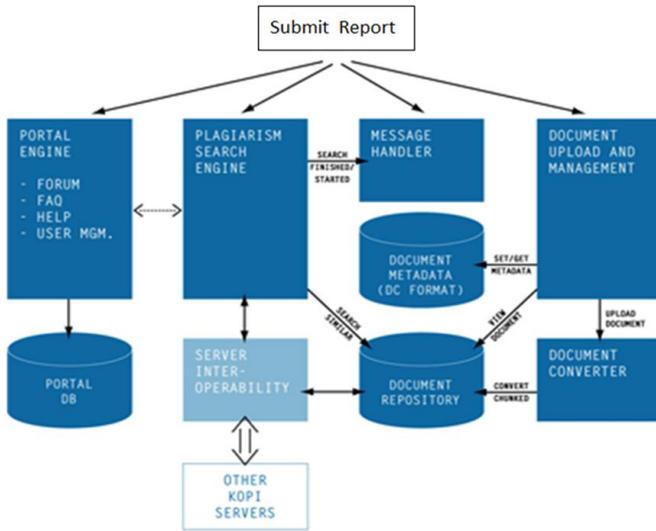


Figure 2: System Design

## 5. Results

It is the important phase in software development life cycle where project implementation phase where logical planning and requirements comes into reality. It describes about the tools, programming languages and data base management of the project. In order to implement proposed system we uses the Net bean IDE, languages like Java and python for database connectivity as JDBC.

Below figure [3] shows the proposed work of plagiarism,



Figure 3: Proposed work



Fig. 9.4(a): Plagiarized Data

This Report is Automatically Generated by DrillBit Software on 2017/06/28 16:09:28

Software Displays Only Similarity Data Page Wise

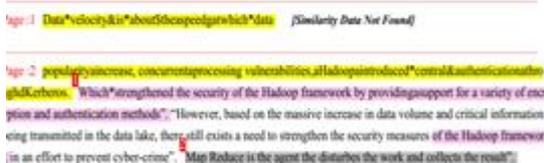


Figure 4: Plagiarized Data

Steps involved in the results part are.,

- User login page where it also login process with already registered user ,if user is new then it follow by registration step. After successful registration then allow the user to login with user name and password credential.
- Next phase is uploading phase where user documents of format doc or pdf are being uploaded. Then followed by scanning of uploaded file by scanner.
- Viewing of report in this phase where data which are copied from internet source or any articles all those things will be highlighted .Based highlighted content planarization of document will reveal to the user.

### Advantage of proposed work:

Features	Unlimited Resource	Plagiarized Links	Unlimited Document Upload	Detect Erroneous Words	Case Sensitive Check	Developed by
VIPER	Yes	Yes	No	No	No	England
Turnitin	Yes	Yes	Yes	No	No	US
DrillBit	Yes	Yes	Yes	Yes	Yes	India

Figure 5: Comparison of Drill Bit Product with Other Product

## 6. Conclusion

DrillBit exclusively takes the lines from the file and checks for similarity over the internet. DrillBit produces a report with complete analysis of the file and even generates the percentage of Originality work and Similarity work. Plagiarism detection is the process of locating instances of plagiarism within a work or document. The widespread use of computers and the advent of the Internet have made it easier to plagiarize the work of others. Most cases of plagiarism are found in academia, where documents are typically essays or reports. Here I am working with a team to support their process. Counterfeiting identification is the way toward finding examples of written falsification inside a work or archive. The far-reaching utilization of PCs and the appearance of the Internet have made it less demanding to appropriate the work of others. Most instances of written falsification are found in the scholarly world, where records are regularly articles or reports. Here working with a group to bolster their procedure.

## References

- [1] R. Iqbal, A. Grzywaczewski, J. Halloran, F. Doctor, and K. Iqbal, "Design Implications For Task-Specific Search Utilities For Retrieval and Reengineering of Code," Enterprise Information Systems, Taylor and Francis, pp. 1751-7575, 2015.
- [2] Iqbal, R., Grzywaczewski, A., James, A., Doctor, F., Halloran, J., "Investigating The Value of Retention Actions As A Source of Relevance Information In The Software Development Environment", in roc. 16th Intl. Conf. on CSCW in Design, IEEE, pp.121-127, 2012.
- [3] O. Alhabashneh, R. Iqbal, F. Doctor, and S. Amin, "Adaptive Information Retrieval System Based on Fuzzy Profiling," in Proc. Of Intl. Conf. on Fuzzy Systems, IEEE, pp.1-8, 2015.
- [4] M. Alsallal, R. Iqbal, S. Amin, and A. James, "Intrinsic Plagiarism Detection Using Latent Semantic Indexing And Stylometry," in Proc. 6th Intl. Conf. on Developments in eSystems Engineering (DeSE), IEEE pp. 145-150, 2013.
- [5] T.A.E. Eisa, N. Salim, and S. Alzahrani, "Existing Plagiarism Detection Techniques: A Systematic Mapping of The Scholarly literature," Online Information Review, vol. 39, pp.383-400, 2015.
- [6] R. Iqbal, F. Doctor, M. Romero, and A. James, "Activity-Led Learning Approach And Group Performance Analysis Using Fuzzy Rule-Based Classification Model," in Proc. 17th Intl. Conf. on CSCW in Design, IEEE, pp. 599-606, 2013.
- [7] R. Iqbal, A. James, R. Gatward, "A Practical Solution to the Integration of Collaborative Applications In Academic Environment," in Proc. 5th Intl. Workshop on Collaborative Editing Systems, hosted by the ECSCW'03, Helsinki, Finland, 2003.
- [8] R. Iqbal, A. James, and R. Gatward, "A Framework For Integration of CSCW," in Proc. 7th Intl. Conf. on CSCW in Design, IEEE, pp. 43-48, 2002.
- [9] M. Zurini, "Stylometry Metrics Selection For Creating A Model For Evaluating the Writing Style of Authors According To Their Cultural Orientation", Informatica Economica, vol.19, p.107, 2015.