# Implementation of the Plagiarism detection software used in universities

Xây dựng phần mềm phát hiện đạo văn dùng trong các trường đại học

**Tran Anh Vu [1], Nguyen Phan Kien [1], Nguyen Tuan Anh [1], Pham Thi Viet Huong [2]**

[1] School of Electronics and Telecommunications, Hanoi University of Science and Technology, Hanoi, Viet Nam
[2] University of Engineering and Technology, Vietnam National University Hanoi, Hanoi, Vietnam

**Abstract**

Plagiarism is a serious problem not only in Vietnam but also in almost all other countries. Many software packages have been proposed to detect plagiarism on the content of the text. However, most of these packages are from foreign countries, it is difficult to check the documents using Vietnamese language. The plagiarism detection software named BKCheck is designed and developed to solve this problem. This software has the ability to provide the similarity among documents based on user's database. The results are reported as the percentage of similarity among documents. Besides, the software also shows the position of similar sentences or paragraphs, which allows users to easily verify the results. It has been tested and provides reliable results in School of Electronics and Telecommunications (SET), Hanoi University of Science and Technology (HUST).

Keywords: Software, Plagiarism, Database, Checking, Thesis

## 1. Introduction

In order to graduate from universities, students are required to submit their disssertation to the commitee. Students are expected to put their motivation, knowledge and effort on it. Hence, the dissertation is said to be students' achievement after they are done with years of studying hard in universities. Although it is not very difficult to write a dissertation, it is not easy to produce a good quality, interesting and practical dissertation. It requires students have serious attitude, dilligence and excellent writing skill. Students nowadays are very lucky with the help of the booming Internet. They can easily find some online outstanding dissertation/papers for reference. Good students take advantage of the Internet by browsing and understanding what the world is doing and find some interesting ideas for their dissertations. Some lazy students are not doing that way. They may browse the Internet, find some published papers and copy the content into their dissertation. It saved their time and finally they still have a qualified dissertation in hand. As the resourse are huge and diversed, professors may not realize their sutdents was copying from somewhere else into their dissertation. Hence, the clarification and evaluation of students may not be exact. The misclarification happens not only in universites but also in the whole education system, which can cause serious problems. For example, bad students may be misclarified as good, then they are assigned a good position at work. When facing the real work, bad students of course not do well, which

can causes damage or harm to the companies. Hence, we need to prevent such a phenomena of copying without crediting the source. The plagiarism software are proposed in this paper to serve this purpose. According to the Merriam-Webster online dictionary [1] , plagiarism is:

- to steal and pass off the ideas or words of another as one's own
- use another's production without crediting the source
- present as new and original an idea or product derived from an existing source

In general, plagiarism is using ideas, sentences or products of another as one's own as default. Or in other way, use another's ideas and content without crediting the source.

Nowadays, when the Internet and computers become popular in all over the world, plagiarism has increased at fast speed and become a big challenge for the education. We should say that, plagiarism helps students meet requirements of the courses or thesis while they do not put their effort on it. Meanwhile, we cannot clarify between students who study hard and those who copy another's product. They are evaluated equally. Hence, it forms the bad habit of stealing and passing off another's content as one's own. We should have solution to solve this problem completely; otherwise, it may cause bad effect to all other students.

So what should we do? Here comes the main content of this paper. As said before, in the booming of technology, plagiarism has been so easy, we need to utilize new technology to get rid of it. Plagiarism software is proposed with the main function of

[1] Corresponding author: Tel.: (+84) 944.639.471
Email: vu.trananh@hust.edu.vn

comparing documents and finding how many percent they look alike. This paper has the following sections:

- The first part introduces some existing plagiarism software in the world as well as their advantages and disadvantages.
- The second part talks about the reality and solutions against plagiarism in Vietnam.
- The third part gives requirements for the nowadays plagiarism detection software.
- The forth part implements the plagiarism detection software which has been tested in the School of Electronics and Telecommunications, Hanoi University of Science and Technology.

## 2. Requirements of the Plagiarism detection software

### 2.1. The existing Plagiarism detection software evaluation

There are a lot Plagiarism detection software packages which are being used in the world. Following are some popular ones.

*Plagiarism Checker [3]*: This software allow users to compare documents by providing websites which have exactly the same content with documents. However, it is difficult to compare with Vietnamese documents, because it can not handle the Vietnamese font display. Hence, it is not effective to use this software in Vietnam

*Plagiarism.net [3]*: It is an online tool to check the source of documents. However, users have difficulty in account registration and search. Otherwise, it is not a free website, users have to pay in order to use it. Moreover, the result is not clear for users to have the proof of plagiarism and to evaluate the level of plagiarism.

*Turnitin [4]*: It is the popular plagiarism detection software package in the world. It can search and compare 300 million essays and 110,000 documents over the Internet. The Lotus University has applied this software into the plagiarism checking in the university. This software can be integrated with the Cloud services like Google Drive and Dropbox. Students can submit their work to the Tunirin Original Check. However, due to the high cost and the complicated use, it is not been widely used.

There are some other similar software packages such as Plagiarism Detector [4], Writecheck [5], Plagium [6], Dupli Checker [7], etc. They are all required account registration and costly so users are hesitate to use them.

From the above examples, we can see the main disadvantage of those software packages is that they can only be applied to English documents; it is difficult to check and test in Vietnamese. Besides, the

evaluation of the percentage of similarity among documents is not presented; the validity is not exact; the method is not easy to use especially for those who are not good at English. Moreover, some software requires users to pay, which is not effective to use.

### 2.2. The reality and solutions against plagiarism in Vietnam universities

The conference on "academic integrity" was held in Lotus University in May 29th, 2015. More than 20 universities have taken part in this conference. They share their research experiences on plagiarism prevention. Fig. 1 shows the Duy Tan University's survey result on students' plagiarism.
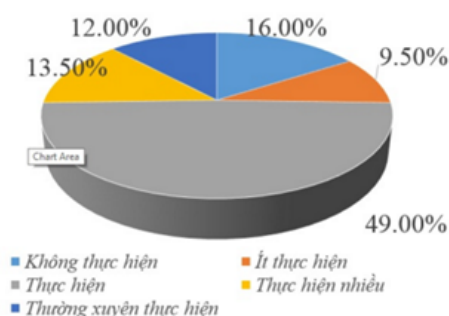


**Fig. 1.** Survey on students' plagiarism [8]

After the conference, most of universities have agreed to build a plagiarism detection software package. However, until now, two years from the conference, no plagiarism detection software has been officially proposed.

This work is expected to be the solution to the Vietnamese documents plagiarism detection. The software should handle with Vietnamese documents. It should not be costly to use, so a large number of people could use it. Moreover, the output should be clear, easy to read, so users can be able to evaluate the level of plagiarism in each document.

## 3. Introduction of BKCheck Plagiarism detection software

With the urgent demand of a Plagiarism detection software package, the BKCheck has been researched and developed, which can be satisfied all the above requirements. The software package works as follow: all the theses/dissertations are collected into a library. This work is done by an administrator, who has the right to add or remove files from library. Users are not allowed to access to the library. When a user inputs a document to check, the software will compare it with all files in the library and return the percentage of similarity between the input and each of the documents in the library.

The main program can check file by file in the library, but it takes time to open and close each file.
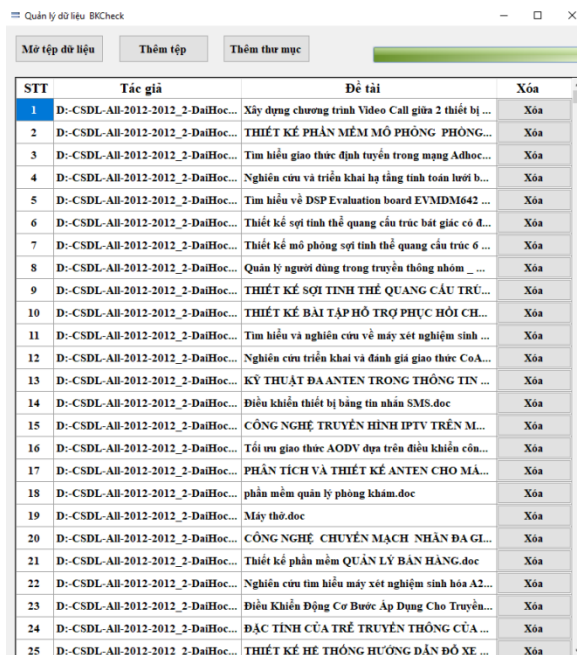
Thus, the total time for checking is increased due to the time of reopening and reclosing the files, especially when the number of files in the database is large. To prevent multiple file opening and closing, all separate files in the library are simply recorded in one database text file. The locations of each word in each file are recorded to reduce the searching time for the main program. This work is done by the database management program.

The database management program is used to manage the database text file. The interface of the database management program is as in figure 2, which represents the location and title of each file in the library.

To manage the database, administrators can add or delete files. The usable file formats are pdf, doc/docx and txt. Administrators can add each file separately, or the whole folder at the same time.

The main checking program finds and compares users' input file with each file in the library, then return the percentage of similarity.

The algorithm of checking is presented by the flow chart as in figure 3. In which:



**Fig. 2.** Interface of the database management program

*A, B* is the word array of the tested document and database document, respectively. *A.Count, B.Count* is the number of words of the tested and the database document.

*M* is the total number of copied words.

*g* is the default number of copied words, it depends on users' selection

*Algorithm explanation*

The algorithm works as follow: Start from the first word of the tested document, search for it in the database. There are two cases:

*Case 1*.	There is no such word in the database => move to second word of the tested document and continue to search for the similar word in the database. The algorithm works in this way until it reaches the last word of the tested document.

*Case 2*.	Find the similar word in the database, then:

*Step 1*: Count the similar words in the two documents, in the place we find them

*Step 2*: Continue the search process and count until reach the end of tested document

*Step 3*: Compare the numbers and output the maximum

Step 4:

➢ If the maximum is greater than the default number of copied words (chosen to be 10), then the number of copied words is assigned to the maximum. The algorithm moves to the next word and continue.

➢ If the maximum is smaller or equal to the default number of copied words, the algorithm moves to the next word and continue.

When the search is done, the software will calculate the number of the copied words of the tested document. The percentage will be calculated as:

Percentage of similarity= ( number of copied words)/(total words of the document)×100%

The document in the library which has the highest percentage of similarity will be updated continuously and presented at the top of the result table. Users can stop the program at any time to check more detail at any file they want.

In order to increase the chance to find the document which has similar content, the order of searching will prioritize documents which have the title's keywords similar to those of the input.

With the above algorithm, the result is optimized; the copied words will be found disregards of their places or formats.
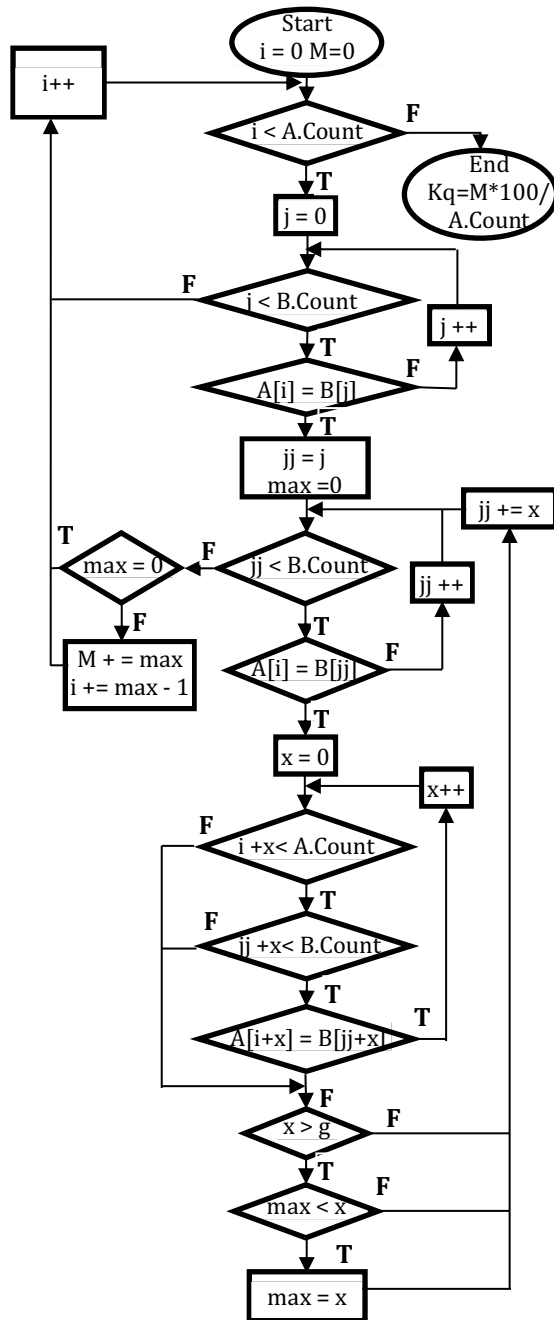
38

**Fig. 3.** The algorithm flowchart



**Fig. 4.** Display of the result



**Fig. 5.** Result extraction to Excel file

## 4. Results

All documents in the library will be uploaded to the database of the software by administrators. The content of the database can be added or remove by the administrator. When a document needs to be checked, it is loaded to the program, the software will then automatically compare it with the database. It will take only about 30 second to check a document with a database of 1300 documents. The result appears in figure 4.

The fast checking brings users the similarity between the tested documents and those in the library. The result after checking can be exported in Excel format (figure 5). Users will know how many percent the tested document looks alike with the document in the database with the order of highest similarly percentage at top position.

In some cases, users want to have the exact information about positions of the similar paragraphs or sentences between tested documents and those in the database; the program will allow users to change to the direct comparison mode. With this mode, the position of each sentence or paragraph will be marked and colored for a convenient checking. Results are illustrated in figure 6.
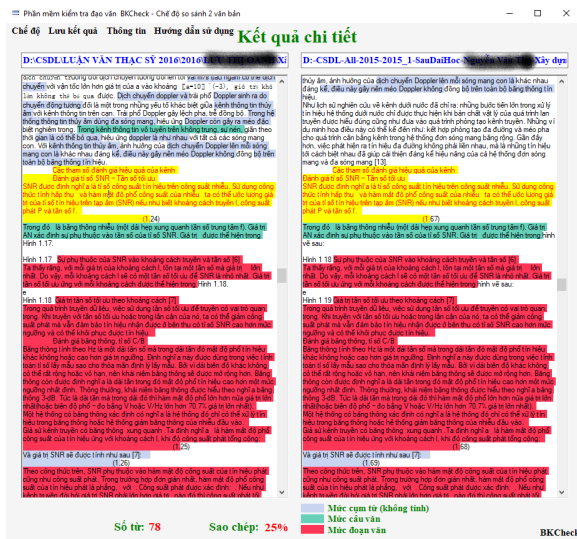
**Fig. 6.** Diplay file comparison in detail

The software also allow users to define parameters for the similarity, for example, paragraph-level similarity is marked as red (the default number of words > 50), sentence-level similarity is marked as blue (the default number of words > 10) and word-level similarity is marked as gray (this part will not be counted toward the similarity percentage).

In order for users to check conveniently, the software will mark the same position of the two files in which users are observing. If users click on any position in one file, the same position in the other file will be colored accordingly. This position is illustrated in yellow.

An improvement of this software compared to the existing ones is the result validation. By visualizing the position and defining the length of similar words, sentences or paragraphs, users are able to validate the results.

Hence, the BKCheck has following advantages:

✓ The interface is easy to use (in Vietnamese). Results are designed visually and conveniently.
✓ The speed of testing is fast (approximately 30 seconds to check a document for a database over 1300 documents)
✓ The accuracy is high. The results are displayed in detail which is easy to evaluate tested documents. It has the ability of searching similar paragraphs even when their formats are changed (words adding, sentences stop, line stop, marks changing, etc)
✓ The software is able to handle with Vietnamese documents (and English, of course) effectively.

✓ The software is able to handle with different document formats: .doc, .pdf, .txt
✓ The database is built by users; hence it is easy to manage.
✓ However, there are some disadvantages we need to get over:
- It is not able to compare images or graphs
- Due to the fact that the database is built by users, so if we want to use the software in a large scale, we need the integration of the database from different parties.

## 5. Conclusion

The BKCheck plagiarism detection software package has been researched and tested in School of Telecommunications and Electronics, Hanoi University of Science and Technology. The software satisfies the requirements such as the ability to test Vietnamese documents, friendly interface, easy database management, especially the ability to change the default number of copied words depending on users. The software has been tested with the database of more than 1300 documents, and the result is validated. In the coming days, authors will put the software online to make it easier for the testing process.

## References

[1] https://www.merriam-webster.com/dictionary/plagiarize

[2] www.grammarly.com

[3] Plagiarisma.net

[4] Turnitin.com

[5] plagiarism-detector.com

[6] WriteCheck.com

[7] www.plagium.com

[8] www.duplichecker.com

[9] http://dantri.com.vn/giao-duc-khuyen-hoc/nan-dao-van-ngay-cang-gia-tang-1433666610.htm