

Spelling Correction Application with Damerau-Levenshtein Distance to Help Teachers Examine Typographical Error in Exam Test Scripts

Viny Christanti Mawardi^{1,*}, Fendy Augustian¹, Jeanny Pragantha¹, and Stéphane Bressan²

¹Faculty of Information Technology, Tarumanagara University, Jl. Letjen. S. Parman no. 1, Jakarta 11440, Indonesia

²School of Computing, National University of Singapore, 13 Computing Drive, 117417 Singapore

Abstract. This research was intended to create Spelling Correction Application to help teachers examine questions scripts with the capability to found typographical error and give suggestion for non-real word error. This application is built with simple Damerau-Levenshtein Distance method to detect errors and give word suggestions from the typo word. This application can be used by the teacher to examine documents in the form of short answer, essay and multiple choices then save them back in the form of original documents. This application is built using a dictionary lookup consist of 41 312 words in Indonesian. The first test result is the application can detect non-real word errors from 50 sentences that have non-real word error in each sentence and produce an accuracy of 88 %. The second test is try to detect typographical error in exam test script that consist of 15 sample questions, consisting of five essay questions, five short answer, and five multiple choices.

Keywords: Examine document, proofreading, spelling correction tools, typing errors, typo word.

1 Introduction

Schools have the duty to provide the best education for their students. The teacher is one part of the school that has a role in the teaching process. One of the teacher's tasks is to assess student learning outcomes [1]. The assessment process can be done through giving a test in the form of a question. The teacher has the task of making questions to assess students' abilities. Questions can be given for training, daily assessments, midterm test or final test. As usual making exam test is done by typing questions using computer. After the questions have been made, re-examination of the questions needs to be done to re-evaluate the material.

After the teacher has typed, the teacher will re-read the questions that have been typed. Re-examination is done to maintain the quality of the questions to be given to students. Commonly, they found some errors such as typing errors that accidentally occur when typing script. Such as excess letters "saya" become "sayaa", lack of letters such as "makan" to "mkan" or letters that are swapped like "bisa" become "bias". That kind of mistakes can

* Corresponding author: viny@untar.ac.id

be found when typing such as lack of letters, excess letters or letters that are accidentally swapped when typing questions [2].

There are still some typing errors in test script when manually evaluated because of large number of questions made from grades 1 to 6 elementary school makes the re-examination process difficult and less than optimal. Based on that problem, this study build applications that can be used to correct spelling of exam test. Although spelling correction tools are widely available in various applications. Like the most commonly used document processing applications (word processors) such as Microsoft Office and Open Office, it already has the ability to check typographical error or correct the typo error when typing. The spell checker application can also be found on devices for typing email (email client), electronic dictionaries and search engines. But the application commonly used for English.

Spelling correction is one of tools proofreading that commonly found in embedded system like word processing system. Therefore, existence of automatic systems such as spell and grammar-checker or correctors can help to improve the quality of produce electronic document [3]. But the different between each language make the research of spelling correction must develop for every specific language. Where, Faili et al. [3] develop proofreading tool for Persian language that consist of detecting spelling, grammatical, and real-word errors.

The specific needs of each problem in produce electronic documents also require their own development tools to improve the quality of document. Therefore, spelling correction can develop in embedded system or standalone application. Yulianto, Arifudin and Alamsyah [4] develop spelling correction for suggest, autocomplete and spell checking word in library data searching. This research focus in the teacher needed to found typographical error after create exam test. More specific teacher needs to create a system that can simplify the process of re-examining exam questions that have different types and can directly save them back in the same form. So that the process of checking typography errors does not occur when typing is in progress but after produce exam test document and the evaluation can be done with other teacher.

2 Research method

This study build applications that can help check the typographical error after question script made. Differences in the form of manuscripts made such as essay, multiple choices or short answer will affect the application process that is built. The process of re-examining the question script is expected to be done after all the manuscripts have been typed, then the text will be re-examined by another teacher or the principal. The application that can be used to help re-examine the question script, and then give annotated on the misspelled words than the teacher will be easier in correcting the words. To build this application it is necessary to develop the spelling correction method and the typed of exam test script. This research used Damerau-Levenshtein Distance and three types of question script.

2.1 Exam test script

In education, assessment is one way to measure students' abilities and knowledge. Teacher will create assessment for measure the learning process. Assessment data can be obtained from directly examining student work to assess the achievement of learning outcomes or can be based on data from which one can make inferences about learning [5]. In general the giving of assessment is done through a written test. Teacher will create script of test for exercise, quiz, exam or anything. Currently, the test script can produce using a computer. The test can be given for student by online/offline, printed or filled directly through a computer.

There are a lot of assessment format like written test, multiple choice, matching type, essay, short answer, true false etc. Each type of assessment test can be used for student exam or exercise. Teacher will plan the type assessment for measure learning process. Teacher will create bank of test that can be used for exam, exercise or quiz. Each type of script test has different format so when teacher produce electronic document they need to be create in different way. This research choose three types of test short answer, essay and multiple choices. The example of exam test can be seen in Table 1.

Table 1. Type of exam script test.

Type of exam	Example
Short answer	1. Katak berkembang biak dengan cara...
	2. Mata akan seketika berkedip ketika terkena debu. Hal itu menunjukkan bahwa makhluk hidup peka terhadap
Essay	1. Sebutkan contoh sistem operasi pada perangkat komputer!
	2. Apa fungsi sistem operasi?
Multiple choice	1. Mana yang bukan untuk pengulangan? a. Switch case b. while c. do while d. for
	2. Mana yang bukan bagian dari kecerdasan buatan? a. Fuzzy Logic b. Image Processing c. Neural Network d. Sistem Pakar

2.2 Damerau-Levenshtein Distance

Spelling correction is the process of correcting the spelling to match the word that should be [3]. The spelling correction process is done by using a spell checker. Spell checker is a program application that is used to mark words that have less precise spelling. Spell checkers can be in the form of their own programs or become one with word processing applications. The spell checker application will check then give a sign directly when each word is typed. So the typed word can be changed automatically or manually changed by humans. Word checking can be built by comparing data in the dictionary or using various methods to predict words that should be correct.

Several studies have been conducted to build spell checkers for Indonesian. The most commonly used method is Levenshtein Distance. LD is method that has operation to calculate the number of operations needed to change a word into another word. LD calculates operations such as deletion, insertion and substitution so that is suitable for spelling correction [6]. The variation of Levenshtein Distance is Damerau-Levenshtein Distance, which adds a variety of transposition operations based on adjacent letters [Damerau in [7]]. Based on and [8], the formula from Levenshtein Distance can be seen in Equation 1, Equation 2, and Equation 3, while Damerau-Levenshtein Distance can be seen in Equation 4.

$$d_{i0} = \sum_{k=1}^i w_{del}(b_k), \text{ for } 1 \leq i \leq m \tag{1}$$

$$d_{0j} = \sum_{k=1}^j w_{ins}(a_k), \text{ for } 1 \leq j \leq m \quad (2)$$

$$d_{ij} = \begin{cases} d_{i-1,j-1} & \text{for } a_j = b_i \\ \min \begin{cases} d_{i-1,j} + w_{del}(b_i) & \text{for } 1 \leq i \leq m, 1 \leq j \leq n \\ d_{i,j-1} + w_{ins}(a_j) & \text{for } a_j \neq b_i \\ d_{i-1,j-1} + w_{sub}(a_j \cdot b_i) \end{cases} \end{cases} \quad (3)$$

$$d_{ij} = \begin{cases} d_{i-1,j-1} & \text{for } a_j = b_i \\ \min \begin{cases} d_{i-1,j} + w_{del}(b_i) & \text{for } i, j > 1, a_i = b_{j-1}, a_{i-1} = b_j \\ d_{i,j-1} + w_{ins}(a_j) & \text{for } a_j \neq b_i \\ d_{i-1,j-1} + w_{sub}(a_j \cdot b_i) \\ d_{i-2,j-2} + w_{tra} \end{cases} \\ \min \begin{cases} d_{i-1,j} + w_{del}(b_i) & \text{for } 1 \leq i \leq m, 1 \leq j \leq n \\ d_{i,j-1} + w_{ins}(a_j) & \text{for } a_j \neq b_i \\ d_{i-1,j-1} + w_{sub}(a_j \cdot b_i) \end{cases} \end{cases} \quad (4)$$

3 Proposed system

This research build applications that can help check the spelling after question script made. Differences in the form of manuscripts made such as essay, multiple choices or short answer will affect the application process that is built. The process of re-examining the question script is expected to be done after all the manuscripts have been typed, then the text will be re-examined by another teacher or the principal. The application that can be used to help re-examine the question script, and then give a sign on the misspelled words than the teacher will be easier in correcting the words.

Previous research has also been done to find out methods used for spelling correction. Previous research has produced spelling correction with Levenshtein Distance [9, 10]. In the previous research, system that made by Mawardi, Rudy and Naga [9], used Levenshtein Distance as the main method that combined with Trie data structure. The system can accept up to three to five words in the input box for spelling correction. While system that made by Mawardi, Susanto and Naga [10], also used Levenshtein Distance as the main method combined with Finite State Automata. The system can accept file in txt format for spelling correction and the result can be downloaded again in txt format.

With further development, this study used better version of Levenshtein Distance named Damerau-Levenshtein Distance. The Damerau-Levenshtein Distance add one more function to Levenshtein Distance to make it more powerful. In this research, the method will show up to 10 top words that will be chosen as the correction. The 10 top words was obtained from words that have the least value of Damerau-Levenshtein Distance. The value of Damerau-Levenshtein Distance was total operation required to change the Indonesian Dictionary that used and the words that want to check the spelling. The smaller the value, then the words are closer to the correct words.

The system that created in this study can accept input in the input box, upload exam script and the most important is can upload word file and get the same result in word file. The system can accept exam script such as multiple choice and essay. It will help typographical correction of test grade 1 to 6 easier and faster. The stages of the system in

this study are as follows:

- i. User will choose the file that consists of exam script test. The interface can be seen at Figure 1.
- ii. After upload the file, the text inside file will appear at correction text box.
- iii. User can choose manual correction or automatic correction. The manual correction will give user some of suggestion so they can choose the best correction. The automatic correction will give user automatically correction of the error word from the top of suggestion word.
- iv. At figure 2, at the correction result box, the annotation of typographical error word can be seen.
- v. If manual correction is selected, the suggestion word can be seen after pointing and clicking the word with blue annotation that can be seen at figure 3.
- vi. At Figure 4. the exam file test can be downloaded after re-examine the typographical error.

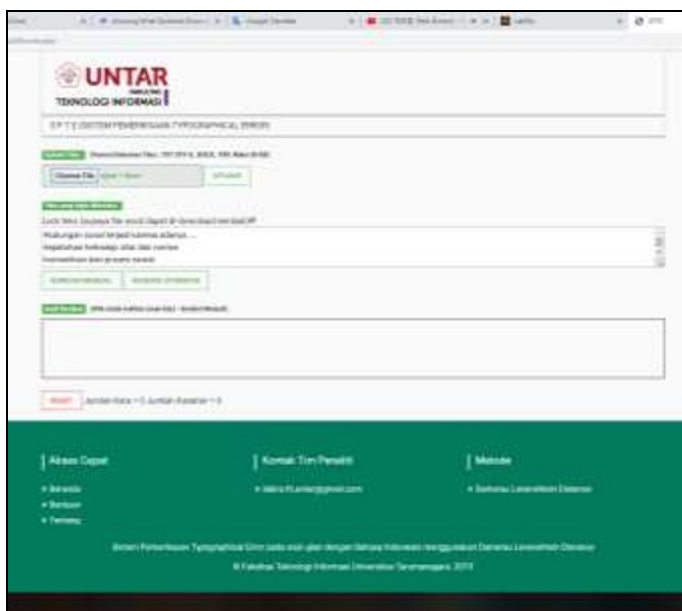


Fig. 1. Interface “Sistem pemeriksaan typographical error untuk soal ujian”.



Fig. 2. Interface the correction result box and the annotation of typographical error word.



Fig. 3. Interface suggestion word from manual correction menu.



Fig. 4. Interface download document.

4 Results and Discussion

Data used in this research are Indonesian language dictionaries, and question documents. The Indonesian language dictionary is used as a comparison in the method of Damerau-Levenshtein Distance. While the test text is used to test the results of the method that have

been made using Damerau-Levenshtein Distance. The Indonesian dictionary is obtained from the indodic.com website. On this website, the dictionary can be downloaded directly on the link provided. The words contained in this dictionary are 41 312 words. However, this dictionary does not cover all of the Indonesian words so it should be re-edited. This study conducted two type of test. The first test is test to detect non-real word errors from 50 sentences. This study created several non-real word errors from each sentence. Table 2 is the example of non-real word error sentence.

The second test is test for exam script document consist of 15 questions in three types format test. Each format test consists of five short answer, five essay and five multiple choices. This testing try to test the result of format document after manual and automatic correction. The format of document from typographical correction result script test must be same after correction done by system. Some non-real word error from question that can be seen at Table 3.

Table 2. Example of sentences with non-real word error.

No	Example of sentence with non-real word error
1	Salah satuny dengan menjaga kecepatan jalan kaki
2	Dari pojok ruangan, merembet ke rsang Hakim Agung Hamdan
3	Tokek Asia Tenggara Terancam Punah
4	Lebih paad sosialisasi kartu apa saja yang bisa dilihat
5	Meninggal sata di perjalanan ke rumah sakit

Table 3. Example of question with typographical error.

	Non-Real word Error
Deletion	Hubungan sosal terjadi karena adanya
Insertion	Katak berkembang biaak dengan cara
Substitution	Sebutkan fungsi rsngka bagi tubuh manusia!
Transposition	Sebutkan 3 pembgaian rangka tubuh manusia!

Some typographical error can be seen at Table 2 and Table 3, the system will detect and give correction. The system was created by using the method of Damerau-Levenshtein Distance in order to correct non-real word error get 88 % accuracy. Results from the system can be seen in Table 4. Manual correction and automatic correction have been tried. The manual method only displays 10 top word suggestions and the automatic method corrects words automatically with the first word suggestions.

Table 4. The Damerau-Levenshtein Distance test results by displaying the top 10 word suggestions.

Test Categories	Manual Correction	Automatic correction
Sentence Accuracy	88 %	70 %
Word Accuracy	84 %	71 %
Average Time Spent (s)	4.3	2.1

After testing done, the result from the automatic correction will choose the first top word suggestion. It can be seen from Figure 5. The first word suggestion for word “biaak” was “biak” so this word is correct. After manually choose the correct word suggestion the word change into the right word as it can be seen at Figure 6. The word “biaak” change to “biak”, “rsngka” change to “rangka” and word “pembgaian” change to “pembagian”.

But sometimes the automatic correction gives false suggestion. This is because the first word suggestion not always the correct word to correct the spelling in the automatic correction. At Figure 7 can be seen word “bagan-bagan” must change to “bagian-bagian”. If manual correction is selected, “bagian-bagian” can be chosen, but when automatic correction is selected the typographical correction will change to “badan-badan” because the first top of suggestion word is “badan-badan”.

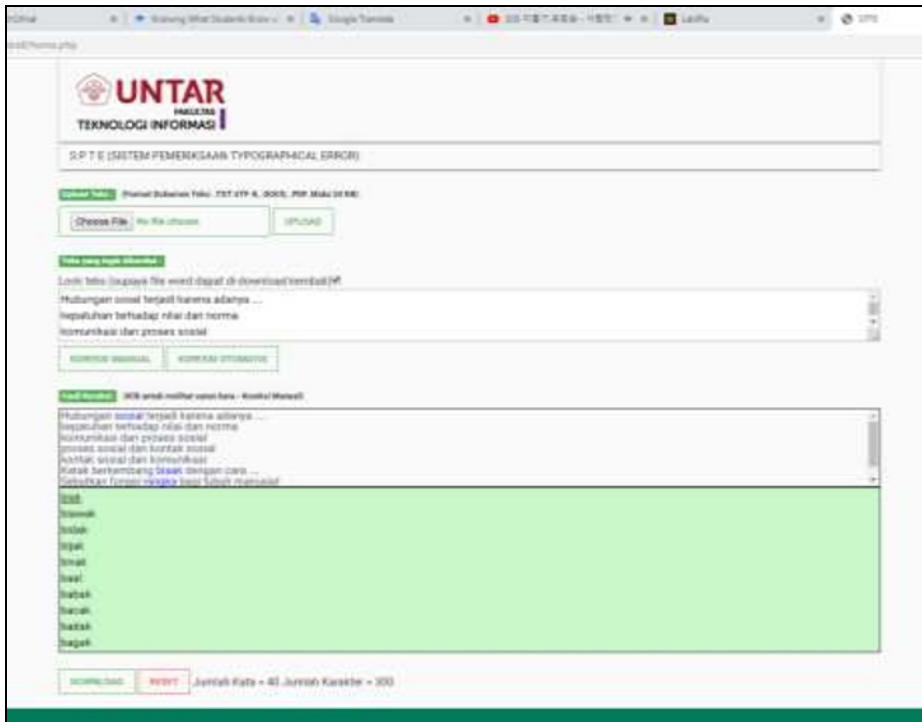


Fig. 5. Word “biak” appear at the first top word suggestion.

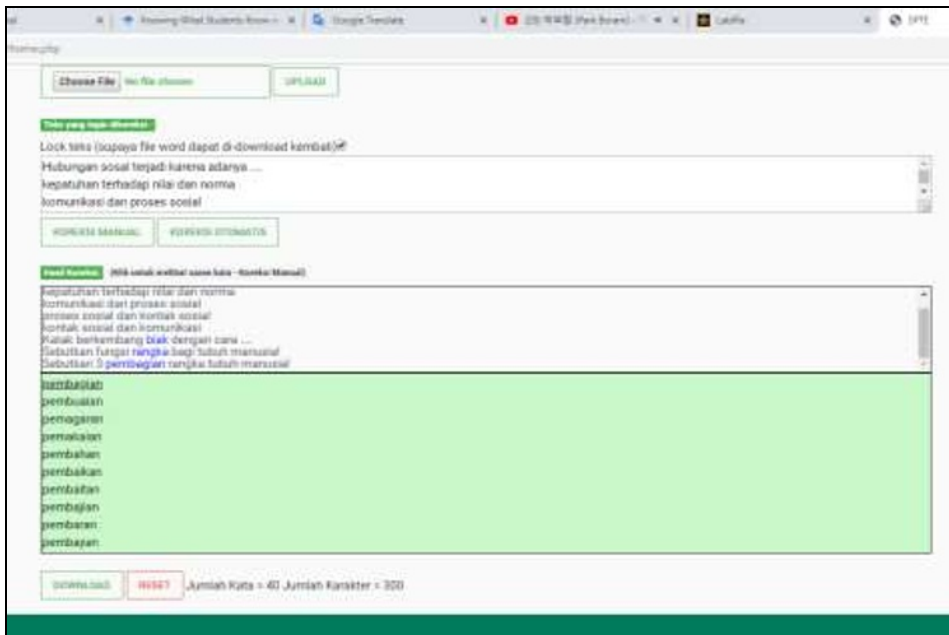


Fig. 6. Manual correction choose word from word suggestion.



Fig. 7. The correct word suggestion not appear as the first word suggestion.

5 Conclusions

The result of this research is a spelling correction application to help teachers typographical error in exam test scripts base on web. This application can correct non-real word error, have been able to produce good accuracy. Conclusions from this research are: (i) This application can detect typographical error word and give annotation for that word; (ii) This application can examine three types of exam script test in .docx format document; (iii) The results obtained by the system using Damerau-Levenshtein Distance with non-real word error, has word accuracy of 88 %; (iv) This application can help teacher to re-evaluate the exam test script and improve quality of exam test.

Suggestions for further improvement in this research are to reduce the time for processing. This can be done by applying finite state automata and the trie data structure, adding Indonesian language dictionary that have been validated by asking for expert help, and using another method that relevant with this research. The important suggestion is to improve this application so this application can accept a lot of type of document and save them to another type of document but keep all format of assessment that produce by teacher.

Thank you to the Ministry of Research, Technology and Higher Education of the Republic of Indonesia (Ristek Dikti) for funding this research based on contract number 225/SP2H/LT/DRPM/2019, 29/AKM/PNT/2019. Thanks to the teachers who have helped indirectly in providing information about how busy when they produce assessment test and difficulties of re-evaluate the script. Thanks also to the main researchers, student assistants who have helped this research. Big thanks to the research assistants who are always ready at any time at the time of this research.

References

1. R.A. Sani. *Inovasi pembelajaran*. Jakarta: Bumi Aksara (2013). [in Bahasa Indonesia]. https://www.researchgate.net/publication/320540340_INOVASI_PEMBELAJARAN
2. A.I. Fahma, I. Cholissodin, R.S. Perdana. *Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer*, **2**,1:53–62(2018). [in Bahasa Indonesia]. <http://j-ptiik.ub.ac.id/index.php/j-ptiik/article/view/690>
3. H. Faili, N. Ehsan, M. Montazery, M.T. Pilehvar. *Digit. Scholarsh. Humanit.* **31**,1:95–117(2014). <https://academic.oup.com/dsh/article-abstract/31/1/95/2605365>

4. M.M. Yulianto, R. Arifudin, A. Alamsyah. Scientific Journal of Informatics **5**,1:67–75(2018). <https://journal.unnes.ac.id/nju/index.php/sji/article/view/67>
5. G.D. Kuh, N. Jankowski, S.O. Ikenberry, J.L. Kinzie. *Knowing what students know and can do: The current state of student learning outcomes assessment in US colleges and universities*. Champaign, Illinois: National Institute for Learning Outcomes Assessment (2014). p. 51. <https://learningoutcomes.web.illinois.edu/wp-content/uploads/2019/02/2013SurveyReport.pdf>
6. L. Barari, B.Q. Zadeh. *CloniZER spell checker adaptive language independent spell checker*. Paper Presented in AIML 2005 Conference CICC (Cairo, Egypt, 2005). AIML, 19–21(2005). https://www.researchgate.net/publication/249876908_CloniZER_Spell_Checker_Adaptive_Language_Independent_Spell_Checker
7. C. Zhao, S. Sahni. BMC Bioinformatic, **20**,Suppl 11:19–46(2019). <https://link.springer.com/content/pdf/10.1186/s12859-019-2819-0.pdf>
8. V.I. Levenshtein. Soviet Physics Doklady, **10**,8:707–710(1966). <https://nymity.ch/sybilhunting/pdf/Levenshtein1966a.pdf>
9. V.C. Mawardi, Rudy, D.S. Naga. TELKOMNIKA **16**,2:827–833(2018). <http://journal.uad.ac.id/index.php/TELKOMNIKA/article/view/6890/4430>
10. V.C. Mawardi, N. Susanto, D.S. Naga. MATEC Web Conf., **164**,01047:1–16(2018). https://www.matec-conferences.org/articles/mateconf/abs/2018/23/mateconf_icesti2018_01047/mateconf_icesti2018_01047.html