

## TURET2.0: Thesis Writing Tutor Aimed on Lexical Richness in Students' Texts

Samuel González-López<sup>1</sup>, Aurelio López-López<sup>2</sup>, Jesús Miguel García-Gorrostieta<sup>2</sup>,  
Indelfonso Rodríguez Espinoza<sup>3</sup>

<sup>1</sup> Instituto Tecnológico de Nogales,  
Sonora, Mexico

<sup>2</sup> Instituto Nacional de Astrofísica, Óptica y Electrónica, Tonantzintla,  
Puebla, Mexico

<sup>3</sup> Universidad Tecnológica de Nogales,  
Sonora, Mexico

samuelgonzalezlopez@gmail.com, {alopez, jesusmiguelgarcia}@inaoep.mx,  
irodriguez@utnogales.edu.mx

**Abstract.** Writing thesis is a process of constant interaction between the student and academic advisor, but writing correctly is a complex task for students even if they have the support of a teacher. The elaboration of the thesis document requires the implementation of a methodology and procedures, which constitute the elements of the format and structure in the thesis. This work shows an intelligent tutoring system (TURET2.0) designed in a web platform and which provides a customized tutoring for students in drafting their writings, specifically to evaluate the lexical richness of seven sections of the thesis. Moreover, as a way to motivate students to achieve their goals, some gamification techniques were implemented. The measures used to assess the lexical richness are lexical variety, lexical density and sophistication.

**Keywords:** E-learning, natural language processing, intelligent tutoring system, lexical richness, gamification.

### 1 Introduction

Writing a thesis is not easy for undergraduate students and even more for academic reviewers, since the document requires several revisions to achieve the essential points stated in most institutional guidelines. This work aims to help undergraduate students to improve the document drafting in terms of Lexical Richness through a tutoring system. TURET2.0<sup>1</sup> includes two game attributes in order to motivate students to use the system. TURET is an updated version of a tutoring system tool previously [1].

---

<sup>1</sup> In Spanish: TURET: Tutor para la Redacción de Tesis.

One essential factor affecting writing is lexical competence, i.e. the writer ability to use properly the vocabulary, leading to considering it a basic reference point for measuring the quality of writing [2]. Universities in Canada take into account the results obtained by students in proficiency exams of different areas. One of them refers to the domain of English; other refers to the domain of mathematics. A study at the University of Calgary for Non-Native English Speaking (NNES) students, aimed to relate the academic success of students with lexical richness [3]. One of their research questions was to compare the lexical richness of NS (Native English Speaking) and NNES (Non Native English Speaking) students with their academic performance. The authors conclude that the results suggest that students with appropriate vocabulary, varied and accurate, have excelled in their studies, while students with a general vocabulary, repetitive, and an uncontrolled set of vocabulary showed a decreased academic performance. This conclusion supports our efforts aimed to improve the writing of students in their research drafts.

Advances in intelligent tutoring systems (ITS) include the use of natural language technologies to analyze student writing and provide feedback as presented in the article by McNamara [4]. Writing Pal (WPal) is an ITS that offers a strategy instruction, practice, and feedback for developing writers. There are also intelligent virtual agents, which are able to answer questions for the student related to an academic subject [5]. A dialogue-based ITS called Guru was proposed in [6], which has an animated tutor agent engaging the student in a collaborative conversation that references a hypermedia workspace, displaying and animating images significant to the conversation. Similarly, our work presented in this document includes the use of Natural Language, but adding two attributes of gamification.

The gamification approach could motivate students to get involved, focus and strive to engage in activities that seem boring, reaching a better performance. The main activities of the game include: information search, selection of information, strategy development, conflict resolution, decision-making exercises, and negotiation [7]. In the work of [8], an intelligent tutor for solving linear equations with elements of gamification is combined with a reward system. Students who used the tutor, were granted with a reward. In subsequent tests (when re-practicing problems), the performance was lower compared to students who did not obtain a reward. In contrast, students who solved new problems to re-practice their skills had better performance. Also, the authors conducted a comparison of the performance of students who used the tutor and a commercial tool. The students that used the commercial tool achieved a lower result in learning.

TURET2.0 is a tutoring system that seeks to support students close to graduating from universities with the need to write a thesis or research project. The document drafting is a difficult activity for the students, as this requires a methodology and procedures to comply properly with the structure that conforms the thesis. TURET2.0 differs from previous version because, it includes elements of gamification with the idea of maintaining student motivation. In addition, we evaluate seven sections of the thesis; in previous work only four sections were assessed. With this version, we seek to support students in the area of IT in Spanish language.

This tutor includes a module for assessing the lexical richness, which is done in terms of lexical density, lexical variety, and sophistication. There are a variety of methods to evaluate the use of vocabulary (lexicon) in text. One of them is to measure the sophistication of some papers using text word lists.

Our proposed system intends to assist the work of the instructor and to facilitate and guide students through this process. The paper is organized as follows. Section 2 describes the lexical richness model, while section 3 details the tutor with examples of draft evaluations. We conclude in section 4, discussing additionally further work.

## 2 Lexical Richness Model

To evaluate the seven elements contained in a thesis, we propose a computational model that will include three lexical dimensions. The first step in the model considers the preprocessing of each element. Each section in this module is processed with the Freeling<sup>2</sup> tool to obtain the word stems, converting the analyzed word in its singular form, grouping similar terms, and allowing a fast lexical analysis.

Another step in the preprocessing of the text was filtering and removing stop words from a list of 325 words provided by the Natural Language Toolkit (Snowball). Stop words include prepositions, conjunctions, articles, and pronouns. After this step, only content words remained, which allowed the calculation of the three dimensions.

**Table 1.** Measures to compute lexical richness.

Dimension descriptions		
Dimension	Labels	Computed as
Variety	LV	$Tlex/Nlex$
Density	LD	$Tlex/N$
Sophistication	LS	$NSlex/Nlex$
Tlex: Unique lexical terms		
Nlex: Total lexical terms		
Nslex: Words out of a list of common terms (SRA)		
N: Total tokens		

The first procedure is computing the lexical variety which seeks to measure student ability to write their ideas with a varied vocabulary. This function is calculated by dividing the unique lexical types (Tlex) between all lexical types (Nlex).

The second module refers to the computation of the lexical density, whose goal is to reflect the proportion of content words with respect to all the words employed, i.e. if the text has a good level of content. This dimension is obtained by dividing the unique lexical types or content words (Tlex) by the total words of the evaluated text (N) i.e. the number of words before removing stop words (see Table 1).

<sup>2</sup> <http://nlp.lsi.upc.edu/freeling/>

Finally, the sophistication method attempts to reveal the knowledge of the technical subject and it is estimated as the proportion of “advanced” words employed. This measure is computed as the percentage of words out of a list of common words (in our case, the 1000 common words, according to SRA).

Each of the measures takes values between 0 and 1, where 1 indicates an acceptable lexical value, and values close to zero mean a poor value of the lexicon of the evaluated section. Together, the three dimensions aim to identify the level of lexical richness of the student writing. The sophistication would be a plus for undergraduate students.

TURET2.0 uses the results computed by the Lexical Richness model to display them to the student, and adding feedback depending on the evaluation result.

### 3 TURET2.0

The results of a pilot test (prior version of tutor) with students of a public university showed positive results. Students who used the tutor had better results when writing their thesis (in terms of Lexical Richness) compared to those who did not use the tutor. The results were detailed in [1].

TURET2.0 was developed under the Python environment, the previous version used PHP and MySQL with XAMPP package to have web access. However, the response time was not as expected because calls were being made to the operating system to use the Freeling tool and Python from PHP.

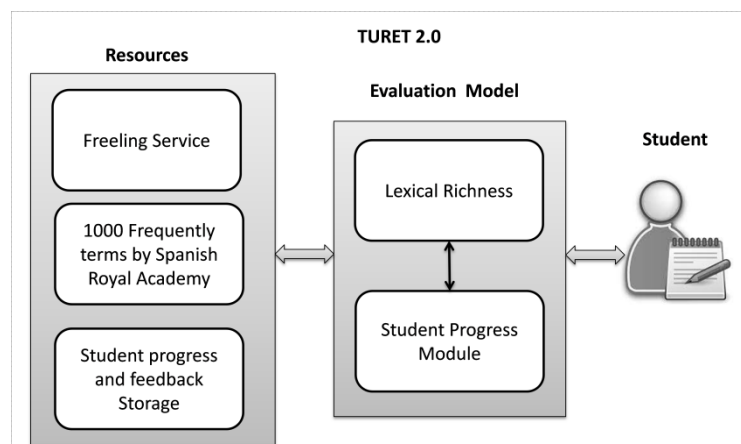


Fig. 1. TURET Scheme System.

Under the Python environment, a Web framework “Django + HTML5” was used to display the interface and results to the student. The use of this environment avoided writing files, system calls and allowed to work in the data memory. Similarly, the open-source relational database management system “MySQL” was used to store the results of each evaluation of students. Finally, Freeling tool was installed as a server, such that

the lemmatization process was performed under the scheme of services, i. e., when a student requests an evaluation in the tutor, the system uses the lemmatization service of Freeling. In Figure 1, we present the system scheme.

The results of the lexical analysis are sent to the Student Progress Module (SPM) to update the student knowledge state. SPM records the student progress in a network. When the student completes the exercises with the Lexical Analyzer, the corresponding node is updated and the SPM estimates the student progress for the parent node using the weights assigned to the measure in turn. Weights were assigned to each node based on instructor's experience. In Table 2, we show the percentages achieved by the student, in case he gets a high score in each of the lexical dimensions. It is worth mentioning that the percentages for each lexical dimension are divided into three equal parts (each part is equal to 1/3). For example, if a student gets high score in all three dimensions in the objective section, he will get  $1/3 + 1/3 + 1/3 = 1$ , which means he has reached 15% of progress. This 15% of progress is the total assigned by the tutor to the objective section.

**Table 2.** Progress percentages for each section of a thesis.

Thesis structure <sup>3</sup>	
Elements	%
Problem statement	15
Objective	15
Justification	15
Methodology	15
Hypothesis	12
Research Questions	13
Conclusion	15

In Figure 2, we can observe the student work environment in TURET2.0. The elements evaluated by the tutoring system are: hypothesis, justification, objectives, problem statement, research questions, methodology, and conclusion. In this section, the student can review his overall progress and observe a section with the overall results of the remaining students who also are using the tutoring system. The aim is that the student is interested in getting the top ranking, similar to a video game. After several iterations of evaluation of the text in the tutor, we expect that its lexical richness improves.

Also we can notice the progress in the objectives section with 68% achieved, since is the only section that the user has been assessed. In this screen, the user can review his last advance, where a progress bar is used to present this progress, this is a feature of games which indicates that for each stage there exists an advancement. Progress bars belong to the category of "Games tasks and challenges" [7].

In the tutor's home page, we provide a description of the Lexical Richness, with this the student can understand the results of the tutor. The levels used to determine the

<sup>3</sup> Suggested by the authors of research methodology books.

assessment are High, Medium and Low; this scale was defined based on the analysis of the corpus of thesis and research proposals [1]. The student can click on one of the sections to write or paste the desired text to evaluate. The tutoring system will produce the result of the analysis in the three dimensions: density, variety and sophistication.

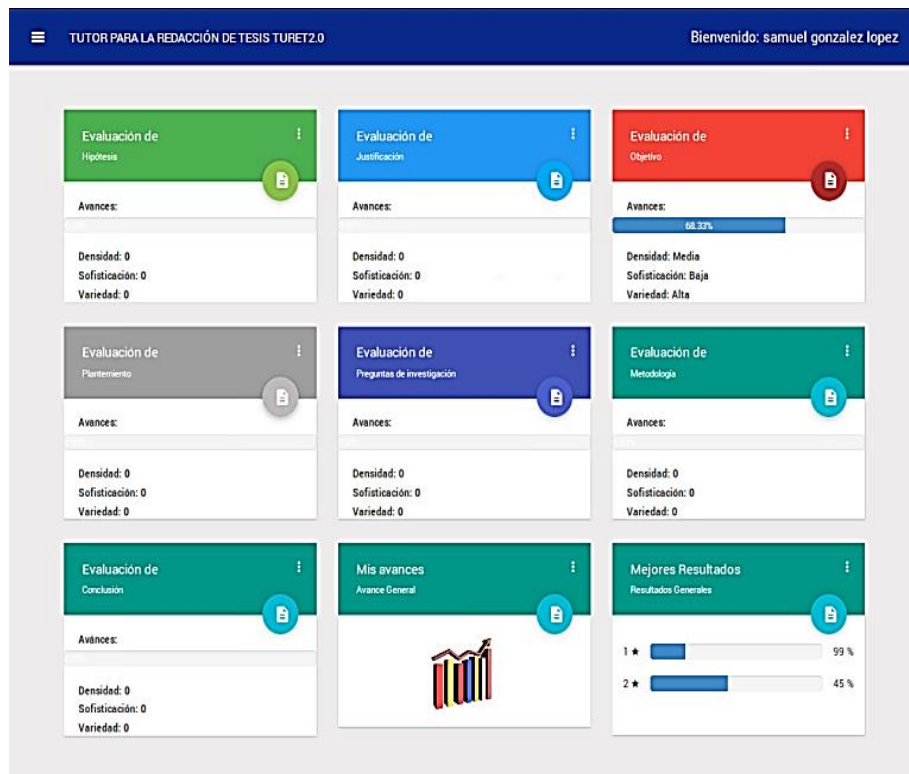


Fig. 2. Main evaluation sections in TURET2.0.

The section where the student can perform the text analysis is presented in Figure 3. Here we can also observe the feedback provided to the student. It can be noticed that two words are marked in red (e.g. “seguridad” in Spanish) in the section of variety assessment. This implies that a content word has been repeated which affects the level of variety.

Textual feedback is also provided, which for now is static. This feedback is de-fined depending on the level achieved by the student. In Figure 3, we observe a medium level of variety assessment, with a textual recommendation for the student to improve his writing (e.g. “Buen trabajo, pero aun nos falta corregir más nuestro texto” in Spanish).

In Figure 4, individual student progress is shown, globally depicting the level reached by the student in all three dimensions, i.e., the student can view the lexical richness of his entire thesis.



Fig. 3. Text evaluation sections in TURET2.0.

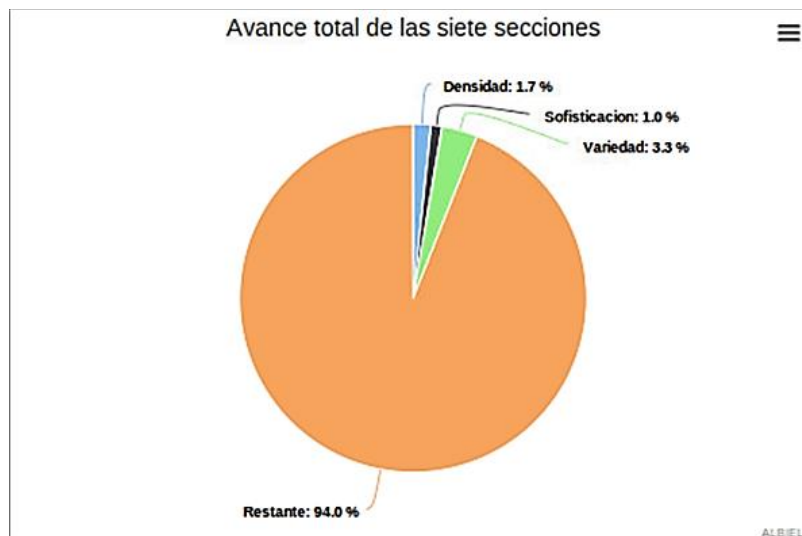


Fig. 4. Individual Global Advance.

Another attribute of gamification considered in TURET2.0 is the score among students using this tool. In Figure 5 we can observe the total score of each student using the tutoring system. The goal is to motivate a competition among them, but also cooperation. This attribute belongs to the category of “Games of collaboration and competition” [7].



Posición	Usuario	Correo	Total de Avances
1	Alex2612	alex@gmail.com	24 %
2	vbrz	vbrz@yahoo.com.mx	6.4 %
3	Carlos	jccp33@gmail.com	0 %

Ver 5 Resultados: 1 - 3 de 3

Fig. 5. Best results report.

## 4 Conclusions

The tool presented in this paper aims to support students to improve their writing in terms of lexical richness, with the possibility of improving the quality of the final document. With this, it would allow the academic advisor to focus on analyzing a higher proportion the content of the thesis rather than vocabulary or structure. TURET2.0 is a tool that aims to support to student and motivate to use it, that is, not just another tool to fulfill a requirement of writing.

We plan to gradually incorporate additional features to assess in the student texts such as coherence or argumentation, adhering to same idea of motivating its use.

In future work, as performed with the previous version of the tutoring system, we seek to pilot test it at different universities. In a first stage as a trial to analyze results and implement improvements, then in a second stage as a released tool.

It is also planned that TURET2.0 can be customized by the student to assess only the sections that he is required to write, since some universities do not ask for all sections, omitting for instance research question or hypotheses. Finally, we will take this tool to the mobile devices field for a higher coverage with students.

**Acknowledgements.** Second author was partially supported by SNI, México, while third author was supported by Conacyt México through scholarship 357381. We thank Jesús Carlos Cardenas Piñuelas and Erwin Alexander Villegas Tun for their collaboration as application developers.

## References

1. González-López, S., López-López, A.: Lexical analysis of student research drafts in computing. *Computer Applications in Engineering Education*, 23(4), pp. 638–644 (2015)
2. Grobe, C.: Syntactic Maturity, Mechanics, and Vocabulary as Predictors of Quality Ratings. *Research in the Teaching of English*, 15(1), pp. 75–85 (1981).
3. Douglas, S. R.: *Non-Native English Speaking Students at University: Lexical Richness and Academic Success*. Doctoral Thesis, University of Calgary, Canada (2010)
4. McNamara, D. S., Raine, R., Roscoe, R., Crossley, S., Jackson, G. T., Dai, J., Cai, Z., Renner, A., Brandon, R., Weston, J., Dempsey, K., Carney, D., Sullivan, S., Kim, L., Rus, V., Floyd, R., McCarthy, P. M., Graesser, A. C.: *The Writing-Pal: Natural language*



- algorithms to support intelligent tutoring on writing strategies. In: Applied natural language processing and content analysis: Identification, investigation, and resolution, Hershey, PA: IGI Global, pp. 298–311 (2012)
5. Rospide, C. G., Puente, C.: Virtual Agent Oriented to e-learning Processes. In: Proceedings of 2012 International Conference on Artificial Intelligence, Las Vegas, Nevada (2012)
  6. Olney, A., D'Mello, S. K., Person, N. K., Cade, W. L., Hays, P., Williams, C., Lehman, B., Graesser, A. C.: Guru: A Computer Tutor That Models Expert Human Tutors. In: Stefano A. Cerri, William J. Clancey, Giorgos Papadourakis & Kitty Panourgia (eds.) ITS, Springer, pp. 256–261 (2012)
  7. Marín, H., Alor, G., Zatarian, R., Barrón, L.: Una Revisión Sistemática de Técnicas de Gamification en Aplicaciones Educativas Inteligentes. In: Congreso Mexicano de Inteligencia Artificial, COMIA (2016)
  8. Yanjin-Long, Y., Alevan, V.: Gamification of Joint Student/System Control Over Problem Selection in a Linear Equation Tutor. In: Proceedings of the 12th International Conference on Intelligent Tutoring Systems, Springer, pp. 378–387 (2014)