

# Construcción de contenido para un Sistema Tutor Inteligente en idiomas: un estudio piloto con el corpus *OneStopEnglish*

## Content Construction for an Intelligent Tutor System in languages: a pilot study on the OneStopEnglish corpus

Adelina Escobar-Acevedo<sup>1</sup> , Josefina Guerrero-García<sup>1</sup> 

<sup>1</sup>Benemérita Universidad Autónoma de Puebla, Puebla, México.  
[adeesa32@gmail.com](mailto:adeesa32@gmail.com), [josefina.guerrero@correo.buap.mx](mailto:josefina.guerrero@correo.buap.mx)

(Recibido: 11 enero 2022; aceptado: 18 abril 2022; Publicado en Internet: 30 junio 2022)

**Resumen.** Durante la adquisición de un idioma extranjero, la lectura representa una de las oportunidades de acercamiento al lenguaje. Sin embargo, los textos inadecuados pueden desencadenar una experiencia contraproducente para un estudiante, por ello, en los cursos regulares, los docentes utilizan su experiencia o la de un equipo editorial para seleccionar las lecturas. En un sistema automático como en un Sistema Tutor Inteligente, es prioritario realizar recomendaciones adecuadas al perfil del alumno. No basta conocer el nivel de idioma del texto. El presente trabajo aplica herramientas para clasificar una muestra de textos extraídos del corpus *OneStopEnglish* conforme al Marco Común de Referencia Europeo, crea grupos temáticos con análisis semántico latente (LSA), y aplica tres métricas populares de lecturabilidad como un referente para recomendar textos a los estudiantes.

**Palabras clave:** Selección de materiales, Métricas de lecturabilidad, Comprensión lectora, Sistemas Tutores Inteligentes.

**Abstract.** During foreign language acquisition, reading represents one of the opportunities to get closer to the language. However, inappropriate texts can cause students to have a negative experience; thus, in regular courses, teachers use their experience or an editorial team to select the readings. In an automatic system, as in an Intelligent Tutor System, making recommendations appropriate to the student's profile is a priority. It is not enough to know the language level of the text. This work uses tools to classify a sample of texts from the OneStopEnglish corpus according to the Common European Framework of Reference for Languages. We create thematic groups based on Latent Semantic Analysis (LSA) and use three popular metrics of readability as a guide to suggest texts to students.

**Keywords:** Material recommendation, Readability metrics, Reading Comprehension, Intelligent Tutor Systems.

**Tipo de artículo:** Artículo de investigación.

## 1 Introducción

La lectura es una actividad compleja que involucra variables cognitivas, lingüísticas, textuales, socioculturales, entre otras. Sin embargo, en el ejercicio de la lectura no basta considerar únicamente las habilidades del lector sino también la complejidad del texto. La ley del mínimo esfuerzo establece que, en una comunicación, tanto el receptor como el emisor tratarán de mantener el mensaje lo menos ambiguo posible para evitar interacciones prolongadas (Zipf, 1949), lo que tiene un impacto en el vocabulario utilizado en los textos.

En el caso de la lectura en idioma extranjero, se utilizan las habilidades y conocimientos del idioma materno para comprender el idioma objetivo. Existe una relación entre la comprensión lectora y el conocimiento de vocabulario en el texto, del cual el estudiante debe conocer alrededor del 90% (Zarobe & Zarobe, 2019). A raíz de esto, se han generado listas etiquetadas de términos por niveles de idioma (Tejada et al., 2015) considerando por ejemplo el Marco Común Europeo de Referencia (CEFR) (Instituto Cervantes, 2002). La Universidad de Cambridge cuenta con listas en inglés británico y americano (Cambridge University Press, 2015).

*Text Inspector* (Bax, 2020) y *Text Analyzer* (roadtogrammar, 2021) son herramientas que permiten determinar el nivel de idioma de los textos utilizando el vocabulario que contienen, además proporcionan información adicional que puede resultar útil para los docentes al momento de seleccionar materiales para la lectura en el aula.

La lectura es una actividad compleja. Para motivar a los alumnos, los docentes suelen utilizar una variedad de estrategias. En la selección de material, uno de los recursos utilizados es proveer material de acuerdo con la edad e interés de los estudiantes. Sin bien parece simple, en ciertas poblaciones la variedad de temas de interés y niveles de idioma de los alumnos es vasta y difusa, como se observa en los grados universitarios. Como dificultad adicional, la mayoría de los repositorios están integrados con textos ajenos a otros niveles, de forma que cambiar el nivel implica cambiar totalmente de temas y materiales.

Una adaptación común en los Sistemas Tutores Inteligentes consiste en proporcionar textos más sencillos (otro texto de nivel inferior) o versiones simplificadas (el mismo texto de nivel inferior) a estudiantes que no han cumplido con el desempeño esperado. Contar con versiones simplificadas es útil también para plataformas de lectura (Fiction Express Education, 2021).

Los sistemas tutores inteligentes, en lenguaje, deben ser capaces de obtener la mayor cantidad de información tanto del texto como el usuario para poder realizar recomendaciones de lectura adecuadas al usuario. El objetivo de este artículo es presentar los resultados de métricas de lecturabilidad populares sobre uno de los grupos temáticos generados automáticamente en el corpus *OneStopEnglish*. El resto del trabajo está organizado de la siguiente manera: en la sección 2 se aborda la lecturabilidad y simplificación de textos. En la sección 3 se presentan algunas métricas para el análisis de textos. La sección 4 presenta brevemente la utilidad de usar LSA. La sección 5 presenta una revisión del trabajo de otros autores. La sección 6 contiene la evaluación realizada en este trabajo sobre uno de los grupos del corpus *OneStopEnglish*. Finalmente, en la sección 7 se presenta la conclusión.

## 2 Lecturabilidad y simplificación de textos

La simplificación de textos se aborda principalmente mediante simplificación léxica o sintáctica. La primera consiste en reemplazar palabras complejas por sinónimos más simples haciendo uso de recursos como diccionarios o las listas de términos etiquetados por nivel. La segunda consiste en simplificar las estructuras complejas preservando el significado original. Hasta hace poco las soluciones todavía estaban basadas en reglas utilizando etiquetadores y analizadores sintácticos (Al-Thanyyan & Azmi, 2021).

La creación de corpus simplificados sobre todo de forma manual es una tarea costosa, por lo que la mayoría de los corpus existentes no se encuentran libres. Se puede crear un corpus de forma automática con Wikipedia y Wikipedia simplificada (Wilkens et al., 2018). No obstante, Xu et al. (2015) determinan que sólo alrededor del 50% de los textos son verdaderas simplificaciones.

La lecturabilidad determina la relación entre el texto y el esfuerzo del lector para entenderlo, lo que involucra vocabulario, estilo, cohesión, entre otros factores. Existen varios tipos de escalas para asignar el grado. Uno de ellos se basa en el grado de escolaridad requerido por el lector para la comprensión del texto (Liu, 2020). Entre las herramientas para determinar la facilidad de lectura de un texto Coh-Metrix (Graesser et al., 2017) y Textcompare.org (2021) se encuentran disponibles en la web. Coh-Metrix actualmente provee 104 métricas de análisis (Graesser et al., 2004) y es utilizada para análisis de textos.

## 3 Métricas de lecturabilidad y características textuales

Las métricas de lecturabilidad se utilizan como predictores de comprensión a nivel oración y para determinar el tiempo necesario para leer el texto. Las métricas más populares, *Flesch Kincaid Grade Level* y *Flesch Reading Ease*, consideran características superficiales de los textos como el número de sílabas y la longitud de palabras y oraciones (McNamara et al., 2014).

### 3.1 Flesch Kincaid Grade Level

Es una métrica que se expresa en grados escolares conforme al sistema americano. Para ser aplicada adecuadamente se requieren documentos con al menos 200 palabras. La fórmula está definida por la ecuación (1).

$$RDFKGL = (.39 \times ASL) + (11.8 \times ASW) - 15.59 \quad (1)$$

Donde:

$$ASL = \left( \frac{\text{número de palabras}}{\text{número de oraciones}} \right)$$

$$ASW = \left( \frac{\text{número de sílabas}}{\text{número de palabras}} \right)$$

Un resumen de la interpretación de puntajes se aprecia en la [Tabla 1](#). La numeración es ascendente conforme a la complejidad del texto de 5 a 18.

**Tabla 1.** Interpretación de puntajes RDFKGL.

Puntaje	Grado escolar	Descripción
5-8	5to a 7mo grado	Fácil de leer
8-9	8vo y 9no grado	Inglés conversacional
10-12	10mo-12vo grado	Ligeramente difícil
13-15	Universitario	Difícil
16-18	Graduado universitario	Muy difícil
18+	Profesional	Extremadamente difícil de leer

### 3.2 Flesch Reading Ease

A diferencia de la métrica anterior, los puntajes más altos están relacionados con la facilidad del texto con rango de 0 a 100. La fórmula está representada en la ecuación (2).

$$RDFRE = 206.835 - (1.015 \times ASL) - (84.6 \times ASW) \quad (2)$$

La [Tabla 2](#) muestra la interpretación de los puntajes considerando también el promedio de sílabas por cada 100 palabras y la longitud de las oraciones.

**Tabla 2.** Interpretación de puntajes RDFRE.

Puntaje	Grado escolar	Descripción
90-100	5	Muy fácil
80-90	6	Fácil
70-80	7	Relativamente fácil
60-70	8-9	Estándar
50-60	10-12	Relativamente difícil
30-50	Universitario	Difícil
0-30	Graduado universitario	Muy difícil

A pesar de la gran variedad de métricas de lecturabilidad disponibles, estas dos son muy populares y se incluyen comúnmente como referencia sobre todo en entornos monolingües. Sin embargo, actualmente existen métricas más sofisticadas que involucran cohesión, complejidad sintáctica, entre otras

características para determinar el nivel de lecturabilidad específicamente del inglés como idioma extranjero, es el caso de RDL2.

### 3.3 RDL2

Esta métrica, a diferencia de las otras dos fórmulas, considera la superposición y la frecuencia de vocabulario, así como la similitud sintáctica (McNamara et al., 2014). La fórmula mostrada en la ecuación (3), intenta predecir la lecturabilidad para lectores de idioma extranjero.

$$RDL2 = -45.032 + (52.230 \times CRFCWO1) + (61.306 \times SYNSTRUT) + (22.205 \times WRDFRQmc) \quad (3)$$

Donde:

- CRFCWO1 es la proporción media de la superposición de vocabulario de contenido en oraciones adyacentes.
- SYNSTRUT es la proporción media de similitud sintáctica entre oraciones adyacentes.
- WRDFRQmc es el valor medio de la frecuencia mínima de palabras de contenido (en logaritmo), usando la base de datos CELEX.

La lista completa de métricas obtenidas por Coh-Metrix puede consultarse en la documentación de la herramienta.

## 4 Análisis Semántico Latente

*Latent Semantic Analysis* (LSA) es un modelo computacional que utiliza una técnica estadística llamada descomposición en valores singulares (DVS), lo que permite relacionar objetos que aparecen en los mismos espacios (Landauer et al., 2006). En el caso de texto, determinar que palabras se utilizan bajo los mismos contextos lo que conlleva distintas aplicaciones en Procesamiento Natural de Lenguaje. Naturalmente reduce la dimensionalidad en los modelos vectoriales, por lo que se utiliza para efectuar agrupamientos temáticos, clasificación, análisis de escritos, entre otros. Es particularmente útil en los Sistemas Tutores Inteligentes orientados a lenguaje. Las versiones de i-START (*Interactive Strategy Training for Active Reading and Thinking*) utilizan LSA para asignar un puntaje de forma automática al escrito generado por los estudiantes con respecto al texto original (Crossley et al., 2015; Allen et al., 2015; McCarthy et al., 2020).

## 5 Revisión del estado del arte

La información que proporciona Coh-Metrix se ha utilizado en diversos estudios. Uno de los usos es comparar textos e incluso libros completos. Zhang (2016) utilizó Coh-Metrix para analizar dos populares libros de texto para la enseñanza del idioma inglés en China. El objetivo era determinar cuál es más adecuado para estudiantes universitarios. Para esto, seleccionó veintiún medidas, incluyendo RDL2, así concluyó que uno de los libros es más básico en los textos que presenta y por lo tanto servía sólo para estudiantes de menor grado.

Cárcamo Morales (2020), incluye las tres métricas de lecturabilidad proporcionadas por Coh.Metrix en tres libros de texto de décimo, décimo primer, y décimo segundo grados utilizados para la enseñanza del inglés en Chile. No sólo analiza los textos sino también las preguntas relacionadas. Concluye que el programa no es progresivo.

Además de comparar libros de texto, también se han utilizado las métricas en experimentos de seguimiento ocular. Los estudios de movimientos oculares indican que, mientras más complejo es el texto, se incrementa el número de fijaciones y su duración. Nahatame (2020) pretende identificar cuál de las métricas de lecturabilidad es mejor predictor de dificultad de textos ante los estudiantes. Los experimentos

involucraron a cuarenta estudiantes japoneses de nivel universitario y concluye que las métricas son confiables para predecir las fijaciones; siendo RDL2 el mejor predictor.

En la producción escrita, Li et al. (2018) hace un análisis de complejidad en las explicaciones de 293 estudiantes de secundaria. Utiliza 19 valores arrojados por Coh-Metrix, incluyendo las métricas de lecturabilidad RDFKGL y RDL2. Sólo resultó útil como discriminante la primera.

Ninguno de los autores, utiliza una única métrica para la caracterización de un texto. En este sentido, la elección depende del objetivo del análisis. En el caso de la facilidad de lectura además se debe tener en cuenta el grado de narratividad, cohesión, y variedad de vocabulario.

## 6 Método

Este trabajo es parte de la creación de contenido para un Sistema Tutor Inteligente. Idealmente debe existir una gran variedad de materiales que puedan despertar el interés del alumno. En el caso de los textos, se propone considerar dos factores para recomendar al estudiante los más afines a su perfil: la temática y la lecturabilidad.

El Sistema Tutor Inteligente que se está creando tiene la intención de cubrir niveles A2, B1, y B2 conforme al marco europeo CEFR (de elemental a intermedio). Un texto puede ser categorizado en diferente número de niveles dependiendo del referente. Incluso, en algunos repositorios las clases pueden ser más difusas al estar relacionadas con el plan de estudios de casas editoriales y la libertad en el orden de los temas a enseñar. Por ello, además de buscar un corpus etiquetado, se hizo un análisis con métricas de lecturabilidad.

El corpus utilizado es el *OneStopEnglish* (Vajjala & Lucic, 2018), que consta de 189 textos en inglés, cada uno con tres versiones por nivel: avanzado, intermedio y elemental. Los textos fueron simplificados manualmente por docentes de inglés, tratando de conservar número de palabras y oraciones en lo posible. A la fecha es el único repositorio apto para niveles universitarios etiquetado por nivel que se ha encontrado libre para experimentación. El corpus fue creado con un doble objetivo, por una parte, apoyar la creación de modelos de simplificación y por otra la evaluación de lecturabilidad.

Las etiquetas originales no están alineadas con respecto a CEFR e IELTS, lo que se comprobó aplicando *Text Analyzer* en una muestra del corpus conformada por los primeros 17 documentos (ver Figura 1). Se observa que, aunque los documentos fueron etiquetados como elemental, la herramienta los ubica en nivel intermedio. Como observación, los documentos asignados a otro nivel como el documento 8 que es *intermediate* en CEFR, y *upper intermediate* en IELTS, reflejan los marcos de referencia diferentes.

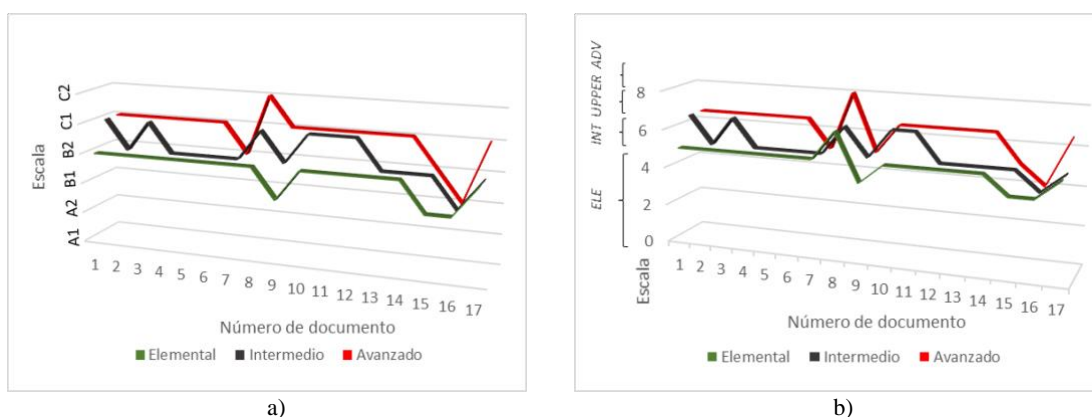


Figura 1. Niveles a) CEFR b) IELTS de los primeros 17 documentos del corpus *OneStopEnglish*.

Ya que el corpus no se encuentra etiquetado de forma temática, se utilizó LSA para forzar agrupamientos sobre el conjunto de archivos correspondiente al nivel avanzado. Para el preprocesamiento se retiraron signos de puntuación, números, y palabras vacías en inglés, se aplicó tokenización y *Porter stemmer* con NLTK. El modelo LSA se creó utilizando *gensim* con pesado *tfidf*. El diccionario del corpus está conformado con 10,762 tokens. En este estudio piloto se formaron 10 grupos aprovechando la

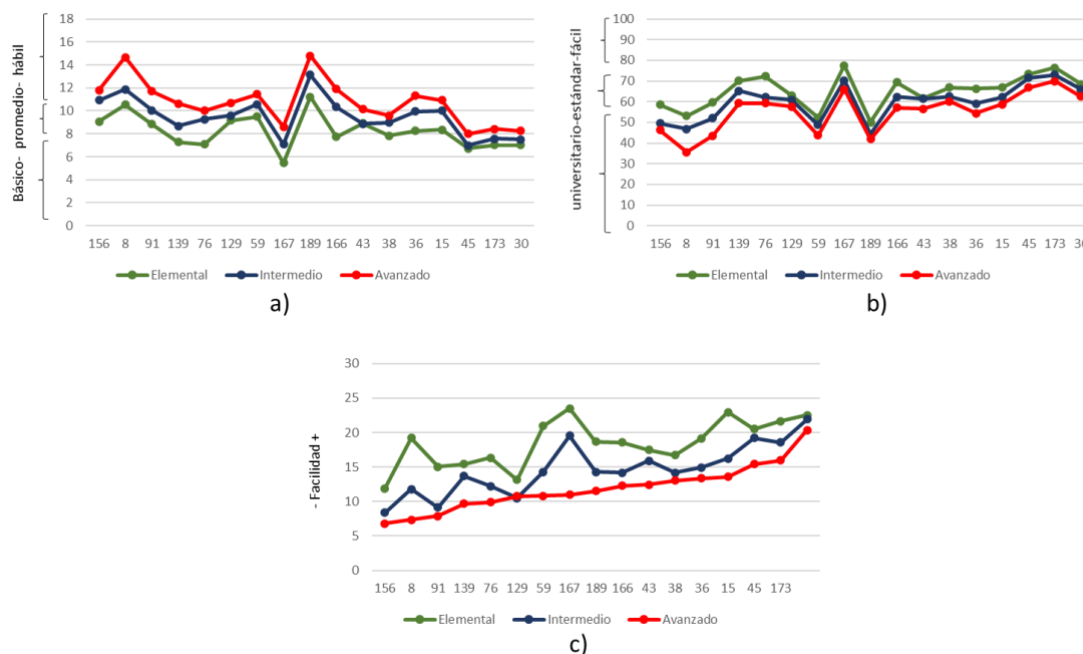
representación basada en temáticas (latentes, no explícitas). La distribución resultante puede apreciarse en la [Tabla 3](#).

**Tabla 3.** Distribución de documentos del corpus *OneStopEnglish* agrupado en 10 temáticas.

Clase	Cantidad de documentos
0	0
1	32
2	23
3	32
4	21
5	27
6	10
7	17
8	17
9	10

El objetivo de realizar esta agrupación temática es recomendar al usuario tanto documentos de su interés como de su nivel. A forma de ejemplo se consideró la clase 7 que consta de 17 documentos.

El segundo paso implica determinar las métricas de lecturabilidad para el conjunto de documentos. Se usaron las tres que proporciona Coh-metrix con las tres versiones de los textos (ver [Figura 2](#)). Se graficó el grupo ordenando de forma ascendente los documentos de nivel avanzado con respecto a la métrica RDL2. Por facilidad se utiliza como título la posición alfabética original del corpus *OneStopEnglish* en vez del título textual.



**Figura 2.** Métricas de lecturabilidad a) RDFKGL, b) RDFRE, c) RDL2.

Se seleccionaron 3 documentos para ver el detalle. El documento 167 cuyo título original es “*WNL The last*” trata del último rinoceronte blanco macho del planeta. Contiene vocabulario mayormente descriptivo y es una muestra del comportamiento esperado. Se observa la simplificación progresiva de forma homogénea en todas las columnas de la [Tabla 4](#).

**Tabla 4.** Detalle de resultados de métricas.

Número de documento	Etiquetado original	RDFKGL (5-18)	RDFRE (0-100)	RDL2 (0-30)*
167	Avanzado	8.60	65.99	10.95
	Intermedio	7.10	70.28	19.52
	Elemental	5.46	77.56	23.51
30	Avanzado	8.27	62.54	20.36
	Intermedio	7.50	65.93	21.92
	Elemental	7.02	68.44	22.51
129	Avanzado	10.698	57.584	10.764
	Intermedio	9.599	60.928	10.508
	Elemental	9.178	63.103	13.163
* Rango en el corpus completo				

En el documento 30 “*Japan*”, cuenta con vocabulario orientado a la tecnología. La variación entre versiones es pequeña, por ejemplo, en RDFKGL es de 0.48 puntos entre la versión intermedia y la elemental. Por otra parte, el documento 129 “*Mystery shopper*” que trata sobre la profesión, en las simplificaciones los puntajes reflejados no son consecutivos para RDL2. Sin embargo, se determina la relación de las métricas conforme al etiquetado original.

## 7 Conclusiones

Con este trabajo se identifica que el etiquetado inicial por niveles de lectura es insuficiente, se requiere un análisis del texto que al menos incluya métricas de legibilidad. Otros datos importantes que pueden considerarse son: narratividad, uso de pronombres, variedad léxica, entre otros, para indicar si el texto se acopla a los objetivos lingüísticos perseguidos, antes de presentarlo a los estudiantes.

Otros autores ya han probado experimentalmente con alumnos la efectividad de hacer análisis más profundos de los textos para predecir la dificultad y el tiempo de la lectura. Este proceso es útil para que un sistema de recomendación sea capaz de tomar decisiones informadas sobre el texto más orientados para el usuario. Incluir la agrupación temática pretende motivar la lectura presentando materiales del interés del usuario. Incorporar material variado e identificar su pertinencia es particularmente importante en la formación del módulo dominio para los Sistemas Tutores Inteligentes en el área de lenguaje.

Como trabajo futuro se considera definir grupos con temáticas transparentes al usuario aplicando otros métodos de agrupamiento y la clasificación de los ejercicios relacionados con los textos para que el estudiante alcance los objetivos lingüísticos esperados.

## Declaración de conflicto de intereses

Los autores declaran no tener conflicto de intereses con respecto a la investigación, autoría o publicación de este artículo.

## Financiación

Apoyo otorgado por el CONACYT con la beca 235249.

## ORCID iD

Adelina Escobar-Acevedo  <https://orcid.org/0000-0003-0574-0932>

Josefina Guerrero-García  <https://orcid.org/0000-0002-3393-610X>

## Referencias

- Al-Thanyyan, S. S., & Azmi, A. M. (2021). “Automated Text Simplification: A survey”. *ACM Computing Surveys*, 54(2), 1–36.
- Allen, L. K., Snow, E. L., & McNamara, D. S. (2015). Are you reading my mind? Modeling students’ reading comprehension skills with natural language processing techniques. *ACM International Conference Proceeding Series, 16-20-Marc*, 246–254. <https://doi.org/10.1145/2723576.2723617>
- Bax, S. (2020). *Text Inspector*. <https://textinspector.com/>
- Cambridge University Press. (2015). *English Profile, The CEFR for English*. <https://www.englishprofile.org/wordlists/evp>
- Cárcamo Morales, B. (2020). “Readability and types of questions in Chilean EFL high school textbooks”. *TESOL Journal*, 11(2), 1–15.
- Crossley, S., Allen, L. K., Snow, E. L., & McNamara, D. S. (2015). Pssst... textual features... there is more to automatic essay scoring than just you! *Proceedings of the Fifth International Conference on Learning Analytics And Knowledge - LAK '15*, 203–207. <https://doi.org/10.1145/2723576.2723595>
- Fiction Express Education. (2021). *Fiction Express*. <https://en.fictionexpress.com>
- Graesser, A. C., McNamara, D. S., Louwse, M. M., & Cai, Z. (2004). “Coh-Metrix : Analysis of text on cohesion and language”. *Behavior Research Methods, Instruments, & Computers*, 36(2), 193–202.
- Graesser, A. C., McNamara, D. S., & Louwse, M. M. (2017). *Coh-Metrix*. <http://cohmetrix.com/>
- Instituto Cervantes (2002). “Marco Común Europeo de Referencia para las Lenguas: aprendizaje, enseñanza, evaluación”. *Instituto Cervantes*.
- Landauer, T. K., McNamara, D. S., Dennis, S., & Kintsch, W. (2006). *Handbook of Latent Semantic Analysis* (Vol. 7, Issue 2). Routledge.
- Li, H., Gobert, J., Dickler, R., & Morad, N. (2018). “Students’ Academic Language Use When Constructing Scientific Explanations in an Intelligent Tutoring System”. *Conference on Artificial Intelligence in Education*, 267–281. [https://doi.org/10.1007/978-3-319-93843-1\\_20](https://doi.org/10.1007/978-3-319-93843-1_20)
- Liu, Y. (2020). *Assessing text readability and quality with language models* [Master Thesis]. University of Helsinki.
- McCarthy, K. S., Watanabe, M., Dai, J., & McNamara, D. S. (2020). Personalized learning in iSTART: Past modifications and future design. *Journal of Research on Technology in Education*, 52(3), 301–321. <https://doi.org/10.1080/15391523.2020.1716201>
- McNamara, D. S., Graesser, A. C., McCarthy, P. M., & Cai, Z. (2014). Coh-Metrix Measures of Text Readability and Easability. En *Automated Evaluation of Text and Discourse with Coh-Metrix* (pp. 78–95). Cambridge University Press.
- Nahatame, S. (2020). Text readability and comprehension processes during L2 reading: A computational and eye-tracking investigation. *Conference of the American Association for Applied Linguistics (AAAL)*.
- Roadtogrammar (2021). *Text Analyzer*. <http://www.roadtogrammar.com/textanalysis/>
- Tejada, M. Á. Z., Gallardo, C. N., Ferradá, M. C. M., & López, M. I. C. (2015). Building a Corpus of 2L English for Automatic Assessment: The CLEC Corpus. *Procedia - Social and Behavioral Sciences*, 198(Cile), 515–525. <https://doi.org/10.1016/j.sbspro.2015.07.474>
- Textcompare.org. (2021). *Textcompare.org*. <https://www.textcompare.org/readability/>
- Vajjala, S., & Lucic, I. (2018). OneStopEnglish corpus: A new corpus for automatic readability assessment and text simplification. *Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Educational Applications*, 297–304. <https://www.aclweb.org/anthology/W18-0535/>
- Wilkins, R., Zilio, L., & Fairon, C. (2018). SW4ALL: a CEFR-Classified and Aligned Corpus for Language Learning. *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, 365–370.
- Xu, W., Callison-Burch, C., & Napoles, C. (2015). Problems in Current Text Simplification Research: New Data Can Help. *Transactions of the Association for Computational Linguistics*, 3, 283–297.
- Zarobe, Y. R. De, & Zarobe, L. R. De (Eds.). (2019). *La lectura en lengua extranjera*. Ediciones Octaedro.
- Zhang, R. (2016). A Coh-Metrix Analysis of Two Textbooks: Successful English for Vocational Colleges and Vocational College English (An Integrated Skills Course). *US-China Foreign Language*, 14(5), 351–356.
- Zipf, G. K. (1949). Introduction and Orientation. *Human behavior and the principle of least effort: an introduction to human ecology*. Addison-Wesley Press.