

Arabic Text Readability Analysis

Project Description:

Motivation: Readability scores are a way to measure whether a written text is likely to be understood by the intended reader. Text that scores a readability grade level of 8 or below will be readable by around 85% of the general public. This can help in improving language skills for language learners, proficiency tests, text simplification, designing of educational materials, etc.

Available Resources: Arabic curricula for Gulf countries (Grades 1 to 6), Automatic processing (Farasa NLP tools) for diacritization, stemming, parsing, etc. huge Arabic corpora from Wikipedia, Aljazeera and other sources.

We already built an initial interface for spotting complex words that didn't appear in a selected grade.

FAVES	GRADE	ISSUES	REACH	WORDS
☆	B	20	100%	304

Readability Grade Levels	
Flesch-Kincaid Grade Level	7.9
Gunning Fog Index	10.3

Readability Scores	
Flesch Reading Ease	66.8

Text Statistics	
Word Count	304
Sentence Count	19
Paragraph Count	6

Language Issues	
Spelling Issues	0.0%

Duties/Activities:

- We will study different readability formulas and apply to Arabic.
- We will build a website to provide readability analysis for an input text and give tips to improve its readability.

Required Skills:

- Java/Python
- Web development
- Understanding Arabic is a plus

Preferred Intern Academic Level:

- Junior/Senior levels

Learning Opportunities:

- This internship will help interns to immerse in coding and using various NLP libraries. It will help as well to get familiar with web development and build full stack profile. Also, you about recent research in reading analysis and apply them in real scenarios.

Expected Team Size: *it is preferable to have team projects*

- 2 to 4

Mentors

Name: Hamdy Mubarak email: hmubarak@hbku.edu.qa
Name: Ahmed Abdelali email: aabdelali@hbku.edu.qa
Name: Kareem Darwish email: kdarwish@hbku.edu.qa

Arabic Text Readability Analysis

Mentors: Hamdy Mubarak, Ahmed Abdelali, Kareem Darwish {hmubarak, aabdelali, kdarwish}@hbku.edu.qa

Motivation: Readability scores are a way to measure whether a written text is likely to be understood by the intended reader. Text that scores a readability grade level of 8 or below will be readable by around 85% of the general public. This can help in improving language skills for **language learners**, **proficiency tests**, **text simplification**, designing of **educational materials**, etc.

Available Resources: Arabic curricula for Gulf countries (Grades 1 to 6), Automatic processing (Farasa NLP tools) for **diacritization**, **stemming**, **parsing**, etc. huge Arabic corpora from **Wikipedia**, **Aljazeera** and other sources.

We already built an initial interface for spotting **complex words** that didn't appear in a selected grade.

- ❑ We will study different **readability formulas** and apply to Arabic.
- ❑ We will build a **website** to provide **readability analysis** for an input text and give tips to improve its readability.

❖ Skills

- ✓ Experience in Java/Python for data processing, calling APIs...
- ✓ Experience in web development: HTML, CSS
- ✓ Passionate about language analysis and education
- ✓ Understanding Arabic is a plus!

Enter a text:

إجراء: الأجراء، أكثر جراً في المغامرة والسفر بدأ الأغنياء، القيام برحلات تحويلية تستند إلى فلسفة اكتشاف الذات. وفي الوقت نفسه، يؤكد الخبراء أنه كلما كان الشخص أكثر ثراء، فكل قضا، إجازات مثيرة.

Grade: Grade 3

SUBMIT

Result:

خبراء: الأجراء أكثر جراً في المغامرة والسفر بدأ الأغنياء القيام برحلات تحويلية تستند إلى فلسفة اكتشاف الذات. وفي الوقت نفسه، يؤكد الخبراء أنه كلما كان الشخص أكثر ثراء، فضل قضاء إجازات مثيرة.

Readability Analysis



* Show possible word forms of a given lemma (stem) with pronunciation (TTS) and examples of usage:

Ex: term "reading" in Qatar (Grade 6)

11	القراءة	12	قراءة	16	قراءة	23	القراءة	43	القراءة	130	قراءة	قراءة
1	القراءة	2	قراءة	3	قراءة	6	قراءة	6	القراءة			
1	قراءة	1	بقراءة	1	القراءة	1	القراءة	1	القراءة			
						1	للقراءة	1	قراءتك			

جَلَسَا لِقِرَاءَةِ الْقُرْآنِ بِتَدْبِيرٍ



- * Show **statistics** about each grade in each country: #words, #lemmas, development with time...
- * Show top terms that are **common** across all countries (for a certain grade), and terms appear in one country but not in the others.
- * Augment text data with **images** (top results from Google image search)
- * Determine **grade level** for any input text, **highlight** complex words, and suggest **simpler** ones.

اللغة السريانية.. تمتد آلاف السنين وتبقى الأثر

السريانية هي عمق التاريخ الذي لا يمكن لشخص أو فئة -صغرت أو كبرت- التجرد منها، فهي نسب من نوع آخر، فموسيقى أجدية حروفها عزف حطات تمتد آلاف السنين.

Grade 6+

- ❖ **Skills**
- ✓ Experience in Java/Python for data processing, calling APIs...
- ✓ Experience in web development: HTML, CSS
- ✓ Passionate about language and culture studies, and preparing resources for language learners (native and non-native speakers)
- ✓ Understanding Arabic is a plus!