



مركز زاي لبحوث اللغة العربية
ZAI Arabic Language
Research Center

ZUZAI

مركز أبوظبي
للغة العربية
Abu Dhabi Arabic
Language Centre



جامعة نيويورك أبوظبي

 NYU | ABU DHABI

 مختبر كامل
CAMEL Lab

RLRL 2023 *Spotlight Talk*
5 September 2023

بارق: المكنز العربي المتوازن لتقييم الانقرائية

BAREC: The Balanced Arabic Readability Evaluation Corpus

Nizar Habash

New York University Abu Dhabi

Hanada Taha

Zayed University

BAREC Project



- The project is a collaboration between **New York University Abu Dhabi's Camel Lab** and **Zayed University's ZAI Arabic Language Research Center**.
- The 2.5 years effort is funded by the Abu Dhabi Arabic Language Center (ALC), starting September 2023.
- The overarching objective of the BAREC project is to develop a comprehensive reference resource to facilitate the study and evaluation of Arabic readability across the Arab world.
- <http://barec.camel-lab.com/>



مجموعة القراءة المتدرجة

هذه برنامج صرح بإشراف معهد تعليم القراءة البريطاني، ووفق معايير، وشروط. تتألف مجموعة القراءة المتدرجة من ثمانية مستويات (أربع مراحل)، لكل منها لون يحدد درجة صعوبته من حيث مضمون القصة، ومستواها اللغوي، وعدد كلماتها، وتنوع مفرداتها. جدول الأنواع المرفقة يوضح تدرج مستوى القراءة في هذه المجموعة. صُنفت هذه الكتب ضمن مجموعات تهدف إلى تعليم اللغة العربية بالقراءة المتدرجة.



SAMER Project Lexicon (Al Khalil & Habash): 36,000 lexical entries

Level	Grade	Age	Examples
Level I	Grade 1	6	بيْت، شَجَرَة، أُرْزُق، كَبِير، صَنَعَ، أَكَل، فَرِحَ، عَلَى، لَكِن house, tree, rabbit, blue, big, to make, to eat, to be happy, on, but
Level II	Grades 2-3	7-8	جَزِيرَة، ذَهَب، سَنَة، دَاكِن، أَشْطَوَانِي، صَنَب، خَدَعَ، كَانَأ، قُرْب، إِذَا island, gold, year, dark, cylindrical, difficult, to cheat, to reward, near, if
Level III	Grade 4-5	9-10	رَنَة، مُتَّخَف، مُعَادَلَة، مُمَكِن، مُوَحَّد، أَغْرَى، نَدَّر، لَدَى، كَي، مَا إِنْ...حَتَى lung, museum, equation, possible, united, to entice, to be rare, with, for, no sooner... than...
Level IV	Grades 6-8	11-14	إِقْتِصَاد، نُسُج، طَمَأْنِينَة، رَاقِي، مُثَبَّت، نَكَتْ، أَغْضَى، إِتَان، إِتْمَا، لَنْ economy, sap, tranquility, sophisticated, proven, to breach, to overlook, during, whereas, if (were)
Level V	Specialist	15 -	أَدْمَة، قَنْطَرَة، هَيْضَة، مَظْيَاف، لَوْذَع، شُعْبِي، لَحَا، طَمَعَن، لَدُنْ، أُنَى epidermis, catheterization, cholera, spectroscope, witty, bronchial, to denounce, to depart, with (≈ chez in French), wherever

دليل عربي 21

دليل إلكتروني للكتب المصنفة
وفق معايير تصنيف عربي 21
وهنادا طه



تم تصنيف أكثر
من 9,000 كتاب!
9,000+ books
classified!

دليل عربي 21

دليل إلكتروني للكتب المصنفة
وفق معايير تصنيف عربي 21
وهنادا طه



تم تصنيف أكثر
من ٩٠٠٠ كتاب!
9,000+ books
classified!

	Hanada Taha Level	Grade	ACTFL	SAMER	Age
1	أ	1	مبتدئ أدنى	I	6
2	ب	1	مبتدئ أدنى	I	6
3	ج	1	مبتدئ متوسط	I	6
4	د	1-2	مبتدئ متوسط	I-II	6-7
5	هـ	1-2	مبتدئ أعلى	I-II	6-7
6	و	1-2	مبتدئ أعلى	I-II	6-7
7	ز	2	متوسط أدنى	II	7
8	ح	2-3	متوسط أدنى	II	7-8
9	ط	2-3	متوسط أوسط	II	7-8
10	ي	3	متوسط أوسط	II	7-8
11	ك	3-4	متوسط أعلى	II-III	8-9
12	ل	4	متقدم أدنى	III	9
13	م	4-5	متقدم أوسط	III	9-10
14	ن	4-6	متقدم أعلى	III-IV	10-11
15	س	6-8	متقن أدنى	IV	11-13
16	ع	6-8	متقن أوسط	IV	11-13
17	ف	7-8	متقن أعلى	IV	12-14
18	ص	9	متفوق	V	15
19	ق	10-12	متميز	V	16-18



نورة والكورونا

قصة للأطفال

تأليف: سارة عبدالله الشامي

رسوم: أسماء الربيش



SAMER Readability Analysis

Analyze Readability

Doc Level Word Level Clear

Modify Markup Show Hide Minimize Delete

Doc Level Analysis

Current Doc %

Names	0%
Level 1	80%
Level 2	10%
Level 3	6.7%
Level 4	0%
Level 5	3.3%
Unknown	0%

Total word count: 30

Target Level 5

At Level 5 | 1 words | 3.3%

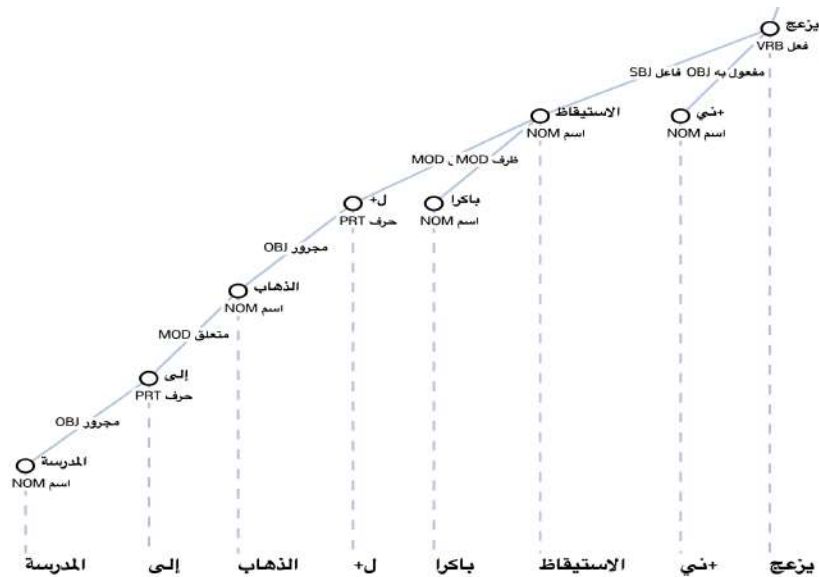
Below Level 5 | 29 words | 96.7%

By Token By Type

جامعة نيويورك أبوظبي



https://github.com/CAMEL-Lab/camel_tools



Camel Treebank Style Analysis

كامليرا: المحلل الآلي متعدد اللهجات للغة العربية

إدخال

إرسال

الجملة

يزعجني كثيراً النوم والاستيقاظ باكراً للذهاب إلى المدرسة، أما العطلة الصيفية فهي أحد ألامي الوردي؛ أحب النوم والاسترخاء، ولا أحب البتة الاعتماد على نفسي، وأنا دائماً التذمر من كل شيء.

الشكل: اسم الكلمة المفرد المعجمي

الجملة: العربية الفصحى

يزعجني كثيراً النوم والاستيقاظ باكراً للذهاب إلى المدرسة، أما العطلة الصيفية فهي أحد ألامي الوردي؛ أحب النوم والاسترخاء، ولا أحب البتة الاعتماد على نفسي، وأنا دائماً التذمر من كل شيء.

الجملة	الشكل	اسم الكلمة المفرد المعجمي
يزعجني	فعل	يزعج
كثيراً	اسم	كثير
النوم	اسم	نوم
والاستيقاظ	حرف	والاستيقاظ
باكراً	اسم	باكراً
للذهاب	حرف	للذهاب
إلى	حرف	إلى
المدرسة	اسم	المدرسة
،	حرف	،
أما	حرف	أما
العطلة	اسم	العطلة
الصيفية	اسم	الصيفية
فهي	حرف	فهي
أحد	حرف	أحد
ألامي	اسم	ألامي
الوردي	اسم	الوردي
؛	حرف	؛
أحب	فعل	أحب
النوم	اسم	النوم
والاسترخاء	حرف	والاسترخاء
،	حرف	،
ولا	حرف	ولا
أحب	فعل	أحب
البتة	اسم	البتة
الاعتماد	اسم	الاعتماد
على	حرف	على
نفسي	اسم	نفسي
،	حرف	،
وأنا	حرف	وأنا
دائماً	اسم	دائماً
التذمر	اسم	التذمر
من	حرف	من
كل	حرف	كل
شيء	اسم	شيء
.	حرف	.

الجملة	الشكل	اسم الكلمة المفرد المعجمي
المدرسة	اسم	المدرسة
إلى	حرف	إلى
الذهاب	حرف	الذهاب
+	حرف	+
باكراً	اسم	باكراً
الاستيقاظ	حرف	الاستيقاظ
+	حرف	+
يزعج	فعل	يزعج

<http://camelira.camel-lab.com/>

BAREC Planned Deliverables

- A **10M word corpus** of text excerpts distributed across readability levels, genres, topics, and text origins.
 - Fair-use of text samples.
 - Simulate the wide range of read materials in the world.
- **Human annotations** for sub-samples of the 10M corpus
 - 1M words - sentence segmentations & sentence level readability leveling in the Taha system.
 - Guidelines to map the system from book-level to sentence-level.
 - 100K words – syntactic annotations in CATiB dependency tree style.
 - 40K lemmas – annotating a morphological lexicon with Taha’s 19 levels.
- **Open-source AI tools**
 - Models for automatic readability using the annotated data.
 - Automatically annotate the 10M corpus.
 - Automatic readability reporting on new texts to help authors make informed decisions.

Q+A



- As this is a new project, may we ask you for...
 - Any feedback on the plan?
 - Pointers to resources?
 - Possible collaborations?
 - General advice from previous and similar experiences?
- Any questions you may have?

