



مركز زاي لبحوث اللغة العربية
ZAI Arabic Language
Research Center
ZUZAI

مركز أبوظبي
لغة العربية
Abu Dhabi Arabic
Language Centre



Linguistic Annotation Workshop (LAW-XIX-2025)

Guidelines for Fine-grained Sentence-level Arabic Readability Annotation

**Nizar Habash,[†] Hanada Taha-Thomure,[‡] Khalid N. Elmadani,[†]
Zeina Zeino,[‡] Abdallah Abushmaes^{††}**

[†]New York University Abu Dhabi

[‡]Zayed University ^{††}Abu Dhabi Arabic Language Centre

Introduction

- Readability refers to how easy or difficult to read and understand a piece of text.

Arabic	Translation
كُرَّة	Ball
غُرْفَةُ النَّوْمِ	The bedroom
سُلُوكِي مَسْئُولِيَّتِي	My behavior is my responsibility
كانت الحديقة واسعة، تطل على شاطئ النيل،	The garden was spacious, overlooking the Nile shore.
تعريف أصول الفقه	Definition of the Principles of Islamic Jurisprudence
بين طعن القنا وحقق البُنُود	Between the thrusts of lances and the fluttering of ensigns

Introduction

- Readability refers to how easy or difficult to read and understand a piece of text.

Grade	Arabic	Translation
KG	كُرَّة	Ball
1st	غُرْفَةُ النَّوْمِ	The bedroom
2nd	سُلُوكِي مَسْئُولِيَّتِي	My behavior is my responsibility
4th	كانت الحديقة واسعة، تطل على شاطئ النيل،	The garden was spacious, overlooking the Nile shore.
8th	تعريف أصول الفقه	Definition of the Principles of Islamic Jurisprudence
Uni	بين طعن القنا وحقق البُنُود	Between the thrusts of lances and the fluttering of ensigns

Introduction

- Readability refers to how easy or difficult to read and understand a piece of text.

RL	Grade	Arabic	Translation
1	KG	كُرَّة	Ball
3	1st	غُرْفَةُ النَّوْمِ	The bedroom
6	2nd	سُلُوكِي مَسْئُولِيَّتِي	My behavior is my responsibility
10	4th	كانت الحديقة واسعة، تطل على شاطئ النيل،	The garden was spacious, overlooking the Nile shore.
14	8th	تعريف أصول الفقه	Definition of the Principles of Islamic Jurisprudence
17	Uni	بين طعن القنا وحقق البُنُود	Between the thrusts of lances and the fluttering of ensigns

Introduction

- **Why Readability?**
 - Helps aligning texts with students' reading abilities.
 - Impacts understanding, retention, reading speed and engagement.
 - Developing readability models is crucial for improving literacy, language learning, and academic performance.
- **BAREC: the Balanced Arabic Readability Evaluation Corpus**
 - Over 69K sentence (1M words) annotated into **19 readability** levels inspired by the Taha/Arabi21 book-leveling approach.

Why 19 levels?

Taha-Thomure (2017)'s 19-level Arabic text leveling framework

- Inspired by Fountas and Pinnell (2006)'s 26 levels.
- 11 of the 19 levels covering up to 4th grade to support teachers in matching books to students' reading abilities.
- Ten qualitative and quantitative criteria: from text genre, and sentence structure, to book production quality, and word count.
- The Arab Thought Foundation adopted this framework and funded the leveling of over 9,000 children's books (**Arabi21**)
- *We mapped and restricted the features to apply on sentences.*

Arabic Orthographic Ambiguity

- Arabic script uses optional diacritical marks
 - 1.5% of newspaper words have some diacritical marks
 - Standard Arabic has 6.8 diacritizations and 2.7 lemmas/word

وَلَعَيْنَ

وَلَعَيْنَ وَلَعَيْنَ وَلَعَيْنَ

Infatuated (m.pl) # and for an eye/spring # and cursed

Arabic Morphological Richness

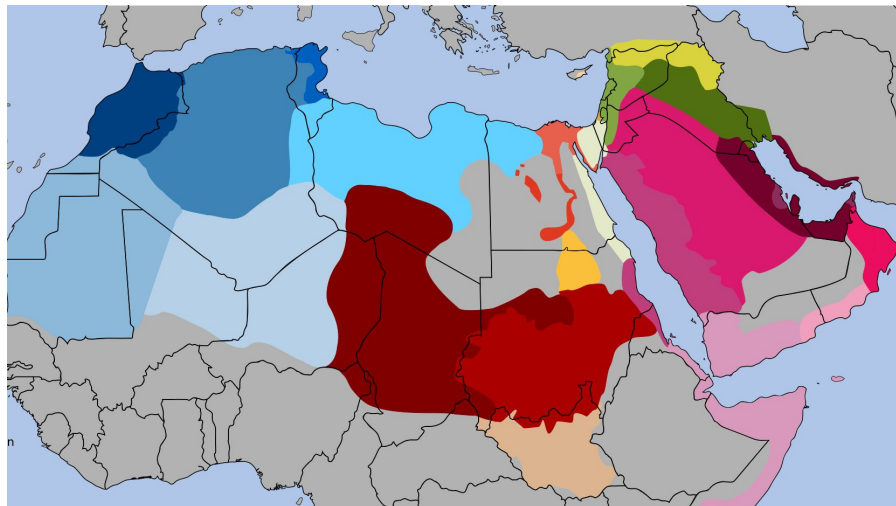
- Arabic has a very rich inflectional+cliticizational system
 - For example, Arabic verbs have 5,400+ forms
 - Whereas English verbs have 6!

وسنقولها
/wasanaqūluhā/
و + س + ن + ق + و + ل + ه + ا
wa+sa+na+qūl+u+hā
and+will+we+say+it
And we will say it

قال، قالت، قالوا، قلت،
قلت، قلتما، قلتهم، قلتن،
يقول، يقول، يقل، تقول، تقول،
تقل، تقولين، تقولي،
... فقال، فقالت، فقالا ...
... وسأقولها، وسنقولها، ...

Arabic Dialectal Variation

- Classical Arabic
 - Quranic Arabic, Historical texts
- Modern Standard Arabic
 - Official language
 - Language of news & media
 - Standard writing & grammar
 - The National Language
- Dialectal Arabic
 - Predominantly spoken
 - No official standardization
 - The Mother Tongue
 - Increasing use on social media



Annotation Guidelines

- **Textual Features:** Levels are assessed in six key dimensions.
 - a. **Spelling:** Word length & syllable count.
 - b. **Word Count:** The number of unique words.
 - c. **Morphology:** Inflection & derivation.
 - d. **Syntax:** Sentence structure & syntactic relations.
 - e. **Vocabulary:** The complexity of word choices.
 - f. **Content:** The required prior knowledge and abstraction levels.

BAREC Pyramid

The BAREC Pyramid illustrates the relationship across BAREC levels and linguistic dimensions, three collapsed variants (3 levels, 5 levels and 7 levels), and educational grades.

						Specialist	V	3-3	5-5	7-7	ق qaf-19			
					Uni 3 + 4	ص sad-18								
					Uni 1 + 2	ف fa-17								
					12	ع ayn-16								
						10-11	IV		5-4	7-6	س sin-15			
						8-9					ن nun-14			
						6-7	III	3-2	5-3	7-5	م mim-13			
						5					ل lam-12			
						4	II	3-1	5-2	7-4	ي ya-10	ك kaf-11		
						3				7-3	ح ha-8	ط ta-9		
						2	I		5-1	7-2	ه ha-5	و waw-6	ز zay-7	
						KG+1				7-1	أ alif-1	ب ba-2	ج jim-3	د dal-4
Spelling	Word Count	Morphology	Syntax	Vocabulary	Content	Grades	SAMER	BAREC-3	BAREC-5	BAREC-7	BAREC-19 Levels			

BAREC Annotation Guidelines

مستوى بارق	صف	ACTFL	عدد كلمات	تهجئة وإملاء	تصنيف واشتقاق	تراكيب نحوية	مفردات	فكرة ومحتوى
أ	روضة-1	مبتدئ أدنى	1	• كلمات من مقطع واحد أو مقطعين	• الفعل المضارع المفرد	• كلمة واحدة	• اسم جنس • اسم علم (متداول بسيط تركيبيا) • ضمير منفصل • مفردات متطابقة مع العامة - سامر I • الأرقام (العربية أو الهندية) 1-10	• فكرة مباشرة وصريحة وحسية. • لا رمزية في النص.
ب		مبتدئ أدنى	≤2	• كلمات من 3 مقاطع		• جملة اسمية (هو يلعب) • إضافة حقيقية (باب البيت) • صفة وموصوف (باب كبير)	• فعل • صفة • مفردات متشابهة مع العامة - سامر I • العدد الأصلي بالأحرف • الأسماء الخمسة: أبو، أخو	
ج	1	مبتدئ متوسط	≤4	• كلمات من 3 مقاطع	• سوابق: ال التعريف • سوابق: واو العطف • لواحق: ضمير المتكلم المفرد المتصل	• بدل كل: (صديقي أحمد) • بدل إشارة: (هذا البيت)	• مفردات فصيحة شائعة - سامر I • اسم الإشارة المفرد • الأرقام (العربية أو الهندية) 10-100	
د		مبتدئ متوسط	≤6	• كلمات تستخدم مذ الألف (أ)	• الفعل المضارع الجمع • سوابق: حروف جر متصلة • ظرف متون	• جملة فعلية بدون مفعول به • جار ومجرور	• حروف الجر	
هـ		مبتدئ أعلى	≤8	• كلمات من 4 مقاطع	• لواحق: ضمير متصل مفرد أو جمع • المثنى (في الأسماء والصفات) • جمع المؤنث السالم	• جملة فعلية مع مفعول به واحد اسم • جمل معطوفة • أدوات استفهام أساسية: ماذا، متى، من، أين، ما، كيف • صيغة التعجب "ما أفعل"	• العدد الترتيبي • الأرقام (العربية أو الهندية) 101-1,000 • اسم إشارة مثنى، جمع	• المحتوى من حياة القارئ. • لا رمزية في النص.
و	2	مبتدئ أعلى	≤9	• كلمات من 5 مقاطع	• الفعل الماضي المفرد والجمع • جمع مذكر سالم	• جملة فيها فعلين (مثلا جملة فعلية مفعولها أن المصدرية)	• مفردات فصيحة - سامر I	
ز		متوسط أدنى	≤10	• كلمات من 6+ مقاطع • أفعال/أسماء معتلة الآخر	• الفعل الماضي المثنى • الفعل المضارع المثنى • فعل الأمر المفرد • لواحق: ضمير المثنى المتصل • جمع التكسير • واو القسم (والله)	• مفعول فيه (ظروف زمان ومكان) • حال • أداة الاستفهام هل	• مفردات فصيحة شائعة - سامر II	• بعض الرمزية أو عدم التصريح المباشر بكل المقصود في الجملة

BAREC Annotation Guidelines

BAREC Level	Grade	ACTFL	Word Count	Spelling/Pronunciation	Morphology	Syntax	Vocabulary	Idea/Content
1-alif	Pre1-1	Novice Low	1	• One-syllable and two-syllable words	• Singular imperfective verb	• One word	• Common noun • Proper noun (frequent and simple) • Personal pronouns (non-clitics) • Vocabulary identical to dialectal form - SAMER I • Numbers (Arabic or Indo-Arabic) 1-10	<ul style="list-style-type: none"> • Direct, explicit, and concrete idea. • No symbolism in the text.
2-ba	1	Novice Low	≤2	• Three-syllable words			• Verb • Adjective • Vocabulary similar to dialectal form - SAMER I • Spelled cardinal numbers • The five nouns: <i>Abw</i> (father), <i>Axw</i> (brother)	
3-jim		Novice Mid	≤4		• Prtocolitic: Definite article <i>Al</i> + • Proclitic: Conjunction <i>wa</i> + • Enclitic: First Person Singular pronoun	• Apposition (full) • Demonstratives	• Common MSA vocabulary - SAMER I • Singular demonstrative pronoun • Numbers: 11-100	
4-dal		Novice Mid	≤6		• Plural imperfective verb • Prepositional proclitics • Nunated adverbials	• Verbal sentence w/o direct object • Preposition and object	• Prepositions	
5-ha	2	Novice High	≤8	• Four-syllable words	• Enclitic: Singular and Plural pronouns • Dual (in nouns and adjectives) • Sound feminine plural	• Verbal sentence with one nominal direct object • Conjoined sentences • Basic interrogative particles: what, when, who, where, how • Exclamatory form: how <comparative adjective>	• Ordinal numbers • Numbers: 101-1,000 • Dual and plural demonstrative pronoun	<ul style="list-style-type: none"> • Content is from the reader's life. • No symbolism in the text.
6-waw		Novice High	≤9	• Five-syllable words	• Singular and plural perfective verb • Sound masculine plural	• Sentence with two verbs (e.g., a verbal sentence a clausal direct object introduced with <i>Masdar 'an</i> [-to/that])	• MSA vocabulary - SAMER I	

BAREC Annotation Guidelines

BAREC Level	Grade	ACTFL	Word Count	Spelling/Pronunciation	Morphology	Syntax	Vocabulary	Idea/Content
1-alif	Pre1-1	Novice Low	1	• One-syllable and two-syllable words	• Singular imperfective verb	• One word	<ul style="list-style-type: none"> • Common noun • Proper noun (frequent and simple) • Personal pronouns (non-clitics) • Vocabulary identical to dialectal form - SAMER I • Numbers (Arabic or Indo-Arabic) 1-10 	<ul style="list-style-type: none"> • Direct, explicit, and concrete idea. • No symbolism in the text.
2-ba	1	Novice Low	≤2	• Three-syllable words			<ul style="list-style-type: none"> • Verb • Adjective • Vocabulary similar to dialectal form - SAMER I • Spelled cardinal numbers • The five nouns: <i>Abw</i> (father), <i>Axw</i> (brother) 	
3-jim		Novice Mid	≤4		<ul style="list-style-type: none"> • Prtcolitic: Definite article <i>Al</i>+ • Proclitic: Conjunction <i>wa</i>+ • Enclitic: First Person Singular pronoun 	<ul style="list-style-type: none"> • Apposition (full) • Demonstratives 	<ul style="list-style-type: none"> • Common MSA vocabulary - SAMER I • Singular demonstrative pronoun • Numbers: 11-100 	
4-dal		Novice Mid	≤6	• Words with an elongated Alif (e.g. /Zāsiṭ/)	<ul style="list-style-type: none"> • Plural imperfective verb • Prepositional proclitics • Nunated adverbials 	<ul style="list-style-type: none"> • Verbal sentence w/o direct object • Preposition and object 	• Prepositions	
5-ha	2	Novice High	≤8	• Four-syllable words	<ul style="list-style-type: none"> • Enclitic: Singular and Plural pronouns • Dual (in nouns and adjectives) • Sound feminine plural 	<ul style="list-style-type: none"> • Verbal sentence with one nominal direct object • Conjoined sentences • Basic interrogative particles: what, when, who, where, how • Exclamatory form: how <comparative adjective> 	<ul style="list-style-type: none"> • Ordinal numbers • Numbers: 101-1,000 • Dual and plural demonstrative pronoun 	<ul style="list-style-type: none"> • Content is from the reader's life. • No symbolism in the text.
6-waw		Novice High	≤9	• Five-syllable words	<ul style="list-style-type: none"> • Singular and plural perfective verb • Sound masculine plural 	<ul style="list-style-type: none"> • Sentence with two verbs (e.g., a verbal sentence a clausal direct object introduced with <i>Masdar</i> 'an [-to/that]) 	• MSA vocabulary - SAMER I	

BAREC Annotation Guidelines

BAREC Level	Grade	ACTFL	Word Count	Spelling/Pronunciation	Morphology	Syntax	Vocabulary	Idea/Content
7-zay		Intermediate Low	≤10	<ul style="list-style-type: none"> • Six-syllable or more words • Verbs/nouns with weak final letters 	<ul style="list-style-type: none"> • Dual perfective verb • Dual imperfective verb • Singular imperative verb • Enclitics: dual pronoun • Broken plurals • Waw of oath 	<ul style="list-style-type: none"> • Adverbial accusative (time and place adverbs) • Circumstantial accusative • Interrogative particle <i>hal</i> 	<ul style="list-style-type: none"> • High frequency MSA vocabulary - SAMER II 	<ul style="list-style-type: none"> • Some symbolism, or not everything is stated directly in the sentence.
8-ha	3	Intermediate Low	≤11		<ul style="list-style-type: none"> • Plural imperative verb • Feminine plural suffix (<i>nun</i>) in nouns and verbs • Other proclitics: future <i>sa+</i>, continuation <i>wa+</i>, conjunction <i>fa+</i> • Conjunctions (e.g., then, until, or, whether, but, as for) 	<ul style="list-style-type: none"> • Absolute object (emphasizing the verb) • Object of purpose • Object of accompaniment • Verbal sentence with two direct objects 	<ul style="list-style-type: none"> • MSA vocabulary - SAMER I and II • Negation particles • Numbers: 1,001-1,000,000 	<ul style="list-style-type: none"> • Some symbolism that requires the reader to seek help to understand the idea.
9-ta		Intermediate Mid	≤12		<ul style="list-style-type: none"> • Dual imperative verb • Interrogative Hamza • Ba of oath • Oath: The particle of oath, the object of the oath, and the answer to the oat 	<ul style="list-style-type: none"> • Vocative 	<ul style="list-style-type: none"> • Vocabulary describing positive and negative emotional and mood states like joy, happiness, anger, regret, sorrow 	<ul style="list-style-type: none"> • Some symbolism at the event level in the sentence that the reader understands through prior knowledge.
10-ya	4	Intermediate Mid	≤15		<ul style="list-style-type: none"> • Passive voice 	<ul style="list-style-type: none"> • <i>Inna</i> and its sisters (particles introducing a subject) • <i>Kana</i> and its sisters (past tense verbs) • Preposed predicate, postponed subject • Chain of narration • <i>rubba</i> preposition construction • Relative clauses • Circumstantial and object clauses 	<ul style="list-style-type: none"> • Singular relative pronouns • Verbal particles <i>qad</i> and <i>laqad</i> • Preposition-Conjunctions: <i>mimma</i>, <i>fima</i>... 	
11-kaf		Intermediate High	≤20		<ul style="list-style-type: none"> • Acting derivatives (e.g., the active participle) 	<ul style="list-style-type: none"> • Nominal sentence with a nominal predicate • False idafa (tall in stature) 	<ul style="list-style-type: none"> • Dual and plural relative pronouns 	<ul style="list-style-type: none"> • A degree of symbolism and a need for prior knowledge to understand the meaning of the sentence.
12-lam	5	Advanced Low			<ul style="list-style-type: none"> • Diminutive form 	<ul style="list-style-type: none"> • Parentheticals (explanation, blessing) • Exception • Exclusivity • Apposition (e.g., partitive or containing) • Specification (<i>tamiyiz</i> construction) 	<ul style="list-style-type: none"> • MSA vocabulary - Samer III • Frozen Verbs (e.g., <i>Amiyn</i> Amen) • Numbers: > 1,000,000 • Five Nouns: Dhu (possession nominal) • Interjections: <i>bala</i>, <i>Ajal</i>, etc. 	

BAREC Annotation Guidelines

BAREC Level	Grade	ACTFL	Word Count	Spelling/Pronunciation	Morphology	Syntax	Vocabulary	Idea/Content
13-mim	6-7	Advanced Mid			<ul style="list-style-type: none">• Energetic mood (emphatic <i>nun</i>)• Ta of oath	<ul style="list-style-type: none">• Conditional sentences• Jussive particle <i>lamma</i> (not yet)	<ul style="list-style-type: none">• Words describing deep psychological states like depression, loss, psychological alertness• Use of coined, uncommon words• Abbreviations (e.g., LLC)	<ul style="list-style-type: none">• Symbolic ideas and deeper meanings, especially in terms of the psychological dimension of characters/events.• Local cultural expressions that may not be understood by those outside the
14-nun	8-9	Advanced High				<ul style="list-style-type: none">• Semantic emphasis• Praise and dispraise• <i>Masdar 'an</i> clause as a subject• Exclamatory form: <comparative adjective> <i>bih min</i>	<ul style="list-style-type: none">• MSA vocabulary - SAMER IV• General legal, scientific, religious, political vocabulary, etc.• Five Nouns: <i>fw</i>, <i>Hmw</i>	
15-sin	10-11	Superior Low				<ul style="list-style-type: none">• Uncommon constructions that are ambiguous and need diacritization for clarification	<ul style="list-style-type: none">• Specialized vocabulary that requires understanding the concept/idea to comprehend it• Shortening in proper names (e.g., <i>fatim</i> for <i>fatima</i>)	
16-ayn	12	Superior Mid					<ul style="list-style-type: none">• MSA vocabulary - SAMER V• Specialized and highly elevated Arabic vocabulary.• Vocabulary mostly distant from dialects.	
17-fa	University Year 1-2	Superior High					<ul style="list-style-type: none">• Scientific and heritage vocabulary not in use today, but familiar to a novice specialist	
18-sad	University Year 3-4	Distinguished					<ul style="list-style-type: none">• Scientific and heritage vocabulary not in use today, but familiar to a specialist	
19-qaf	Specialist	Distinguished+					<ul style="list-style-type: none">• Scientific and heritage vocabulary not in use today, but familiar to the advanced researcher specialist	
Difficulty	This tag is used when there is difficulty in assessing the level. It is preferred to use this tag so that the team can find a solution (for example, by adjusting the criteria or adding explanatory details).							
Problem	Generally, we use this tag for sentences containing:	<ul style="list-style-type: none">• Spelling mistakes (e.g., Hamzas, Ta Marbuta, Alif maqsura/Ya)• Errors in diacritics• Linguistic awkwardness (illiteracy, colloquialism, poor translation from a foreign language)• Inappropriate topics (racism, bias, bullying, pornography, etc.)• Sentences and phrases mostly written in languages other than Arabic or in non-Arabic script				However, in the following cases, we provide the level and add a note in the comments column: <ul style="list-style-type: none">• Error in Hamzat al-Wasl/Hamzat al-Qat' >> (إ)• Offensive words >> (ع)• Error in diacritics at the beginning of the sentence >> (ن)• Dotted Yaa missing at the end of the word >> (ي)		

Annotation Guidelines

- **Textual Features:** Levels are assessed in six key dimensions.
 - a. **Spelling:** Word length & syllable count.
 - b. **Word Count:** The number of unique words.
 - c. **Morphology:** Inflection & derivation.
 - d. **Syntax:** Sentence structure & syntactic relations.
 - e. **Vocabulary:** The complexity of word choices.
 - f. **Content:** The required prior knowledge and abstraction levels.

RL	Arabic Sentence/Phrase	Translation	Reasoning
1-alif	أرنب <u>Rabbit</u>		One bisyllabic familiar noun
2-ba	ملعب واسع <u>A large playground</u>		Noun-adjective
3-jim	أنا أحب اللون الأحمر. I love <u>the</u> color red.		Definite article
4-dal	الشمس تشرق في الصباح الباكر. The sun rises early <u>in the morning</u> .		Prepositional phrase
5-ha-	القطّة تستريح على السرير وتستمع بأشعة الشمس الدافئة. <u>sunshine</u> .	The cat rests on the bed <u>and enjoys the warm</u>	A conjoined sentence
6-waw	سلوكي <u>مُسؤوليتي</u> My behavior is <u>my responsibility</u>		Five syllable word

Annotation Process

- **Annotation Team:** Five annotators, all of whom are experienced Arabic language educators with experience in Arabi21 project.
- They receive folders containing annotation sets for individual annotation.
- Each set contains 100 randomly selected sentences.
- An average of 3 hours per set (~2min/sentence).
- **Result:** 1,039,371 words - 69,441 sentences - 1,922 documents.
- 19 shared sets for measuring inter-annotator agreement (IAA).

Annotation Team

	A0 ^P	A1	A2	A3	A4	A5 ^L
Native Language	Arabic	Arabic	Arabic	Arabic	Arabic	Arabic
Other Language	En, Fr	En	En, Fr	En, Fr	En, Fr	En, Fr
Nationality	Syrian	Lebanese	Lebanese	Lebanese	Lebanese	Lebanese
Residence	USA	Lebanon	Lebanon	Lebanon	UAE	Lebanon
Gender	Female	Female	Female	Female	Female	Female
Background	Muslim	Muslim	Muslim	Muslim	Christian	Muslim
Degree	MA	BA	BA	MA	MA	B MA
Major	Applied Ling.	Arabic Lit.	Geography	Arabic Lit.	Arabic Lit.	Arabic Lit.
Experience	CT, LA, RA	PT, LA	PT, LA	CT, LA	CT, LA	CT, LA, RA
School	Private	-	-	Public&Private	Private	Public
Level	University	Elementary	Elementary	Secondary	Secondary	Secondary
Students	L2	L1	L1	L1	L1	L1
Years	16	16	22	22	8	25

Annotation Interface

Sentence/Phrase	Length	Level	Word Count	Spelling/Pronunciation	Morphology	Syntax	Vocabulary	Idea/Content	Notes
الجملة \ العبارة	عدد الكلمات	المستوى	عدد الكلمات	تهجئة/إملاء	تصريف واشتقاق	تركييب نحوية	مفردات	فكرة / محتوى	ملاحظات
خَيْرُ	1	و (صف 2)	6-waw	٩ هو أعلى عدد كلمات مطبعية غير متكررة بدون علامات الترقيم	• كلمات من ٥ مقاطع (بدون حساب حركات الإعراب)	• الفعل الماضي المفرد والجمع • جمع مذكر سالم	• جملة فيها فعلين (مثلا جملة فعلية مفعولها أن المصدرية)	• مفردات فصيحة - سامر ١	• المحتوى من حياة القارئ. • لا رمزية في النص.
جودي يقربي	2	ز (صف 2)	7-zay	١٠ هو أعلى عدد كلمات مطبعية غير متكررة بدون علامات الترقيم	• كلمات من ٦ مقاطع أو أكثر (بدون حساب حركات الإعراب) • أفعال/أسماء معتلّة الآخر	• الفعل الماضي المثني • فاعل المضارع المثني • فعل الأمر المفرد • جمع التذكير • واو القسم (والله)	• مفعول فيه (ظروف زمان ومكان) • حال • أداة الاستفهام هل	• مفردات فصيحة شائعة - سامر ٢	• بعض الرمزية أو عدم التصريح المباشر بكل المقصود في الجملة
بيروت في يوليو ١٩٦٦	4	ح (صف 3)	8-ha	١١ هو أعلى عدد كلمات مطبعية غير متكررة بدون علامات الترقيم	• فعل الأمر الجمع • نون النسوة في الأسماء والأفعال (انتظرن دورهن) • سوأيق أخرى: سين الاستقبال ، واو الاستئناف، فاء العطف • ثم ، حتى ، أو ، أم ، لكن ، أمّا	• المفعول المطلق • المفعول لأجله • المفعول معه • جملة فعلية تتعدى إلى مفعولين	• مفردات فصيحة - سامر ١ و سامر ٢ • أحرف النفي • الأرقام (العربية أو الهندية) 1,000,000-1,001	• بعض الرمزية يحتاج معها القارئ إلى مساعدة من يشرح له المقصود من الفكرة	
كتابة خطِّة لمشروع الوحدة	4	ك (صف 4)	11-kaf	٢٠ هو أعلى عدد كلمات مطبعية غير متكررة بدون علامات الترقيم	• المشتقات على أنواعها (تتركز على المشتقات العاملة لاسيما اسم الفاعل واسم المفعول)	• جملة اسمية خبرها جملة اسمية (فيها مبتدآن) • إضافة خيالية (لفظية) • طويل القامة	• أسماء الوصل المثني • متلازمات لفظية مثل شارد الذهن، وارف الظلال	• هناك درجة من الرمزية وحاجة للمعرفة السابقة كي يفهم المقصود من الجملة	
اجتمع أهل في العيد.	4	و (صف 2)	6-waw	٩ هو أعلى عدد كلمات مطبعية غير متكررة بدون علامات الترقيم	• كلمات من ٥ مقاطع (بدون حساب حركات الإعراب)	• الفعل الماضي المفرد والجمع • جمع مذكر سالم	• جملة فيها فعلين (مثلا جملة فعلية مفعولها أن المصدرية)	• مفردات فصيحة - سامر ١	• المحتوى من حياة القارئ. • لا رمزية في النص.
ولا خالطنا عجز ولا حور	4	ل (صف 5)	12-lam	لا حد لعدد الكلمات المطبعية	• التصغير	• جمل اعترافية (تفسير - دعاء...) • استثناء • أمين - خي - هائم - هاك • حبس • هيا - هيت - هلم إلى - • رويك • بعد اسم الإشارة • تمييز	• مفردات فصيحة - سامر ٣ • اسم الفعل : إيه - مة - • أمين - خي - هائم - هاك • هيا - هيت - هلم إلى - • رويك • الأرقام (العربية أو الهندية < 1,000,000 • ذو • بل - بلى - أجل)	• هناك درجة من الرمزية وحاجة للمعرفة السابقة كي يفهم المقصود من الجملة	

Inter Annotator Agreement

Stage	#Sets	Distance	Acc ¹⁹	± 1 Acc ¹⁹	QWK
Pilot 3	1	1.69	37.5%	58.5%	79.3%
Phase 1	2	1.38	48.4%	64.4%	80.2%
Phase 2A	6	1.21	49.4%	67.4%	72.4%
Phase 2B	10	0.80	67.6%	78.3%	78.8%
Overall / Macro	19	1.04	58.2%	72.3%	76.9%
Phase 2 / Macro	16	0.96	60.8%	74.2%	76.4%
Phase 2 / Micro	16	0.95	61.1%	74.4%	81.8%

Average pairwise inter-annotator agreement (IAA) across different annotation stages.

Inter Annotator Agreement

Sentence (Arabic)	A1	A2	A3	A4	A5	UL	MM	Comments
<p>أبي.. أبي..</p> <p><i>Dad .. Dad .. [lit. my father .. my father ..]</i></p>	2	2	2	3	3	3	1	First person singular pronoun is level 3.
<p>اِحْتِضَانُ الْأُمِّ لَهُمْ.</p> <p><i>The mother's embrace for them.</i></p>	9	12	5	5	5	5	7	Disagreement over اِحْتِضَان 'embrace': standard or dialect aligned.
<p>أشعر بالتعب والجوع..</p> <p><i>I feel tired and hungry..</i></p>	9	9	9	9	4	9	5	Vocabulary describing emotions (level 9).
<p>يتم ضمان حيادية الإدارة بموجب القانون.</p> <p><i>Administrative neutrality is guaranteed by law.</i></p>	12	12	12	14	12	12	2	Disagreement over حيادية 'neutrality': general advanced or specialized.

Table 8: Examples of Annotator Disagreements with Unified Levels (UL) and Max-Min Differences (MM)

Selected examples of IAA errors.

Corpus Selection

- 1,922 documents, from 30 different resources.

- Manually categorized into **3 domains**

- **Arts & Humanities**
- **Social Sciences**
- **STEM**

- And **3 readership groups**

- **Foundational**: up to grade 4
- **Advanced**: individuals with average adult reading abilities
- **Specialized**: from grade 9

	#Documents		#Sentences		#Words	
Arts & Humanities	1,367	71%	50,442	73%	652,995	63%
Social Sciences	375	20%	14,319	21%	275,731	27%
STEM	180	9%	4,680	7%	110,645	11%
Foundational	633	33%	27,781	40%	314,068	30%
Advanced	731	38%	22,696	33%	381,660	37%
Specialized	558	29%	18,964	27%	343,643	33%
	1,922	100%	69,441	100%	1,039,371	100%

Corpus Observations

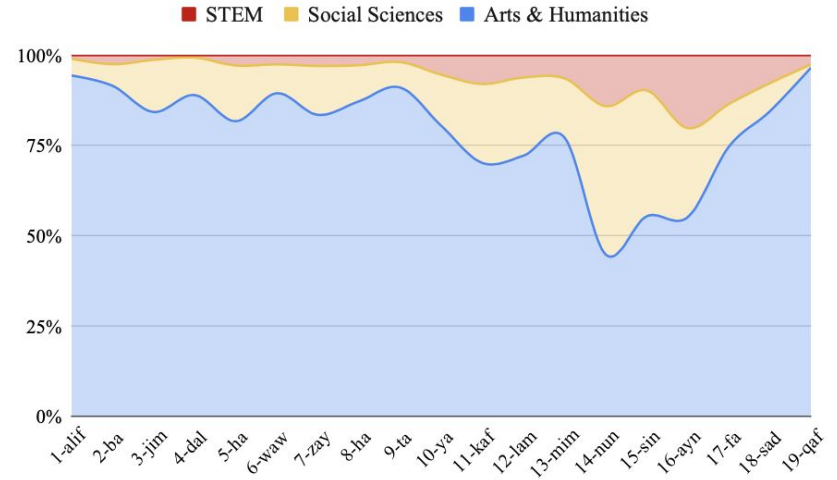
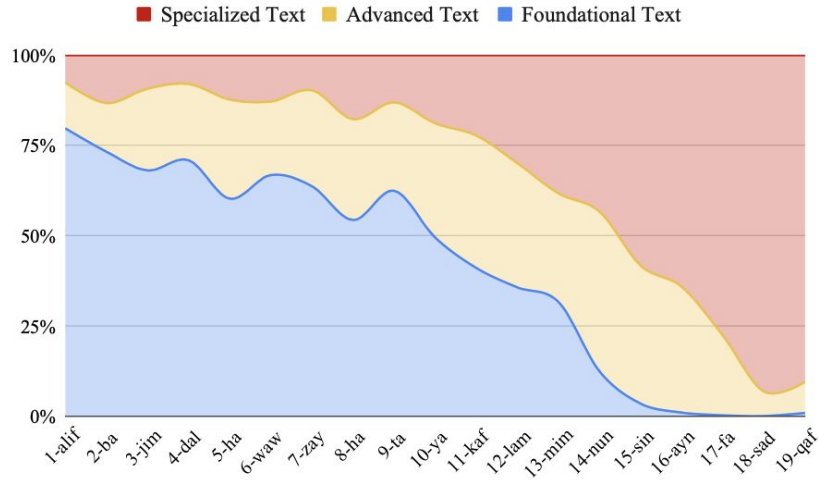


Figure 4: The relative distribution of readership groups and domains across **BAREC** levels.

Automatic Assessment Evaluation

- We define Readability Assessment as an ordinal classification task.
- Metrics: Distance, Accuracy, ± 1 Accuracy, Quadratic Weighted Kappa
- Fine-tuned AraBERTv02 (Antoun et al., 2020)

Train	Distance	Acc¹⁹	± 1 Acc¹⁹	QWK	Acc⁷	Acc⁵	Acc³
12.5%	1.35	45.0%	61.3%	77.2%	56.8%	63.0%	71.3%
25.0%	1.33	46.9%	63.0%	77.6%	58.8%	64.3%	72.3%
50.0%	1.16	52.4%	68.1%	80.7%	62.9%	67.6%	74.0%
100.0%	1.09	55.8%	69.4%	81.0%	64.9%	69.1%	74.7%

Table 9: Performance at different training data sizes across multiple evaluation metrics.



A Large and Balanced Corpus for Fine-grained Arabic Readability Assessment

Khalid Elmadani | Nizar Habash | Hanada Taha

Findings of the Association for Computational Linguistics: ACL 2025

BAREC Analyzer

BAREC Analyzer determines the readability level of Arabic texts using artificial intelligence.

سمير وحبّ الكتب
سمير طفلٌ صغير، يحبّ القراءة كثيرًا. كلّ يوم يقرأ كتابًا جديدًا قبل النوم. في المدرسة، يختار سмир كتبًا صعبة ليفهم كلمات جديدة. أصبح يكتب قصصًا قصيرة ويقرأها لأصدقائه في الصف. بفضل حبّه للكتب، صار يفكر بذكاء ويحلّ المسائل بسهولة.

ANALYZE

Results:

Readability Levels	Expert			Proficient			Advanced		Intermediate				Beginner						
	ق 19	ص 18	ف 17	غ 16	س 15	ن 14	م 13	ل 12	ك 11	ي 10	ط 9	ح 8	ز 7	و 6	هـ 5	د 4	ج 3	ب 2	أ 1

Level 10 Ya		Text Level
سمير وحبّ الكتب		Level 7 Zay
سمير طفلٌ صغير، يحبّ القراءة كثيرًا.		Level 3 Jim
كلّ يوم يقرأ كتابًا جديدًا قبل النوم.		Level 10 Ya
في المدرسة، يختار سмир كتبًا صعبة ليفهم كلمات جديدة.		Level 8 Ha
أصبح يكتب قصصًا قصيرة ويقرأها لأصدقائه في الصف.		Level 10 Ya
بفضل حبّه للكتب، صار يفكر بذكاء ويحلّ المسائل بسهولة.		Level 10 Ya



BAREC Shared Task 2025

Arabic Readability Assessment

The Third Arabic Natural Language Processing Conference (ArabicNLP 2025) @ EMNLP 2025

Suzhou, China



Conclusion & Future Work

- We introduced **BAREC**, the largest manually annotated dataset for Arabic Readability Assessment.
- We achieve a high level of agreement between annotators.
- We plan to expand the corpus, enhancing its size and diversity to cover additional genres and topics.
- We also aim to add annotations related to vocabulary leveling and syntactic treebanks to study less-explored genres in syntax.
- We look forward to readability leveling being used for other NLP tasks and as an analytical tool - e.g. educational systems, summarization, reasoning, etc.

Thank You

شكراً

Website



Data



Questions?