

Interpretable Arabic Readability Assessment Using Linguistic Rules and Expert-Guided Annotations

Anonymous ACL submission

Abstract

This work focuses on explaining Arabic readability decisions by explicitly modeling the linguistic phenomena underlying text complexity. Building on the Balanced Arabic Readability Evaluation Corpus (BAREC), we introduce BAREC-X, a new dataset of 1,793 sentences annotated with expert-provided linguistic justifications aligned with the BAREC guidelines, enabling direct evaluation of explanation quality. We further propose a fully interpretable, rule-based readability model grounded in linguistically motivated features spanning morphology, syntax, vocabulary, syllabic structure, and content complexity. The model mirrors the BAREC annotation process, produces structured human-readable explanations, and supports both readability level prediction and linguistic reasoning generation. As a final contribution, we present the first reasoning-annotated Arabic readability dataset, achieving an average inter-annotator agreement of 93.3% in terms of at least one shared justification. We also report the first results for an automated Arabic readability reasoning system, which attains 65.8% agreement with human annotators under the same criterion.

1 Introduction

Readability assessment plays a key role in education, content creation, and language learning. While extensively studied in English (Fountas and Pinnell, 2006; Heilman et al., 2008; Vajjala and Lučić, 2018), Arabic readability at fine-grained and interpretable levels remains underexplored. Most existing approaches focus on predicting readability levels (Al-Khalifa and Al-Ajlan, 2010; Nassiri et al., 2018; Elmadani et al., 2025b), often at the expense of transparency. In educational settings, however, understanding why a text is assigned a particular level is as important as the prediction itself. This challenge is especially pronounced for Arabic due to its rich morphology, syntactic vari-

Sentence	عالم غريب المعالم world _{ms} strange _{ms} features _{mp} 'A world of strange features'
Readability	Level 11-kaf (Grade 4)
Linguistic Reason	إضافة لفظية False Idafa

Table 1: Readability Assessment Example based on Linguistic Features

ation, and the limited availability of explanation-oriented resources.

This work addresses this gap by introducing an interpretable, feature-based approach to Arabic readability assessment grounded in the Balanced Arabic Readability Evaluation Corpus (BAREC) (Elmadani et al., 2025a). Linguistically motivated features are designed in accordance with the BAREC annotation framework (Habash et al., 2025). It covers morphology, syntax, vocabulary, syllabic structure, and content complexity. The approach explicitly models the linguistic phenomena underlying readability decisions and generates human-readable explanations alongside its predictions, as shown in table 1.

The contributions of this work are as follows:

- We introduce the first reasoning-annotated Arabic readability dataset (BAREC-X) aligned with the BAREC framework, comprising 1,793 sentences with expert-provided linguistic justifications alongside readability labels.
- We propose a fully interpretable, feature-based readability classifier grounded in rule-based linguistic features aligned with the BAREC framework, enabling transparent and linguistically grounded explanations.

The paper is organized as follows. After reviewing related work (§2), we present the linguistic

background and BAREC guidelines (§3), describe the proposed dataset (BAREC-X) and model (§4), and detail the experimental setup (§5). We then report and discuss the results (§6, §7, §8) and conclude in §9.

2 Related Work

Automatic readability assessment has been widely studied in NLP, resulting in a variety of datasets, models, and evaluation frameworks across languages. Recent advances in readability assessment have been driven by deep learning. Neural architectures, including recurrent models and hierarchical attention networks, have shown strong performance by learning representations directly from text (Sun et al., 2020). Transformer-based pretrained language models, such as BERT (Devlin et al., 2019), further improved accuracy by capturing contextualized linguistic information (Deutsch et al., 2020; Liu et al., 2025). Despite their effectiveness, these models are typically black-box systems that provide limited insight into the linguistic factors underlying their predictions.

As readability models became more complex, a critical question emerged: how to explain why a text is assigned a particular readability level. This has motivated research on computational methods to improve the interpretability of readability classifiers. One line of work applies post-hoc explanation methods that treat the underlying model as a black box. Model-agnostic approaches such as LIME (Ribeiro et al., 2016) and SHAP (Lundberg and Lee, 2017) have been used to explain individual predictions by estimating feature importance scores. While these methods can provide local explanations for complex models, their outputs are often difficult for non-expert users to interpret and may not align cleanly with linguistically meaningful categories (Gilpin et al., 2018). Moreover, such explanations are indirect and depend on the behavior of the underlying model rather than reflecting explicit linguistic reasoning.

An alternative direction focuses on incorporating linguistically motivated features directly into readability models. Traditional readability formulas are inherently interpretable but limited in expressiveness, while hybrid approaches combine neural representations with handcrafted linguistic features to balance performance and transparency (Deutsch et al., 2020). Although these methods demonstrate that linguistic features can be beneficial, particu-

larly in low-resource settings, their explanatory capacity is often implicit and rarely evaluated against human-provided reasoning.

Arabic readability research has additional challenges due to the language’s rich morphology, flexible word order, and orthographic variation. Machine learning approaches extended Arabic readability modeling by incorporating morphological, lexical, and syntactic features. For example, Al-Khalifa and Al-Ajlan (2010) demonstrated the effectiveness of SVM-based classifiers using handcrafted features, while more recent work showed that pretrained models such as AraBERT (Antoun et al., 2020) achieve higher predictive accuracy (Berrichi et al., 2024; Liberato et al., 2024). Nevertheless, these neural approaches largely prioritize performance and offer limited transparency into the linguistic drivers of readability decisions.

The development of annotated datasets has played a crucial role in advancing Arabic readability research. Existing resources include CEFR-based datasets for Arabic as a foreign language (Khallaf and Sharoff, 2021; Naous et al., 2023) and school-oriented corpora such as DARES (El-Haj et al., 2024). BAREC (Elmadani et al., 2025a) represents the largest and most fine-grained Arabic readability dataset to date, comprising 19 levels spanning early childhood to postgraduate comprehension, with high inter-annotator agreement. BAREC provides detailed linguistic guidelines for annotation (Habash et al., 2025), but its labels do not explicitly encode the reasoning behind each decision.

In contrast to prior work that emphasizes prediction accuracy or relies on post-hoc explanations, this study centers interpretability as a primary objective. We build on BAREC by introducing a reasoning-annotated subset in which expert annotators provide explicit linguistic justifications aligned with the BAREC guidelines. We further propose a fully interpretable, rule-based readability model that mirrors the human annotation process and produces structured, linguistically grounded explanations alongside predictions. By directly evaluating agreement between system-generated explanations and human reasoning, our work establishes a transparent baseline for interpretable Arabic readability assessment.

3 Linguistic Background

This section outlines the linguistic foundations of Arabic readability assessment as defined in BAREC guidelines (Habash et al., 2025). BAREC adopts a fine-grained, pedagogically motivated readability framework inspired by Taha-Thomure (2017), in which sentences are assigned to one of 19 ordered readability levels corresponding to progressive educational stages, from early literacy to postgraduate comprehension (Habash et al., 2025). These levels follow the Abjad order of Arabic letters (1-alif, 2-ba, 3-jim, through to 19-qaf) and are designed to provide fine distinctions at lower levels, where readability variation is greatest.

Readability in BAREC is modeled as the interaction of multiple linguistic dimensions rather than a single surface-level measure. Each sentence is assigned a readability level based on the most complex linguistic phenomenon it contains, reflecting a 'most complex feature' principle. As a result, the presence of a single high-level feature is often sufficient to determine the final readability level, regardless of simpler features present elsewhere in the sentence.

Annotators evaluate sentences across six core linguistic dimensions: spelling and phonology, word count, morphology, syntax, vocabulary, and content. These dimensions jointly define the progression of difficulty across the 19 levels, with different dimensions becoming dominant at different stages of readability development. The same dimensions guide both the annotation of our newly introduced dataset and the design of the interpretable models proposed in this work.

- **Spelling and Phonology:** Word length and syllabic structure influence difficulty, with more complex syllable patterns associated with higher levels. Diacritics are not considered. This dimension primarily constrains readability up to Level 7.
- **Word Count:** The number of unique words serves as an upper bound on readability at early levels and is mainly relevant up to Level 11.
- **Morphology:** Morphological complexity reflects inflectional and derivational patterns, including clitics, verb forms, plurality, and less frequent constructions. Morphology influences readability up to approximately Level 13.

- **Syntax:** Sentence structure and syntactic relations contribute to difficulty, with non-canonical and complex constructions appearing at higher levels. Syntactic complexity becomes a key determinant up to Level 15.
- **Vocabulary:** Lexical difficulty is driven by word familiarity and usage, with higher levels introducing rarer, technical, or classical vocabulary. Vocabulary spans the full readability range up to Level 19.
- **Content:** Conceptual difficulty reflects the level of abstraction and required prior knowledge, ranging from concrete everyday concepts to abstract or domain-specific content. Content applies across all levels and extends to Level 19.

The complete BAREC annotation guidelines are provided in Appendix A.

4 Methodology

Our methodology combines data-centric annotation with interpretable modeling approaches for Arabic readability assessment. Building on the Balanced Arabic Readability Evaluation Corpus (BAREC), we introduce a newly annotated reasoning-aware subset, BAREC-X, in which sentences are labeled with both readability levels and explicit linguistic justifications aligned with the BAREC guidelines. This dataset enables the direct study and evaluation of interpretability in readability modeling.

4.1 Reasoning-Annotated Readability Dataset

While BAREC provides a strong foundation for sentence-level readability classification, its annotations are limited to readability labels and do not explicitly capture the linguistic reasoning behind each decision. To support interpretable readability modeling, we introduce BAREC-X, a newly annotated subset of 1,793 sentences derived from BAREC. These sentences were selected from the inter-annotator agreement (IAA) portion of the corpus to ensure high-quality readability labels.

Three annotators, native Arabic-speaking educators with prior experience in Arabic readability annotation from a relevant project, were assigned approximately 800 sentences in total, with around 300 sentences shared across all annotators to enable the measurement of inter-annotator agreement on interpretive reasoning. The data was organized into 8 batches of approximately 100 sentences to

268 facilitate the annotation process. The exact sen- 315
269 tence counts per batch and their distribution across 316
270 readability levels are provided in Appendix A.3. 317
271 Each annotator was given 5 unique batches and 3 318
272 batches were shared among all 3 annotators. They 319
273 were asked to provide explicit linguistic justifica- 320
274 tions for each sentence underlying the given level 321
275 from the BAREC corpus. Selecting from a fixed 322
276 set of categories defined by the BAREC guidelines, 323
277 each sentence may be associated with multiple lin- 324
278 guistic reasons, reflecting the presence of different 325
279 factors contributing to readability. An example of 326
280 the annotation process sheet used by annotators, 327
281 along with the structured output format, is pro- 328
282 vided in Appendix A.3. This annotation process 329
283 results in aligned sentence-level readability labels 330
284 and human-provided explanations, enabling direct 331
285 evaluation of model interpretability. 332

286 4.2 Feature-Based Readability Model 333

287 We propose a fully interpretable, feature-based 334
288 readability model designed to mirror the decision 335
289 process followed by human annotators under the
290 BAREC framework. Rather than learning implicit
291 patterns from data, the model relies exclusively on
292 handcrafted linguistic features and deterministic
293 rules to assign readability levels and generate ex-
294 plicit explanations. The model supports two modes
295 of operation: (i) as a readability classifier, where
296 a sentence is provided as input and the model pre-
297 dicta readability level together with its underlying
298 linguistic justifications, and (ii) as an explanatory
299 system, where a sentence and a specified readabil-
300 ity level are provided, and the model outputs the
301 linguistic reasons associated with that level.

302 **Extraction of Linguistic Features** Linguistic 336
303 features corresponding to the dimensions described 337
304 in Section 3 are extracted using a combination 338
305 of morphological analysis with CAMEL Tools¹ 339
306 (Obeid et al., 2020), syntactic parsing with Camel- 340
307 Parser² (Elshabrawy et al., 2023), surface-level 341
308 measurements, and external lexical resources. In 342
309 total, the model operates over 66 linguistically mo- 343
310 tivated features. These include 25 morphological 344
311 features, 18 syntactic features, and 19 vocabulary- 345
312 based features, in addition to surface-level features 346
313 such as word count and syllable count, as well as a 347
314 vocabulary and content-level feature. 348

¹https://github.com/CAMEL-Lab/camel_tools, ver- 349
sion 1.2.0 350

²https://github.com/CAMEL-Lab/camel_parser, ver- 351
sion 2.0 352

Morphological features capture inflectional and 315
derivational properties such as affixes, clitics, verb 316
tense, voice, and plural forms. Syntactic features 317
encode the presence of specific dependency struc- 318
tures associated with increased sentence complex- 319
ity. Surface-level features include word count and 320
syllable count, computed in accordance with the 321
BAREC guidelines. Vocabulary difficulty is mod- 322
eled through a level-based lexicon constructed from 323
the BAREC training data, where lemma-POS pairs 324
are associated with the earliest readability level at 325
which they appear. To mitigate annotation noise, 326
relaxed frequency thresholds are applied when as- 327
signing vocabulary levels. This lexicon is fur- 328
ther augmented using the SAMER readability lex- 329
icon (Al Khalil et al., 2020) and a supplementary 330
BAREC dialectal lexicon, enabling the model to ac- 331
count for curriculum-aligned vocabulary difficulty 332
and lexical variation between Modern Standard 333
Arabic and dialects. Detailed feature definitions 334
and extraction rules are provided in Appendix B. 335

Readability Modeling In the first mode of oper- 336
ation, where the model functions as a readability 337
classifier, extracted reasons for each sentence are 338
mapped to their associated readability levels, with 339
each activated reason contributing a candidate level. 340
The final readability prediction is determined using 341
a max-based aggregation strategy, whereby the sen- 342
tence is assigned the highest level triggered by any 343
reason. Table 2 illustrates this reason competition 344
mechanism by showing an example sentence with 345
its activated reasons and corresponding readabil- 346
ity levels; only triggered reasons are considered, 347
and the highest assigned level defines the final pre- 348
diction. In addition to predicting a readability la- 349
bel, the model generates a structured set of linguis- 350
tic reasons corresponding to the activated features. 351
These reasons constitute the model’s explanation 352
output, ensuring full transparency and traceability 353
between predictions and linguistic evidence. 354

In the second mode of operation, where the 355
model is used as an explanatory system, the model 356
restricts its output to the subset of extracted reasons 357
whose associated levels match the given input level, 358
returning only the linguistic features that justify 359
that level for the sentence. When no explicit lin- 360
guistic features are activated at the specified level, 361
the model falls back to vocabulary and content- 362
based features to provide a minimal yet informative 363
explanation. 364

Category	Activated Reason	Assigned Level
WC	unique word count = 3	2
SP	syllable count = 3	2
M	definite article prefix (Al ل)	3
	broken plural	7
V	adjective	2
	vocabulary difficulty	5
S	noun + adjective	2
	false idafa	11

Table 2: Example of feature competition for 'عالم غريب المعالم'. Only activated features are shown. WC=Word Count, SP=Spelling, M=Morphology, V=Vocabulary, S=Syntax.

Dataset	Train	Dev	Test
BAREC	54,845	7,310	7,286
BAREC-X	1,409	196	188

Table 3: Dataset statistics for BAREC and the reasoning-annotated subset BAREC-X across train, development, and test splits.

5 Experimental setup

5.1 Data Splits

The readability prediction experiment uses the official data split provided with the BAREC dataset. To evaluate the explanatory system, we augment the IAA portion of BAREC with reasoning annotations, resulting in BAREC-X (see section 4.1). Since the IAA portion spans across Train, Dev and Test splits, BAREC-X follows the same splits although we did not train the model for this task explicitly. In this paper, we report on the three splits separately. However, for future work that would use this data for training, we recommend following the splits shown in table 3.

5.2 Evaluation Metrics

We evaluate interpretability from two complementary perspectives: **Reason+Category (RC)** and **Category-only (C)**. The RC setting assesses whether a model correctly identifies the *exact linguistic reasons* within each category (e.g., specific morphological or syntactic phenomena) provided by human annotators. In contrast, the C setting evaluates performance at a coarser level by considering only the broader linguistic category (e.g., morphology, syntax, vocabulary), regardless of the specific reason selected within that category.

For both evaluation settings, we report three met-

rics to capture different aspects of alignment with human-provided explanations. First, **Exact Match** measures whether the model’s predicted set of reasons exactly matches the annotator-provided set. Second, **Jaccard Similarity** is used to quantify the degree of overlap between predicted and gold sets, providing a graded measure of interpretability rather than a binary outcome. Finally, **At Least One** measures whether the model correctly predicts at least one gold reason, capturing partial but meaningful interpretive overlap.

We adopt the official evaluation metrics defined by BAREC to assess readability classification performance:

- **Quadratic Weighted Kappa (QWK)** - the primary evaluation metric, which penalizes larger misclassifications more heavily (Cohen, 1968; Doewes et al., 2023).
- **Accuracy ($\text{Acc}^{19} / \text{Acc}^7 / \text{Acc}^5 / \text{Acc}^3$)** - classification accuracy computed at different levels of granularity by collapsing the original 19 readability labels into 7, 5, or 3 categories.
- **Adjacent Accuracy ($\pm 1 \text{Acc}^{19}$)** - an off-by-one tolerance measure that counts predictions within one readability level of the gold label as correct.
- **Average Distance** - the average absolute difference between two sets of labels.

5.3 Model Architecture and Hyperparameters

The feature-based model does not require traditional training for most of its components, particularly for morphological and syntactic features. Each linguistic pattern is explicitly defined by a unique code and linked to one or more readability levels according to the BAREC guidelines; these patterns are treated as discrete codes within the system. During inference, all relevant features are extracted from a sentence, and each activated code contributes its associated readability level. A complete list of linguistic codes and their corresponding readability levels is provided in appendix B.1.

In its default configuration, the model assigns a final readability level using a max-based aggregation strategy, selecting the level associated with the most complex activated feature. In addition, the same feature extraction pipeline can be used in an explanatory mode, where the model filters the activated features to return only those corresponding

to a specified readability level. This design enables both forward prediction and targeted explanation using a shared, fully transparent mechanism.

Figure 1(a) shows the pipeline used for readability level prediction, while Figure 1(b) depicts the explanatory mode, in which a readability level is provided as input and the model returns the corresponding linguistic reasons.

For the content-based feature, which captures thematic or conceptual complexity, we rely on a learned approach. We fine-tune a transformer-based model to classify sentences into one of the eight content levels defined by the BAREC guidelines. The model is trained on the training split of the BAREC dataset. We fine-tune AraBERTv02 (Antoun et al., 2020),³ configured using AdamW ($\beta_1=0.9$, $\beta_2=0.999$, $\epsilon=1e-08$) with a learning rate of 5×10^{-5} , linear learning rate decay, batch sizes of 64. All experiments are evaluated one a single run.

The output of the content-level classifier is treated like any other feature and contributes to the final readability level using the same max-based logic.

6 Results

6.1 Reasoning Dataset IAA Evaluation

To assess the quality and consistency of the reasoning annotations in the newly introduced interpretable subset, we compute inter-annotator agreement (IAA) on the 300 selected sentences. Since each sentence may be associated with multiple reasoning labels, agreement is evaluated using the same interpretability metrics described in Section 5.2. Agreement is computed pairwise across annotators.

We report IAA results under both the **RC** and **C** evaluation settings, capturing agreement at fine-grained and coarse-grained levels of interpretation. The resulting scores provide an empirical measure of annotation reliability and validate the quality of the reasoning labels. Detailed agreement statistics are presented in Table 4.

In addition, we apply a unification procedure to derive a single consolidated set of reasoning labels per sentence. For each sentence, unified reasons are defined as those agreed upon by all three annotators; in cases where no such consensus exists,

Annotator Pair	Exact		Jaccard		At Least 1	
	RC	C	RC	C	RC	C
A1-A2	68.4	79.5	68.0	78.0	91.3	94.3
A2-A3	82.8	85.5	80.2	84.4	95.3	97.0
A1-A3	70.4	80.8	70.9	79.5	93.3	95.0
Average	73.9	81.9	73.0	80.6	93.3	95.4

Table 4: Pairwise inter-annotator agreement (IAA). All values are reported as percentages(%)

reasons agreed upon by at least two annotators are retained.

The inter-annotator agreement results demonstrate strong consistency in the reasoning annotations, supporting the reliability of the proposed dataset. Agreement is consistently higher under the coarse-grained **C** setting than under the finer-grained **RC** setting, as expected. Substantial Exact Match and Jaccard scores are observed across annotator pairs, while the At Least One metric exceeds 90% in all cases, indicating near-universal agreement on at least one key linguistic factor. These results confirm that the annotation guidelines enable consistent interpretive judgments and that the dataset is suitable for evaluating interpretable readability models.

6.2 Feature-based System Evaluation

When the proposed feature-based model is used for readability level prediction, its performance varies systematically across the readability spectrum. In particular, evaluating the model separately on early-to-intermediate levels (1-11) yields substantially stronger results than treating all 19 levels as a single prediction task. This pattern reflects fundamental differences in how readability is defined within the BAREC framework. Levels 1-11 are governed by fine-grained, linguistically explicit guidelines that rely primarily on surface, morphological, and syntactic features, which are well suited to rule-based modeling. In contrast, higher levels (12-19) depend increasingly on vocabulary rarity and abstract content complexity, which are inherently more subjective and less amenable to deterministic approaches. Motivated by this distinction, we report results both for the full 19-level setting and for a focused evaluation on levels 1-11. As shown in Table 5, performance on the full level range remains moderate, reflecting the challenges posed by higher-level semantic abstraction. However, restricting evaluation to levels 1-11 leads to consistent and substantial im-

³<https://huggingface.co/aubmindlab/bert-base-arabertv02>

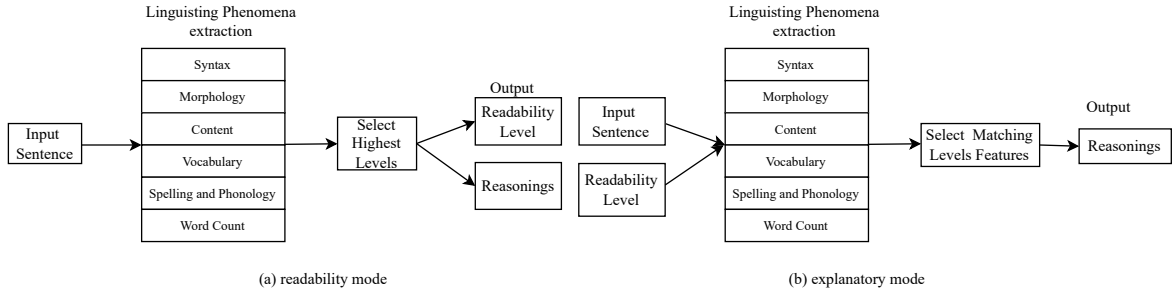


Figure 1: Feature-Based Custom Model Pipeline

	1-19		1-11	
	FBM	BL	FBM	BL
Count	7,310	7,310	3,949	3,949
Acc¹⁹	34.9	56.6	58.2	64.1
± 1 Acc¹⁹	47.9	69.9	73.9	77.0
Avg. Distance	2.2	1.1	1.0	0.9
QWK	48.2	80.0	69.6	73.6
Acc⁷	42.3	65.9	69.9	73.8
Acc⁵	49.9	70.3	84.0	82.3
Acc³	58.3	76.5	100.0	91.4

Table 5: BAREC dev set readability prediction performance of the feature-based model (FBM) compared with the BAREC baseline (BL). all metrics are in %

527 improvements across all metrics. While the BAREC
528 baseline achieves higher raw accuracy and QWK
529 on levels 1-11, the performance gap is relatively
530 modest. Crucially, the proposed model produces
531 fully transparent, linguistically grounded explana-
532 tions alongside its predictions, making it partic-
533 ularly well suited for educational and diagnostic
534 applications where interpretability and actionable
535 feedback are primary requirements.

536 The primary objective of the feature-based sys-
537 tem is interpretability rather than maximizing pre-
538 dictive performance. Accordingly, we evaluate the
539 system by comparing its generated explanations
540 against the human-provided reasoning annotations
541 in the BAREC-X subset.

542 Table 6 reports agreement between the system
543 and human annotators across data splits inherited
544 from BAREC. The development split is not used
545 for model training or tuning in this work and is re-
546 ported solely for completeness. The model operates
547 in the same inference-only manner across the train-
548 ing, development, and test splits, and exploration
549 of split-specific training or evaluation strategies is
550 left for future work.

Split	Exact		Jaccard		At Least 1	
	RC	C	RC	C	RC	C
Train	43.0	59.3	40.7	50.3	65.4	73.0
Dev	36.2	51.5	35.4	49.1	65.8	77.0
Test	37.8	59.6	39.0	53.3	70.7	80.3

Table 6: Agreement between the feature-based system and human annotations across data splits. All values are percentages (%).

551 Within this evaluation setting, while Exact
552 Match performance is lower in the fine-grained
553 RC condition, agreement improves substantially at
554 the category level, indicating that the system gen-
555 erally identifies the correct linguistic dimension
556 influencing readability. Moreover, the relatively
557 high At Least One scores suggest that the system
558 frequently captures at least one core feature suffi-
559 cient to justify a given readability level, even when
560 it does not fully align with the annotators’ complete
561 set of reasons.

562 7 Error Analysis

563 Table 7 presents a category-level analysis of the sys-
564 tem’s reasoning prediction performance. Substan-
565 tial variation is observed across linguistic dimen-
566 sions, reflecting differences in both feature explic-
567 itness and model coverage. Categories grounded
568 in surface-level and structurally explicit cues-most
569 notably Word Count and Morphology, achieve the
570 strongest performance, with F1 scores of 70.4%
571 and 68.2%, respectively. Syntactic features achieve
572 high precision (78.8%) but lower recall (40.7%),
573 suggesting that identified constructions are typi-
574 cally correct, though many valid instances remain
575 undetected. Vocabulary features show higher recall
576 (74.3%) than precision (58.5%), indicating broader
577 lexical coverage at the cost of occasional over-
578 generation, particularly in context-dependent cases.
579 Performance is weakest for Orthography and Con-

Category	Precision (%)	Recall (%)	F1 (%)
Spelling & Phon.	55.8	25.3	34.8
Word Count	61.3	82.6	70.4
Morphology	87.4	55.9	68.2
Syntax	78.8	40.7	53.7
Vocabulary	58.5	74.3	65.4
Content	4.5	27.4	7.7

Table 7: Precision, recall, and F1 scores for reasoning prediction across linguistic categories.

580 tent. Orthographic features are relatively sparse
581 and contribute less consistently to readability de-
582 cisions, resulting in low recall (25.3%). Content-
583 level reasoning performs poorest overall (F1 of
584 7.7%), which is expected given its reliance on se-
585 mantic abstraction, cultural knowledge, and inter-
586 preptive judgment-factors that are difficult to capture
587 through deterministic, rule-based features. Overall,
588 the results highlight a clear trade-off between inter-
589 pretable and semantic coverage. The proposed
590 feature-based system performs best in linguistically
591 explicit categories, while performance degrades for
592 phenomena requiring deeper semantic understand-
593 ing. This analysis reinforces the model’s suitability
594 as an interpretable diagnostic tool for readability
595 assessment and motivates future extensions that
596 incorporate semantic modeling for content-level
597 reasoning.

598 8 Discussion

599 This work demonstrates that Arabic readability can
600 be modeled in an interpretable and linguistically
601 faithful manner when annotation, feature design,
602 and evaluation are tightly aligned. The proposed
603 feature-based model performs strongest at early
604 and intermediate readability levels, where BAREC
605 guidelines rely on explicit morphological, syntactic,
606 and surface-level cues. In contrast, performance
607 declines at higher levels that depend more heavily
608 on abstract vocabulary and content interpretation,
609 highlighting the inherent limitations of determin-
610 istic rules for modeling semantic complexity. The
611 interpretability evaluation shows that, while exact
612 matching of fine-grained reasons remains challeng-
613 ing, the system reliably identifies the correct lin-
614 guistic dimension influencing readability. High At
615 Least One agreement indicates that the model fre-
616 quently captures a core justification aligned with
617 human reasoning, which is particularly relevant for
618 educational and diagnostic use cases where partial
619 but faithful explanations are sufficient.

620 Finally, the introduction of BAREC-X estab-
621 lishes the first reasoning-annotated Arabic readabil-
622 ity dataset, with strong inter-annotator agreement,
623 demonstrating that expert linguistic reasoning can
624 be consistently captured and evaluated. Together,
625 these results position interpretable, feature-driven
626 models as a robust and transparent baseline for
627 readability assessment, and motivate future work
628 on extending semantic coverage through hybrid or
629 learning-based approaches.

630 9 Conclusions and Future Work

631 This work presents an interpretable approach to
632 Arabic readability assessment grounded in the
633 BAREC annotation framework. We introduced
634 BAREC-X, the first reasoning-annotated Arabic
635 readability dataset, in which expert annotators pro-
636 vide explicit linguistic justifications, achieving an
637 average inter-annotator agreement of 93.3% un-
638 der the At Least One criterion. Building on this
639 resource, we developed a fully interpretable rule-
640 based model that mirrors the BAREC annotation
641 process and produces structured, human-readable
642 explanations for its predictions, constituting the
643 first reported results for automated Arabic readabil-
644 ity reasoning, with 65.8% agreement with human
645 annotations under the same criterion.

646 Experimental results demonstrate that the pro-
647 posed model performs competitively at early and in-
648 termediate readability levels (1-11), where linguis-
649 tic guidelines are most explicit, while also provid-
650 ing strong alignment with human reasoning. These
651 findings underscore the continued relevance of lin-
652 guistically informed models for educational NLP
653 tasks and establish a solid baseline for interpretable
654 Arabic readability assessment.

655 Future work includes expanding the feature set
656 to cover additional aspects of the BAREC guide-
657 lines that are not fully captured in the current im-
658 plementation. Further effort will also be devoted to
659 refining and debugging the linguistic tools used in
660 feature extraction to improve accuracy and robust-
661 ness. Finally, enlarging the reasoning-annotated
662 dataset would enable more comprehensive evalu-
663 ation and support the development of hybrid or
664 learning-based models that retain interpretability
665 while improving performance at higher readability
666 levels.

667 Limitations

668 Despite its strengths, this work has several limi-
669 tations. First, while the proposed feature-based
670 model is grounded in the BAREC annotation guide-
671 lines, it does not yet cover the full spectrum of
672 linguistic phenomena described in the guidelines.
673 Certain fine-grained constructions and rare linguis-
674 tic patterns are not explicitly modeled, which may
675 affect performance at higher readability levels. Sec-
676 ond, the accuracy of the extracted features depends
677 on the reliability of the underlying NLP tools, in-
678 cluding morphological disambiguation and depen-
679 dency parsing. Errors or ambiguities in these tools
680 may propagate to the feature representations and,
681 consequently, to both predictions and explanations.
682 Although we mitigate this by relying on well es-
683 tablished resources, further debugging and valida-
684 tion of the feature extraction pipeline are necessary.
685 Third, vocabulary- and content-based features are
686 derived from data-driven resources trained on the
687 BAREC training split. While this design choice is
688 consistent with standard evaluation practices, it lim-
689 its the interpretability evaluation to development
690 and test subsets and may introduce bias toward
691 the training distribution. Finally, the reasoning-
692 annotated dataset remains relatively modest in size.
693 Although it is sufficient for interpretability analy-
694 sis and qualitative evaluation, scaling the dataset
695 would enable more robust statistical analysis and
696 support future learning-based explanation models.

697 Ethical considerations

698 The annotation process is conducted with trans-
699 parency and fairness, with multiple annotators in-
700 volved to mitigate biases and ensure reliability. All
701 annotators are paid fair wages for their contribu-
702 tion.

703 The annotated dataset will be openly accessi-
704 ble to promote transparency, reproducibility, and
705 collaboration in Arabic language research.

706 We used AI writing assistance within the scope
707 of “Assistance purely with the language of the pa-
708 per” described in the ACL Policy on Publication
709 Ethics.

710 References

711 Hend S Al-Khalifa and Amani A Al-Ajlan. 2010. Au-
712 tomatic readability measurements of the arabic text:
713 An exploratory study. *Arabian Journal for Science
714 and Engineering*, 35(2 C):103–124.

- Muhamed Al Khalil, Nizar Habash, and Zhengyang
Jiang. 2020. [A large-scale leveled readability lex-
icon for Standard Arabic](#). In *Proceedings of the
Twelfth Language Resources and Evaluation Confer-
ence*, pages 3053–3062, Marseille, France. European
Language Resources Association. 715
716
717
718
719
720
- Wissam Antoun, Fady Baly, and Hazem Hajj. Arabert:
Transformer-based model for arabic language under-
standing. In *LREC 2020 Workshop Language Re-
sources and Evaluation Conference 11–16 May 2020*,
page 9. 721
722
723
724
725
- Wissam Antoun, Fady Baly, and Hazem Hajj. 2020.
[AraBERT: Transformer-based model for Arabic lan-
guage understanding](#). In *Proceedings of the 4th Work-
shop on Open-Source Arabic Corpora and Process-
ing Tools, with a Shared Task on Offensive Language
Detection*, pages 9–15, Marseille, France. European
Language Resource Association. 726
727
728
729
730
731
732
- Safae Berrichi, Naoual Nassiri, Azzeddine Mazroui, and
Abdelhak Lakhouaja. 2024. Exploring the impact
of deep learning techniques on evaluating arabic ll
readability. In *Artificial Intelligence, Data Science
and Applications*, pages 1–7, Cham. Springer Nature
Switzerland. 733
734
735
736
737
738
- Jacob Cohen. 1968. Weighted kappa: Nominal scale
agreement provision for scaled disagreement or par-
tial credit. *Psychological bulletin*, 70(4):213. 739
740
741
- Tovly Deutsch, Masoud Jasbi, and Stuart Shieber. 2020.
Linguistic features for readability assessment. *arXiv
preprint arXiv:2006.00377*. 742
743
744
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
Kristina Toutanova. 2019. [BERT: Pre-training of
deep bidirectional transformers for language under-
standing](#). In *Proceedings of the 2019 Conference of
the North American Chapter of the Association for
Computational Linguistics: Human Language Tech-
nologies, Volume 1 (Long and Short Papers)*, pages
4171–4186, Minneapolis, Minnesota. Association for
Computational Linguistics. 745
746
747
748
749
750
751
752
753
- Afrizal Doewes, Nughthoh Arfawi Kurdhi, and Akрати
Saxena. 2023. [Evaluating quadratic weighted kappa
as the standard performance metric for automated es-
say scoring](#). In *Proceedings of the 16th International
Conference on Educational Data Mining*, pages 103–
113, Bengaluru, India. International Educational Data
Mining Society. 754
755
756
757
758
759
760
- Mo El-Haj, Sultan Almujaivel, Damith Premasiri,
Tharindu Ranasinghe, and Ruslan Mitkov. 2024.
Dares: Dataset for arabic readability estimation of
school materials. In *Proceedings of the Workshop
on DeTermIt! Evaluating Text Difficulty in a Multi-
lingual Context@ LREC-COLING 2024*, pages 103–
113. 761
762
763
764
765
766
767
- Khalid Elmadani, Nizar Habash, and Hanada Taha.
2025a. A large and balanced corpus for fine-grained
arabic readability assessment. In *Findings of the As-
sociation for Computational Linguistics: ACL 2025*,
pages 16376–16400. 768
769
770
771
772

773	Khalid N Elmadani, Bashar Alhafni, Hanada Taha, and Nizar Habash. 2025b. Barec shared task 2025 on arabic readability assessment. In <i>Proceedings of The Third Arabic Natural Language Processing Conference: Shared Tasks</i> , pages 239–252.	826
774		827
775		828
776		829
777		830
778	Ahmed Elshabrawy, Muhammed AbuOdeh, Go Inoue, and Nizar Habash. 2023. CamelParser2.0: A State-of-the-Art Dependency Parser for Arabic. In <i>Proceedings of The First Arabic Natural Language Processing Conference (ArabicNLP 2023)</i> .	831
779		832
780		833
781		834
782		835
783	Irene C Fountas and Gay Su Pinnell. 2006. <i>Leveled books (k-8): Matching texts to readers for effective teaching</i> . Heinemann Educational Books.	836
784		837
785		838
786	Leilani H Gilpin, David Bau, Ben Z Yuan, Ayesha Bajwa, Michael Specter, and Lalana Kagal. 2018. Explaining explanations: An overview of interpretability of machine learning. In <i>2018 IEEE 5th International Conference on data science and advanced analytics (DSAA)</i> , pages 80–89. IEEE.	839
787		840
788		841
789		842
790		843
791		844
792	Nizar Habash, Hanada Taha-Thomure, Khalid N. Elmadani, Zeina Zeino, and Abdallah Abushmaes. 2025. Guidelines for fine-grained sentence-level Arabic readability annotation . In <i>Proceedings of the 19th Linguistic Annotation Workshop (LAW-XIX-2025)</i> , pages 359–376, Vienna, Austria. Association for Computational Linguistics.	845
793		846
794		847
795		848
796		849
797		850
798		851
799	Michael Heilman, Kevyn Collins-Thompson, and Maxine Eskenazi. 2008. An analysis of statistical models and features for reading difficulty prediction. In <i>Proceedings of the third workshop on innovative use of NLP for building educational applications</i> , pages 71–79.	852
800		853
801		854
802		855
803		856
804		857
805	Go Inoue, Bashar Alhafni, Nurpeiis Baimukan, Houda Bouamor, and Nizar Habash. 2021. The interplay of variant, size, and task type in Arabic pre-trained language models . In <i>Proceedings of the Sixth Arabic Natural Language Processing Workshop</i> , pages 92–104, Kyiv, Ukraine (Virtual). Association for Computational Linguistics.	858
806		859
807		860
808		861
809		862
810		863
811		864
812	Nouran Khallaf and Serge Sharoff. 2021. Automatic difficulty classification of arabic sentences. <i>arXiv preprint arXiv:2103.04386</i> .	865
813		
814		
815	Juan Piñeros Liberato, Bashar Alhafni, Muhamed Al Khalil, and Nizar Habash. 2024. Strategies for arabic readability modeling. <i>arXiv preprint arXiv:2407.03032</i> .	
816		
817		
818		
819	Fengkai Liu, Tan Jin, and John SY Lee. 2025. Automatic readability assessment for sentences: neural, hybrid and large language models. <i>Language Resources and Evaluation</i> , pages 1–32.	
820		
821		
822		
823	Scott M. Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions . <i>CoRR</i> , abs/1705.07874.	
824		
825		
	Tarek Naous, Michael J Ryan, Anton Lavrouk, Mohit Chandra, and Wei Xu. 2023. Readme++: Benchmarking multilingual language models for multi-domain readability assessment . <i>arXiv preprint arXiv:2305.14463</i> .	
	Naoual Nassiri, Abdelhak Lakhouaja, and Violetta Cavalli-Sforza. 2018. Modern standard arabic readability prediction. In <i>Arabic Language Processing: From Theory to Practice</i> , pages 120–133, Cham. Springer International Publishing.	
	Ossama Obeid, Nasser Zalmout, Salam Khalifa, Dima Taji, Mai Oudah, Bashar Alhafni, Go Inoue, Fadhil Eryani, Alexander Erdmann, and Nizar Habash. 2020. CAMEL tools: An open source python toolkit for Arabic natural language processing . In <i>Proceedings of the 12th Language Resources and Evaluation Conference</i> , pages 7022–7032, Marseille, France. European Language Resources Association.	
	Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "why should i trust you?": Explaining the predictions of any classifier . In <i>Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16</i> , page 1135–1144, New York, NY, USA. Association for Computing Machinery.	
	Yuxuan Sun, Keying Chen, Lin Sun, and Chenlu Hu. 2020. Attention-based deep learning model for text readability evaluation . <i>2020 International Joint Conference on Neural Networks (IJCNN)</i> , pages 1–8.	
	Hanada Taha-Thomure. 2017. <i>Arabic Language Text Leveling</i> (معايير هنادا طه لتصنيف مستويات النصوص العربية). Educational Book House (دار الكتاب التربوي للنشر والتوزيع).	
	Sowmya Vajjala and Ivana Lučić. 2018. On-eStopEnglish corpus: A new corpus for automatic readability assessment and text simplification . In <i>Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Educational Applications</i> , pages 297–304, New Orleans, Louisiana. Association for Computational Linguistics.	

A.2 English Translation

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	• Direct, explicit, and concrete idea. • No symbolism in the text.
2-ba	1	Novice Low	≤2	• Three-syllable words	• Prtocolitic: Definite article <i>Al+</i> • Proclitic: Conjunction <i>wa+</i> • Enclitic: First Person Singular pronoun	• Apposition (full) • Demonstratives	• 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	• Plural imperfective verb • Prepositional proclitics • Numated adverbials			• Common MSA vocabulary - SAMER I • Singular demonstrative pronoun • Numbers: 11-100	
4-dal		Novice Mid	≤6	• Words with an elongated Alif (e.g. /äsilif/)			• Verbal sentence w/o direct object • Preposition and object	
5-ha	2	Novice High	≤8	• Four-syllable words	• Plural imperfective verb • Prepositional proclitics • Numated adverbials	• 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	• 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 [-to/that]</i>)	• MSA vocabulary - SAMER I	
7-zay	3	Intermediate Low	≤10	• Six-syllable or more words • Verbs/nouns with weak final letters	• Dual perfective verb • Dual imperfective verb • Singular imperative verb • Enclitics: dual pronoun • Broken plurals • Waw of oath	• Adverbial accusative (time and place adverbs) • Circumstantial accusative • Interrogative particle <i>hal</i>	• High frequency MSA vocabulary - SAMER II	• Some symbolism, or not everything is stated directly in the sentence.
8-ha		Intermediate Low	≤11		• 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)	• Absolute object (emphasizing the verb) • Object of purpose • Object of accompaniment • Verbal sentence with two direct objects	• MSA vocabulary - SAMER I and II • Negation particles • Numbers: 1,001-1,000,000	• Some symbolism that requires the reader to seek help to understand the idea.
9-ta	Intermediate Mid	≤12		• Dual imperative verb • Interrogative Hamza • Ba of oath • Oath: The particle of oath, the object of the oath, and the answer to the oath		• Vocative	• Vocabulary describing positive and negative emotional and mood states like joy, happiness, anger, regret, sorrow	• Some symbolism at the event level in the sentence that the reader understands through prior knowledge.
10-ya	4	Intermediate Mid	≤15		• Passive voice	• <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	• Singular relative pronouns • Verbal particles <i>qad</i> and <i>laqad</i> • Preposition-Conjunctions: <i>nimma</i> , <i>fima</i> ...	
11-kaf		Intermediate High	≤20		• Acting derivatives (e.g., the active participle)	• Nominal sentence with a nominal predicate • False <i>idafa</i> (tall in stature)	• Dual and plural relative pronouns	• A degree of symbolism and a need for prior knowledge to understand the meaning of the sentence.
12-lam	5	Advanced Low			• Diminutive form	• Parentheticals (explanation, blessing) • Exception • Exclusivity • Apposition (e.g., partitive or containing) • Specification (<i>tamiyiz</i> construction)	• MSA vocabulary - Samer III • Frozen Verbs (e.g., <i>Amiyin</i> Amen) • Numbers: > 1,000,000 • Five Nouns: <i>Dhu</i> (possession nominal) • Interjections: <i>bala</i> , <i>Ajal</i> , etc.	
13-mim	6-7	Advanced Mid			• Energetic mood (emphatic <i>nun</i>) • Ta of oath	• Conditional sentences • Jussive particle <i>lamma</i> (not yet)	• Words describing deep psychological states like depression, loss, psychological alertness • Use of coined, uncommon words • Abbreviations (e.g., LLC)	• Symbolic ideas and deeper meanings, especially in terms of the psychological dimension of characters/events.
14-nun	8-9	Advanced High				• Semantic emphasis • Praise and dispraise • <i>Masdar 'an</i> clause as a subject • Exclamatory form: <comparative adjective> <i>bih min</i>	• MSA vocabulary - SAMER IV • General legal, scientific, religious, political vocabulary, etc. • Five Nouns: <i>fw</i> , <i>Hmw</i>	• Local cultural expressions that may not be understood by those outside the
15-sin	10-11	Superior Low				• Uncommon constructions that are ambiguous and need diacritization for clarification	• Specialized vocabulary that requires understanding the concept/idea to comprehend it • Shortening in proper names (e.g., <i>fatim</i> for <i>fatima</i>)	• Symbolic, abstract, scientific, or poetic ideas that require prior linguistic and cognitive knowledge to understand.
16-ayn	12	Superior Mid					• MSA vocabulary - SAMER V • Specialized and highly elevated Arabic vocabulary. • Vocabulary mostly distant from dialects.	
17-fa	University Year 1-2	Superior High					• Scientific and heritage vocabulary not in use today, but familiar to a novice specialist	
18-sad	University Year 3-4	Distinguished					• Scientific and heritage vocabulary not in use today, but familiar to a specialist	
19-qaf	Specialist	Distinguished+					• 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' >> (ﻱ) • Offensive words >> (ﻉ) • Error in diacritics at the beginning of the sentence >> (ﻱ) • Dotted Yaa missing at the end of the word >> (ﻱ) 		

A.3 Annotation Interface

Figure 2 illustrates the Google Sheets interface used for annotation. The first two columns contain the sentence and its assigned readability level. The remaining columns correspond to the possible linguistic reasons underlying the assigned level, as defined in the BAREC annotation guidelines (Figure A.2). Annotators select one or more applicable reasons for each sentence. An additional notes column is provided for flagging ambiguous or problematic cases and for recording supplementary comments.

To facilitate the annotation process, the data was organized into batches of approximately 100 sentences, as summarized in Table 9, which reports the distribution of readability levels across all batches. This batching strategy ensured balanced coverage of levels while keeping annotation units manageable for annotators. Batches were then assigned to three expert annotators following a partially overlapping design, as shown in Table 8. Several batches were intentionally annotated by more than one annotator to enable the measurement of inter-annotator agreement (IAA), while the remaining batches were distributed to ensure full coverage of the dataset. Together, these tables illustrate both the level-wise composition of the annotation batches and the annotation assignment strategy used to support reliable and reproducible reasoning annotations.

B Modeling linguistic Phenomena

Feature extraction was conducted using a combination of CAMEL Tools (Obeid et al., 2020), regular expressions, external lexicons, and custom Python scripts. Below is a breakdown of each feature group and the extraction methods used:

• Morphological Features

- **Number of prefixes, suffixes, and clitics:** Extracted using CAMEL Tools’ morphological disambiguator. Each token was decomposed into its base form and affixes, and counts were aggregated per sentence.
- **Verb tense and voice (e.g., passive, active):** Identified using the POS tags and morphological features provided by CAMEL Tools.

- **Use of different forms (broken plurals, feminine plurals):** Detected using morphological patterns and specific tag combinations (e.g., singular form and plural num) from CAMEL analysis.

Table 10 shows the specific rules for all morphological features.

• Syntactic Features

A set of rule-based syntactic features was developed using dependency parsing outputs using the CamelParser (Elshabrawy et al., 2023). A dependency parse was first used to construct syntactic trees for each sentence, allowing for the identification of grammatical relations between words. From these structures, a set of binary features was extracted to reflect the presence or absence of key syntactic phenomena. Table 11 shows the specific rules for all Syntactic features.

• Content-Based Features

To automatically estimate this content complexity, a sentence-level classifier was developed by fine-tuning an AraBERT model (Antoun et al.). The model was trained to predict one of the eight content levels defined in the guidelines, treating this as a multi-class classification task. These predicted levels were then included as features in the broader feature set used for readability prediction.

• word/syllable counts

- **Word count:** Computed as the number of unique words in a sentence, ignoring repetitions, or punctuation.
- **syllable count:** The number of syllables in each word is computed by incorporating morphological and phonetic information. The CAPHI (consonant-vowel pattern) representation, the diacritized form of the word, and morphological prefix annotations are used for a more accurate count of syllables. The CAPHI string is tokenized and scanned for vowel segments, each indicating a potential syllable. Specific linguistic rules are applied to refine the syllable count:
 - * The final vowels are excluded if it is a diacritic (حركات الإعراب).

970	* Morphological prefixes such as the	difficult) level among all matched pairs. Experi-	1018
971	definite article (ال التعريف) and con-	ments were conducted using all three thresh-	1019
972	junction (واو عاطفة) are excluded, as	old variants, and the version yielding - 1% er-	1020
973	they do not contribute to the core syl-	ror margin- the best performance was selected	1021
974	labic structure of the main word.	for use in the final model.	1022
975	* In the absence of CAPHI informa-	In addition to the data-driven vocabulary ex-	1023
976	tion, syllables are counted by identi-	tracted from the training set, it was observed	1024
977	fying diacritic characters correspond-	that expanding the lexical coverage further	1025
978	ing to short vowels within the dia-	improved performance. The BAREC annota-	1026
979	critized form of the word.	tion guidelines specifically reference certain	1027
980		levels from the SAMER readability lexicon	1028
981	• Vocabulary-based Features	(Al Khalil et al., 2020) as indicative of vo-	1029
982	Vocabulary was handled in three different	cabulary difficulty. To incorporate this, the	1030
983	ways to estimate the lexical difficulty of sen-	SAMER lexicon was used to augment the ex-	1031
984	tences. Firstly, to estimate the vocabulary dif-	isting vocabulary-level dictionary. Lemma-	1032
985	iculty of a sentence, a level-based vocabu-	POS pairs from SAMER were assigned levels	1033
986	lary scoring system was constructed using the	in accordance with the BAREC guidelines,	1034
987	training set of the BAREC dataset. The pro-	thereby enriching the vocabulary feature set	1035
988	cess began by extracting all lemma-Part of	with structured, curriculum-aligned informa-	1036
989	Speech (POS) pairs from the training data us-	tion.	1037
990	ing the cameltools disambiguator. For each	To introduce dialectal sensitivity-also high-	1038
991	pair, the number of occurrences was counted	lighted in the BAREC guidelines-a supple-	1039
992	across all 19 BAREC readability levels. This	mentary lexicon from the BAREC project was	1040
993	allowed for identifying the earliest level at	utilized. This lexicon consists of approxi-	1041
994	which each lemma-POS pair appeared in the	mately 5,000 annotated words, each marked	1042
995	corpus.	with a dialectal match indicator. Although this	1043
996	To account for annotation noise or occasional	represents a relatively small subset of the over-	1044
997	use of advanced vocabulary in lower levels,	all vocabulary, it introduces an important di-	1045
998	three variants of vocabulary level assignment	mension of variation and adds a foundational	1046
999	were considered:	layer of dialectal awareness to the feature set.	1047
1000	– Strict : The lowest level at which the	This layer enables the model to distinguish	1048
1001	lemma-POS pair appeared.	between vocabulary that overlaps across Mod-	1049
1002	– Relaxed (1%) : The lowest level where	ern Standard Arabic and dialects versus vo-	1050
1003	the pair appeared, allowing for a 1% er-	cabulary that exists only in dialectal usage.	1051
1004	ror margin of frequency across levels.	Words that are common across both MSA	1052
1005	– Relaxed (2%) : Similar to the above but	and dialects-such as 'chair' (كرسي), which ap-	1053
1006	with a 2% margin.	pears consistently in both-are typically intro-	1054
1007	These thresholds introduced flexibility, ensur-	duced at earlier reading levels and thus ranked	1055
1008	ing that a few early occurrences of complex	lower in complexity. In contrast, words like	1056
1009	vocabulary in lower-level sentences did not	'window', which differ in MSA and dialect-	1057
1010	skew the overall difficulty estimation.	al forms (e.g., نافذة vs. شباك), are treated as	1058
1011	Once each lemma-POS pair was associated	more complex and are ranked at higher read-	1059
1012	with a level, a vocabulary-level dictionary was	ability levels. Incorporating this information	1060
1013	constructed containing all lemma-POS pairs	allows the model to better reflect the lexical	1061
1014	from the training data along with their as-	difficulty that dialectal divergence introduces,	1062
1015	signed difficulty levels. New input sentences	especially for learners who are trained primar-	1063
1016	were transformed into lists of lemma-POS	ily on MSA vocabulary.	1064
1017	pairs, and for each sentence, the vocabulary	In addition to the above features, the barec	1065
	level was defined as the highest (i.e., most dif-	dataset specifies certain closed groups of	1066
		vocabs that can be identified using the	1067

A1	A2	A3
1	7	13
2	2	2
3	8	14
9	9	9
4	10	15
16	16	16
5	11	17
6	12	18

Table 8: Assignment of annotation batches to annotators A1, A2, and A3.

Cameltools disambiguator, these are shown in table 12.

B.1 Linguistic codes and their associated readability levels

Table 13 presents the complete set of linguistic features defined in the BAREC framework and used by the proposed feature-based model. Each feature is represented by a unique code and is associated with a specific readability level, reflecting the guideline followed by human annotators. The table serves as a reference for the deterministic mapping between linguistic phenomena and readability levels employed by the system.

C License

We list below the licenses of the data and tools used in this work, all of which are employed in accordance with their intended use.

- SAMER Lexicon (Al Khalil et al., 2020): Non Commercial Use - NYUAD License.⁴
- BAREC Corpus (Elmadani et al., 2025a): Creative Commons Attribution Share Alike 4.0
- CAMEL Tools (Obeid et al., 2020) and CAMELBERT (Inoue et al., 2021): MIT License

⁴<https://camel.abudhabi.nyu.edu/samer-readability-lexicon/>

Sentence	word count	Spelling and Phonology		Morphology		Syntax		Vocabulary	Content	comment
		عدد الكلمات	تهجئة/إملاء	تصريف واشتقاق	تركيب نحوية	مفردات	فكرة / محتوى			
في رأيك ماذا كان خُطبتُ داني فأر المُرُوج؟	4 (صف 4)					كان وأخواتها				
يَتَعَرَّفُ الحُرُوكَةُ الفُصَيْرَةَ (الْفَتْحَةُ) وَيَنْطِقُهَا نَطْقًا-3-صَحِيحًا	5 (صف 5)					جمل اعتراضية (تفسير- دعاء...)				
سلوكي مسؤوليتي	2 (صف 2)		كلمات من ٥ مقاطع (بد...)							
مَنْ هُمْ أَوْلُو العِزْمِ مِنَ الرُّسُلِ؟	5 (صف 5)						مفردات فصيحة - ...			
ماذا ترى في الصُّورِ؟	2 (صف 2)		أفعال/أسماء ممتدة الآخر							
يَسْتَنْتِجُ فَضْلَ العِلْمِ وَأَهْمِيَّةَ القِراءَةِ؟	4 (صف 4)					إضافة خيالية (لفظية) طول القامة				
تحدث الملاك أيضًا إلى الطفل	3 (صف 3)					المفعول المطلق				

Figure 2: Annotation interface

Batch/Level	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	Total
1	1	0	1	0	0	0	1	0	0	0	0	0	0	1	0	0	0	0	4
2	1	0	1	0	1	1	1	1	1	1	1	1	1	0	0	1	0	0	12
3	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	1	1	34
4	1	1	1	1	0	0	1	1	1	1	0	0	0	0	0	0	1	1	10
5	2	3	2	3	2	3	2	3	2	3	3	3	3	3	3	3	3	3	49
6	3	3	4	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	55
7	7	7	7	6	8	7	7	7	8	7	8	8	8	8	8	8	8	7	134
8	7	7	7	7	7	7	7	7	7	7	7	8	8	8	8	8	7	7	130
9	3	3	3	3	4	4	3	3	4	3	3	3	3	3	3	3	5	5	61
10	11	12	11	12	11	12	12	12	11	12	12	12	11	12	11	12	11	12	209
11	5	5	5	5	5	4	5	5	5	5	5	5	5	5	5	5	5	5	89
12	24	24	24	24	23	24	23	23	23	23	23	23	24	24	24	23	24	23	423
13	5	5	6	5	5	5	6	5	5	5	6	5	5	5	5	5	5	5	93
14	19	19	19	18	19	18	18	19	19	19	18	19	19	18	19	19	19	19	337
15	4	4	4	5	5	4	5	4	4	4	4	4	4	4	4	4	4	4	75
16	3	2	2	3	2	3	2	3	3	3	3	2	3	2	2	3	3	3	47
17	1	2	1	2	1	1	1	1	1	1	1	2	1	2	1	1	1	2	23
18	1	0	0	0	0	0	0	0	0	0	1	0	0	0	1	0	0	0	3
19	0	0	0	1	0	1	1	1	0	0	0	0	0	1	0	0	0	0	5
Total	100	99	100	100	98	99	100	100	99	99	100	99	100	101	99	100	100	100	1793

Table 9: Distribution of sentences per annotation batch and readability level.

Feature	Feature (Arabic)	Rule
Singular imperfective verb	الفعل المضارع المفرد	num=s, asp=i, pos=verb
Proclitic: Definite article Al+	سوابق: ال التعريف	prc0=Al_det
Proclitic: Conjunction wa+	سوابق: واو العطف	prc2=wa_conj
Enclitic: First person singular pronoun	لواحق: ضمير المتكلم المفرد المتصل	enc0=1s_pron / 1s_poss / 1s_dobj
Plural imperfective verb	الفعل المضارع الجمع	pos=verb, asp=i, num=p
Prepositional proclitics	سوابق: حروف جر متصلة	prc1=bi_prep / li_prep / ka_prep
Enclitic: singular and plural pronouns	لواحق: ضمير متصل مفرد أو جمع	enc0 in [1p_dobj, ..., 3p_pron]
Dual (nouns and adjectives)	الثنى في الأسماء والصفات	num=d, pos=noun / adj
Sound feminine plural	جمع المؤنث السالم	form_num=p, form_gen=f, pos=noun / adj
Broken plurals	جمع التكسير	pos=noun / adj, form_num=s, num=p
Passive voice	المبني للمجهول	vox=p, pos=verb
Singular and plural perfective verb	الفعل الماضي المفرد والجمع	pos=verb, asp=p, num=s / p
Sound masculine plural	جمع مذکر سالم	form_gen=m, form_num=p, pos=noun / adj
Dual perfective verb	الفعل الماضي الثنى	pos=verb, asp=p, num=d
Dual imperfective verb	الفعل المضارع الثنى	pos=verb, asp=i, num=d
Singular imperative verb	فعل الأمر المفرد	pos=verb, asp=c, num=s
Enclitic: dual pronoun	لواحق: ضمير الثنى المتصل	enc0=[2d_dobj, ..., 3d_pron]
Waw of oath	واو القسم (والله)	prc2=wa_prep and lex ∈ qassam_lex
Plural imperative verb	فعل الأمر الجمع	pos=verb, asp=c, num=p
Nun of feminine plural (nūn al-niswā)	نون النسوة في الأسماء والأفعال	ends with noon and 3FP is present
Conjunctions / connectives	(ثم، حتى، أو، أم، لكن، أما)	lex ∈ connective_lex
Dual imperative verb	فعل الأمر للثنى	pos=verb, asp=c, num=d
Interrogative Alif	أداة الاستفهام: أ (أسمعت؟)	prc3 in [>a_ques, A_ques]
Ba of oath	باء القسم	prc1=bi_prep and lex ∈ qassam_lex
Ta of oath	تاء القسم	prc1=ta_prep and lex ∈ qassam_lex

Table 10: Morphological features and extraction rules derived from the BAREC annotation guidelines.

Feature	Feature (Arabic)	Rule
Nominal sentence	الجملة الاسمية	parent \neq inna and sisters; has a child with dependency relation SBJ; POS \neq VRB
Verbal sentence without direct object	جملة فعلية بدون مفعول به	parent = verb; no dependency relation OBJ
Preposition and object	جار + مجرور	parent POS: PRT; pos=prep in FEATS; has a child with dependency relation OBJ
Verbal sentence with one nominal direct object	جملة فعلية مع مفعول به واحد اسم	parent = verb; has exactly one child with dependency relation OBJ
Sentence with two verbs	جملة فيها فعلين	verb count \geq 2
Verbal sentence with clausal object (Masdar <i>an</i>)	جملة فعلية مفعولها أن المصدرية	token lemma = أن; POS = PRT; has a child with POS = VRB, deprel = OBJ, asp=i
Verbal sentence with two direct objects	جملة فعلية تتعدى إلى مفعولين	parent POS = VRB; two children with dependency relation OBJ
Vocative	النادى	parent POS = PRT; FEATS include pos=part_voc; has child with deprel OBJ
Inna and its sisters	إن وأخواتها	parent lemma in inna set; has a child with dependency relation PRD
Kana and its sisters	كان وأخواته	parent POS = VRB; lemma in kana set; has a child with dependency relation PRD
Preposed predicate, postponed subject	الخبر المقدم والمبتدأ المؤخر	dependency relation = SBJ; child POS \neq VRB; parent index < child index
Nominal sentence with nominal predicate	جملة اسمية خبرها جملة اسمية	sentence does not start with a verb; contains child with deprel TPC
False idafa (attributive)	إضافة خيالية (لفظية)	parent POS = NOM; pos=adj in FEATS; has child with deprel IDF
Exception	استثناء	POS = part_restrict
True idafa	إضافة حقيقية	parent POS = NOM; pos=adj not in FEATS; has child with deprel IDF
Noun + adjective	صفة وموصوف	adjacent NOM followed by ADJ (or noun with pos=adj); agreement in gender/number
Basic interrogatives	أدوات استفهام أساسية	lemma in {كيف, ما, أين, من, متى, ماذا}; POS indicates interrogative
Interrogative particle <i>hal</i>	أداة الاستفهام هل	lemma = هل; POS = part_interrog

Table 11: Syntactic features and extraction rules derived from the BAREC annotation guidelines.

Vocabulary	Vocabulary (Arabic)
Proper noun	اسم علم
Personal pronouns (non-clitics)	ضمير منفصل
Singular demonstrative pronoun	اسم الإشارة المفرد
Dual and plural demonstrative pronoun	اسم إشارة مثنى، جمع
Adjective	صفة
Verb	فعل
Prepositions	حروف الجر
Negation particles	أحرف النفي
Singular relative pronouns	أسماء الوصل المفردة
Dual and plural relative pronouns	أسماء الوصل المثنى والجمع
Father / Brother nouns	أبو - أخو
Ordinal numbers (written in words)	العدد الترتيبي
Cardinal numbers (written in words)	العدد الأصلي بالأحرف
Arabic or Indic digits (1-10)	الأرقام (العربية أو الهندية) من ١ إلى ١٠
Arabic or Indic digits (11-100)	الأرقام (العربية أو الهندية) من ١١ إلى ١٠٠
Arabic or Indic digits (101-1,000)	الأرقام (العربية أو الهندية) من ١٠١ إلى ١٠٠٠
Arabic or Indic digits (1,001-1,000,000)	الأرقام (العربية أو الهندية) من ١٠٠١ إلى ١٠٠٠٠٠٠

Table 12: Vocabulary feature levels in English and Arabic derived from the BAREC guidelines.

Table 13: Linguistic phenomena and their associated code and readability levels used for our feature-based model

Arabic Description	Code	Level
عدد المقاطع	O	–
كلمات تستخدم مد الألف (مثل أكل، أسف)	O4-1	4
أفعال معتلة الآخر	O7-1	7
عدد الكلمات	WC	–
مفرد مضارع (مذكر، مؤنث) (متكلم، مخاطب، غائب)	M1-1	1
سوابق: ال التعريف	M3-1	3
سوابق: واو العطف	M3-2	3
لواحق: ضمير المتكلم المفرد المتصل	M3-3	3
الفعل المضارع الجمع	M4-1	4
سوابق: حروف جر متصلة (ب+ ل+ ك+)	M4-2	4
ظرف منون	M4-3	4
لواحق: ضمير متصل مفرد أو جمع	M5-1	5
المثنى (في الأسماء والصفات)	M5-2	5
جمع المؤنث السالم في الأسماء والصفات	M5-3	5
الفعل الماضي المفرد والجمع	M6-1	6
جمع مذكر سالم	M6-2	6
الفعل الماضي المثنى	M7-1	7
الفعل المضارع المثنى	M7-1	7
فعل الأمر المفرد	M7-3	7
لواحق: ضمير المثنى المتصل	M7-4	7
جمع التكسير	M7-5	7
واو القسم (والله)	M7-6	7
فعل الأمر الجمع	M8-1	8
نون النسوة في الأسماء والأفعال (انتظرنَ دورهنَ)	M8-2	8
سوابق أخرى: سين الاستقبال، واو الاستئناف، فاء العطف	M8-3	8
(ثم، حتى، أو، أم، لكن، أمّا)	M8-4	8
فعل الأمر للمثنى	M9-1	9
أداة الاستفهام: أ (أسمعت؟)	M9-2	9
باء القسم	M9-3	9
القسم: أداة القسم والمقسم به وجواب القسم	M9-4	9
المبني للمجهول	M10-1	10
المشتقات على أنواعها	M11-1	11
التصغير	M12-1	12
نون التوكيد	M13-1	13

Continued on next page

Arabic Description	Code	Level
تاء القسم	M13-2	13
آخر - اشتقاق	M-O	-
كلمة واحدة (لا تراكيب نحوية ممكنة)	S1-1	1
إضافة حقيقية (باب بيت / بيت مريم)	S2-1	2
جملة اسمية	S2-2	2
صفة وموصوف	S2-3	2
بدل كل (صديقي أحمد)	S3-1	3
بدل إشارة (هذا البيت)	S3-2	3
جملة فعلية بدون مفعول به	S4-1	4
جار ومجرور	S4-2	4
جملة فعلية مع مفعول به واحد اسم	S5-1	5
جمل معطوفة	S5-2	5
أدوات استفهام أساسية: ماذا، متى، من، أين، ما، كيف	S5-3	5
صيغة التعجب ما أفعل	S5-4	5
جملة فيها فعلين	S6-1	6
مفعول فيه (ظروف زمان ومكان)	S7-1	7
حال	S7-2	7
أداة الاستفهام هل	S7-3	7
المفعول المطلق	S8-1	8
المفعول لأجله	S8-2	8
المفعول معه	S8-3	8
جملة فعلية تتعدى إلى مفعولين	S8-4	8
النادى	S9-1	9
إن وأخواتها	S10-1	10
كان وأخواتها	S10-2	10
خير مقدم / مبتدأ مؤخر	S10-3	10
النعنة (حدثني ... قال)	S10-4	10
رُبّ (حرف جر شبهه بالزائد)	S10-5	10
جملة الصلة، جملة الصفة، جملة الحال، جملة المفعول به	S10-6	10
جملة اسمية خبرها جملة اسمية	S11-1	11
إضافة خيالية (لفظية)	S11-2	11
جمل اعتراضية (تفسير، دعاء...)	S12-1	12
استثناء	S12-2	12
حصر	S12-3	12
بدل غير بدل الكل	S12-4	12
تمييز	S12-5	12
الجمل الشرطية	S13-1	13

Continued on next page

Arabic Description	Code	Level
لما حرف الجزم	S13-2	13
التوكيد المعنوي	S14-1	14
المدح والذم	S14-2	14
جملة أن المصدرية في محل رفع مبتدأ	S14-3	14
صيغة التعجب أفعل به	S14-4	14
تراكيب غير متداولة تحتاج إلى التشكيل الإعرابي لفكها	S15	15
آخر - تراكيب نحوية	S-O	-
اسم علم (متداول بسيط تركيبياً)	V1-1	1
ضمير منفصل	V1-2	1
اسم جنس	V1-3	1
مفردات متطابقة مع العامية - سامر ١	V1-4	1
الأرقام (العربية أو الهندية) ١٠--١	V1-5	1
أبو - أخو	V2-1	2
مفردات متشابهة مع العامية - سامر ١	V2-2	2
صفة	V2-3	2
فعل	V2-4	2
العدد الأصلي بالأحرف	V2-5	2
اسم الإشارة المفرد	V3-1	3
مفردات فصيحة شائعة - سامر ١	V3-2	3
الأرقام (العربية أو الهندية) ١٠٠-١١	V3-3	3
حروف الجر	V4-1	4
العدد الترتيبي	V5-1	5
الأرقام (العربية أو الهندية) ١٠٠٠-١٠١	V5-2	5
اسم إشارة مثنى، جمع	V5-3	5
مفردات فصيحة - سامر ١	V6-1	6
مفردات فصيحة شائعة - سامر ٢	V7-1	7
مفردات فصيحة - سامر ١ و سامر ٢	V8-1	8
أحرف النفي	V8-2	8
الأرقام (العربية أو الهندية) ١٠٠٠٠٠٠-١٠٠١	V8-3	8
مفردات تصف حالات مزاجية وشعورية إيجابية وسلبية	V9-1	9
أسماء الوصل المفردة	V10-1	10
(قد - لقد)	V10-2	10
(مما - مما - عمّا - عمّ - علام - فيم - إلام - بم)	V10-3	10
أسماء الوصل المثنى والجمع	V11-1	11
مفردات فصيحة - سامر ٣	V12-1	12
اسم الفعل: إيّه، صّه، رويدك ...	V12-2	12

Continued on next page

Arabic Description	Code	Level
الأرقام (العربية أو الهندية) أكبر من ١٠٠٠٠٠٠٠٠	V12-3	12
ذو	V12-4	12
(بل - بلى - أجل - قطع)	V12-5	12
كلمات تصف حالات نفسية عميقة مثل الاكتئاب، الضياع	V13-1	13
استخدام كلمات منحوتة غير متداول...	V13-2	13
الرموز (شم).	V13-3	13
مفردات فصيحة - سامر ء	V14-1	14
مفردات قانونية، علمية، دينية، سياسية، غير متخصصة	V14-2	14
فو - حمو	V14-3	14
المفردات المتخصصة التي لا تكفي معرفة الكلمة وحدها لفهمها...	V15-1	15
الترخيم في أسماء العلم (مثلا أفاطم؟)	V15-2	15
مفردات فصيحة - سامر ه	V16-1	16
مفردات متخصصة ومفردات عربية عالية غير شائعة كثير...	V16-2	16
مفردات علمية وتراثية غير متداولة اليوم	V17	17
مفردات علمية وتراثية غير متداولة...	V18	18
مفردات علمية وتراثية غير متداولة اليوم...	V19	19
آخر - مفردات	V-O	-
لا رمزية في النص	C1-1	1
فكرة مباشرة وصريحة وحسية	C1-1	1
المحتوى من حياة القارئ	C5-1	5
لا رمزية في النص	C5-1	5
بعض الرمزية أو عدم التصريح المباشر...	C7-1	7
بعض الرمزية يحتاج معها القارئ إلى مساعدة...	C8-1	8
هناك شيء من الرمزية على مستوى الحدث ..	C9-1	9
هناك درجة من الرمزية وحاجة للمعرفة السابقة...	C11-1	11
أفكار رمزية ومعنى باطن خاصة على صعيد البعد ...	C13-1	13
تعايير ثقافية محلية قد لا يفهمها من لا يشترك في نفس الثقافة	C13-1	13
أفكار رمزية، مجردة، علمية، أو شعرية...	C15-1	15
آخر - محتوى	C-O	-