



The SAMER Arabic Text Simplification Corpus

Bashar Alhafni, Reem Hazim, Juan Piñeros Liberato, Muhamed Al Khalil, Nizar Habash

نحن يا سيدتي في زمن مضطرب لا يركد عجاجه، ولا تسكن سيوفه في أغمادها، بعد أن انحلَّت أواصر بني العباس، وأصبحت دولتهم أشلاءً ممزقة، يفترسها كل مفترس، ويُغير عليها كلُّ واثب.

We live, my lady, in a troubled age whose dirt-laden air does not settle, nor its swords rest in their scabbards, after the bonds of the house of Abbas dissolved, and their state became like torn limbs preyed upon by every predator and raided by every opportunist.

نحن يا سيدتي في زمن <mark>مضطرب</mark> لا <mark>يركد عجاجه</mark>، ولا تسكن سيوفه في أغمادها، بعد أن انحلَّت أواصر بني العباس، وأصبحت دولتهم أشلاءً ممزقة، يفترسها كل مفترس، <mark>ويُغير</mark> عليها كلُّ واثب.

We live, my lady, in a <u>troubled</u> age whose <u>dirt-laden air</u>does not <u>settle</u>, nor its swords rest in their <u>scabbards</u>, after the <u>bonds</u> of the house of Abbas <u>dissolved</u>, and their state became like torn <u>limbs</u> preyed upon by every predator and <u>raided</u> by every <u>opportunist</u>.

نحن يا سيدتي في زمن مضطرب لا يركد عجاجه، ولا تسكن سيوفه في <mark>أغمادها</mark>، بعد أن انحلَّت أواصر بني العباس، وأصبحت دولتهم أشلاءً ممزقة، يفترسها كل مفترس، <mark>ويُغير</mark> عليها كلُّ واثب.

We live, my lady, in a <u>troubled</u> age whose <u>dirt-laden air</u>does not <u>settle</u>, nor its swords rest in their <u>scabbards</u>, after the <u>bonds</u> of the house of Abbas <u>dissolved</u>, and their state became like torn <u>limbs</u> preyed upon by every predator and <u>raided</u> by every <u>opportunist</u>.

نحن يا سيدتي في زمن صعب لا يهدأ غباره، ولا تسكن سيوفه في جيوبها، بعد أن تقطعت صلات بني العباس، وأصبحت دولتهم قطعا ممزقة، يفترسها كل مفترس، ويهجم عليها كلُّ <mark>عدو</mark>

We live, my lady, in a <u>difficult</u> age whose <u>dust</u> does not <u>clear</u> and whose swords do no rest in their <u>sheaths</u>, after the <u>joints</u> of the house of Abbas <u>were cut off</u>, and their state became like torn body parts preyed upon by every predator and <u>attacked</u> by every <u>enemy</u>.

Text Simplification Approaches:

• Lexical Simplification

 Identifying and replacing difficult words with easier synonyms while preserving the information and meaning of the original text

• Syntactic Simplification

• Making complex syntactic constructions simpler. For example, replacing passive voice constructions with active voice

End-to-End Simplification

• Both lexical and syntactic simplification. Usually defined as a sequence-to-sequence task

Text Simplification Datasets:

- Mostly in English: LexMTurk (Horn et al., 2014), LSeval (De Belder and Moens, 2012), Newsela Corpus (Xu et al. 2015), BenchLS (Paetzold and Specia, 2016), TurkCorpus (Xu et al. 2016)
- No publicly available Arabic text simplification datasets

Text Simplification Datasets:

- Mostly in English: LexMTurk (Horn et al., 2014), LSeval (De Belder and Moens, 2012), Newsela Corpus (Xu et al. 2015), BenchLS (Paetzold and Specia, 2016), TurkCorpus (Xu et al. 2016)
- No publicly available Arabic text simplification datasets

In this work:

- We create the first manually annotated Arabic dataset for **lexical simplification**
- 159K words selected from 15 Arabic fiction novels
- Readability level annotations at both the document and word levels
- Two simplification parallels for each text targeting school-aged learners at two different readability levels

Roadmap

- Motivation
- Arabic Linguistic Facts
- The SAMER Project
- The SAMER Simplification Corpus
- Conclusion

Arabic Linguistics Facts

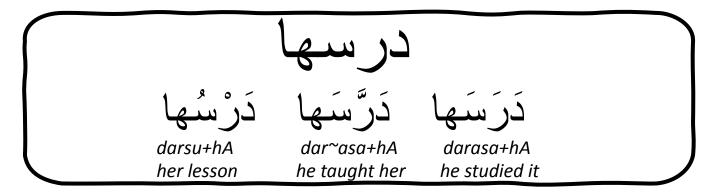
• Morphological richness:

- Inflects for gender, number, person, case, state, aspect, mood and voice
- 22,400 possible POS tags (48 POS tags for English)
- ~3 different core meanings (i.e., lemmas) per word
- Orthographic ambiguity:
 - Optional diacritics to specify short vowels and consonantal doubling
 - ~7 diacritizations per word

Arabic Linguistics Facts

• Morphological richness:

- Inflects for gender, number, person, case, state, aspect, mood and voice
- 22,400 possible POS tags (48 POS tags for English)
- ~3 different core meanings (i.e., lemmas) per word
- Orthographic ambiguity:
 - Optional diacritics to specify short vowels and consonantal doubling
 - ~7 diacritizations per word



Arabic Linguistics Facts

• Morphological richness:

- Inflects for gender, number, person, case, state, aspect, mood and voice
- 22,400 possible POS tags (48 POS tags for English)
- ~3 different core meanings (i.e., lemmas) per word
- Orthographic ambiguity:
 - Optional diacritics to specify short vowels and consonantal doubling
 - ~7 diacritizations per word



Roadmap

- Motivation
- Arabic Linguistic Facts
- The SAMER Project
- The SAMER Simplification Corpus
- Conclusion

The SAMER Project

- Simplification of Arabic Masterpieces for Extensive Reading (SAMER)
 - Funded by New York University Abu Dhabi Research Enhancement Fund
- **Objective**: to create standards and tools for the simplification of Arabic fiction to school-aged learners
- The project has developed the following open-source resources:
 - **Readability Lexicon:** A five-level readability scale at the word-level
 - **Readability Thesaurus:** An Arabic Readability-leveled Thesaurus for Arabic
 - Simplification Interface: A simplification interface platform as a Google Docs Add-on
 - Simplification Corpus: A text simplification corpus of Arabic novels targeting two readability levels

Roadmap

- Motivation
- Arabic Linguistic Facts
- The SAMER Project
- The SAMER Simplification Corpus
- Conclusion

Corpus Selection:

- 15 publicly available Arabic fiction novels from the Hindawi foundation
- 14 novels were published between 1865 and 1955. One famous novel published in the 12th century
- Extracted the first ~10K words from each novel (~159K words in total)
- Segmented chapters into paragraphs of 1500 words (4289 paragraphs)

Book Title	Author	Date	Para.	Words
حي بن يقظان Hayy Ibn Yaqzan	ابن طفیل Ibn Tufail	1150	213	9,962
غابة الحق The Forest of Truth	فر انسیس مر اش Francis Marrash	1865	241	10,081
لادیا <i>س</i> Ladiyas	أحمد شوقي Ahmed Shawqi	1899	300	9,846
المحالفة الثلاثية في المملكة الحيوانية The Tripartite Alliance of the Animal Kingdom	أمين ريحاني Ameen Rihani	1904	315	10,577
الملك كورش Cyrus the Great	زينب فواز Zaynab Fawwaz	1905	235	9,910
الأجنحة المتكسرة Broken Wings	جبر ان خلیل جبر ان Kahlil Gibran	1912	170	11,482
زينب Zaynab	محمد حسین هیکل Mohammed Hussein Heikal	1913	138	9,861
شجرة الدر The Pearl Tree	جرجي زيدان Jurji Zaydan	1914	296	12,230
إبر اهيم الكاتب Ibrahim Al-Katib	إبراهيم عبد القادر المازني Ibrahim Abd Al-Qadir Al-Mazini	1931	295	10,173
ڻورة في جهنم A Revolution in Hell	نقو لا حداد Niqula Haddad	1938	492	11,713
سارة Sara	عباس محمود العقاد Abbas Mahmoud Al-Aqqad	1938	300	10,079
فارس بني حمدان The Knight of Beni Hamdan	علي الجارم Ali Al-Jarem	1945	329	10,820
على باب زويلة On Bab Ziwaila	محمد سعيد العريان Mohammed Saeed Al-Aryan	1951	313	11,701
نماذج بشرية Human Examples	أحمد رضا حوحو Ahmad Rida Huhu	1955	443	10,080
هذا التاج This Crown	واصف البارودي Wasef Al-Baroudi	1955	209	10,750
			4,289	159,265

- **Goal:** simplify the Arabic fiction novels so that they can be targeted towards school-aged learners
- We use the five readability levels that were defined by Al Khalil et al. (2020) when creating the SAMER Lexicon (40K lemmas)

- **Goal:** simplify the Arabic fiction novels so that they can be targeted towards school-aged learners
- We use the five readability levels that were defined by Al Khalil et al. (2020) when creating the **SAMER Lexicon** (40K lemmas)

Level	Grade	Age	Examples
Level I	Grade 1	6	بَيْت، شَحِبَرَة، أَوْنَب، أَزْرَق، كَبير، صَنَعَ، أكَلَ، فَرِحَ، عَلى، لكِن
			house, tree, rabbit, blue, big, to make, to eat, to be happy, on, but
Level II	Grades 2-3	7-8	جَزيرة ، ذَهَب، سَنة ، داكِن، أُسْطَوانِي، صَعْب، خَدَعَ، كافأً، قُرْبَ، إذا
			island, gold, year, dark, cylindrical, difficult, to cheat, to reward, near, if
Level III	Grade 4-5	9-10	رِئة، مُتْحَف، مُعادَلة، مُمْكِن، مُوَحَّد، أَغْرى، نَدُرَ، لَدى، كَيْ، ما إنْحَتّى
			lung, museum, equation, possible, united, to entice, to be rare, with, for, no sooner than
Level IV	Grades 6-8	11-14	اِقْتِصاد، نُسْغ، طُمَأنينة، راقي، مُثْبَت، نَكَتَ، أَغْضى، إبّانَ، إنَّما، لَئَنْ
			economy, sap, tranquility, sophisticated, proven, to breach, to overlook, during, whereas, if (were)
Level V	Specialist	15 -	أَدَمة، قَسْطَرة، هَيْضة، مِطْياف، لَوْذَع، شُعَيّ، لَحا، ظَعَنَ، لَدُنْ، أَنَّى
			epidermis, catheterization, cholera, spectroscope, witty, bronchial, to denounce, to depart, with ($\approx chez$ in French), wherever

- **Goal:** simplify the Arabic fiction novels so that they can be targeted towards school-aged learners
- We use the five readability levels that were defined by Al Khalil et al. (2020) when creating the **SAMER Lexicon** (40K lemmas)
- We consider the document readability level to be equal to the highest readability level found among the words in the document
- Given the nature of the documents, all of them will have a readability of Level 5
- We focus on producing two simplified versions of each document: Level 4 and Level 3

Level	Grade	Age	Examples
Level I	Grade 1	6	بَيْت، شَعَبَرَة، أَرْنَب، أَزْرَق، كَبير، صَنَعَ، أكَلَ، فَرِحَ، عَلى، لكِن
			house, tree, rabbit, blue, big, to make, to eat, to be happy, on, but
Level II	Grades 2-3	7-8	جَزيرة، ذَهَب، سَنة، داكِن، أُسْطَوانِي، صَعْب، خَدَعَ، كافأً، قُرْبَ، إذا
			island, gold, year, dark, cylindrical, difficult, to cheat, to reward, near, if
Level III	Grade 4-5	9-10	رِئة، مُتْحَف، مُعادَلة، مُمْكِن، مُوَحَّد، أَغْرى، نَدُرَ، لَدى، كَيْ، ما إنْحَتّى
			lung, museum, equation, possible, united, to entice, to be rare, with, for, no sooner than
Level IV	Grades 6-8	11-14	اِقْتِصاد، نُسْغ، طْمَأْنينة، راقى، مُثْبَت، نَكَتَ، أَغْضى، إبّانَ، إنَّما، لَئَنْ
			economy, sap, tranquility, sophisticated, proven, to breach, to overlook, during, whereas, if (were)
Level V	Specialist	15 -	أَدَمة، قَسْطَرة، هَيْضة، مِطْياف، لَوْذَع، شُعَيّ، لَحا، ظَعَنَ، لَدُنْ، أَتّى
			epidermis, catheterization, cholera, spectroscope, witty, bronchial, to denounce, to depart, with ($\approx chez$ in French), wherever

Corpus Annotation:

- SAMER Simplification Add-on (Hazim et al. (2020))
- All documents were automatically labeled with their word-level readability by using a Python-API version of the add-on:
 - 1) Morphological disambiguation;
 - 2) Lemma and POS lookup in the SAMER lexicon

९ ५ ८ ८ ५ ५ १ 100% र । / र ∧	SAMER Readability Analysis
the sector of the sector states and the sect	Analyze Readability
• #بالله #١ #عليك #١ #لا #٢ #تطيلي #١ #يا #١ #ليلى	Doc Level Word Level
۲# <mark>فإن</mark> #۱#مما #۳#يثير #٤# <mark>شجون</mark> #۱#النفس، ۱#وِيزيد #۱#في #۲#ألم #۱#الحزين، #۲# <mark>أن</mark>	Modify Markup Show
۱ # يُدفع # ۱ # إلى # ۳ # العزاء # ۱ # والصبر # ۱ # بكلمات	# Success!
٤ # <mark>خاوية</mark> #٤ # <mark>متخاذلة</mark> #٣ # حفظها #١ # الناس	Loading This may take a moment
٣ لينثروها # ۱ # في # ۱ # كل # ٥ # <mark>مأتم.</mark> # ۲ # إن	Doc Level Analysis
	## Current Doc %
١# كل #١ # كلمة #١ #من #١ #هذه #١ #يا #١ #ليلي	# Names 5.7% Level 1 49.1%
in the second se	H Level 2 15.1%
٣#شعلة #٤# <mark>تؤجج</mark> #٣#وجدي، #٥# <mark>وتِضْطِر</mark> م #إ#في	# Level 3 17% Level 4 7.5%
۳ #فؤادى، #۲#إن #۱ #الحزن # • # حرم # • #قد سى	# Level 5 3.8%
	Unknown 1.9%
٣#يجب #٢ #أن #٣ ٣ تخشع #١ #أمامه #١ #الرءوس	# Total word count: 53
۲#بالصمت #۲#والاطراق.	# Target Level 5
	At Level 5 3 words 5.7%
	Below Level 5 50 words 94.3%

Chapters were loaded into 146 Google Docs

Clear

- Word highlights with different colors according to their readability levels
- Word-level markup (#<level>#)
- The add-on employs two additional readability levels: Level 0 (proper nouns), Level 6 (unknown words)

오 5 2 帚 冬 뤃 100% ▼ :	1 - 1 - 5	SAMER Readability Analysis	:
عليك ١٠ ٩ ٩ معيلي ١٠ معيلي	#۱#بالله #۱	Analyze Readability Doc Level Word Level	Clear
١ = في ٢ = ٢ = ألم ١ = ١ = الحزين، ٢ = ٢ = أن	#۱ #ويزيد #	Modify Markup Show Mini	w Hide imize Delete
ع # ٤ # <mark>متخاذلة</mark> # ٣ # حفظها # ١ # الناس	# ٤ # <mark>خاوية</mark> +	Success! Loading This may take a mom Doc Level Analysis	nent —
كلية #١ #مِن #١ #هذه #١ #يا #١ #ليلې	# ۱ # کل ۽ # ۱ #	Current Doc % Names 5.7% Level 1 49.1% Level 2 15.7%	
ري. #٢ #إن #١ # الحزن # • # حَرَّمُ # • # قُدْسي	#۳#فؤادي، #	Level 3 17% Level 4 7.5% Level 5 3.8% Unknown 1.9% Total word count: 53	
۳ <mark>۱۵</mark> #۲۴ محسع #۲۴ #۲۵۵۵ #۲۴ #۲۷ وس ۲۶ #والإطراق.	#۲ #بالصمت	Target Level 5 At Level 5 3 words 5.7% Below Level 5 50 words 94.	.3%
		By Token By Type	

Corpus Annotation:

- Word highlights with different colors according to their readability levels
- Word-level markup (#<level>#)
- The add-on employs two additional readability levels: Level 0 (proper nouns), Level 6 (unknown words)

SAMER Readability Analysis < 5 2 母 ☆ 気 100% ▼ :</p> 1 -# ١ # علىك # ١ # لا # ۲ # تطا 1# 6#1# Level Assia Level 2 #۱ #الحزين Level 3 Assian #*۳*#العزاء # ۱ #والصبر # ۱ # بكلمات 11#1# # ٤ # خاوية Level 4 Assia Level 5 #1# Jh *(verb)* : be inflamed, flare up Synonym (verb) أَحْتَرُق (verb) ٱلْتَهَب : Level 1 Level 2: ٱشْتَعَل (verb), #٣ # فؤادى، #٢ # ان (verb), تَنَجَّج (verb), تَرَقَّد (verb), ٱسْتَعَر (verb), تُحَرَّق #۱#امامه #۱#الرءوس #۲#ان #۳# تخشع Hypernyms غُيُّر ,(verb) تَبَدُّل (verb) بَدُّل (verb) تَحَوُّل :Level 1 (verb). Level 2: اَنْقَلَى (verb).

Support word substitution by displaying suggestions of related words from the Arabic WordNet (Black et al., 2006) (e.g., synonyms, hypernyms, and hyponyms)

Assign al

Assign all

Assign all

Assign all

Assign all

Level 3: آرتد (verb)

- Three professional linguists
- Guidelines:

	وليس جميع الحوادث والأحوال تساوي الدم	Not all events and conditions are worth human
	الإنساني الذي لا يوجد أثمن منه، ولا يجب	blood, for there is nothing more valuable than
Origin	مضارعة أو لنك الشعوب الذين بيادرون الي	blood, for there is nothing more valuable than that blood. Nor should one emulate those people
Loval	شن الذارات مفتلك بمصر معرمة بالصل أقل أرب -	who take to launching raids and slaughtering
Leve	٧ بعتريه، أو أن خرافة ٧ بيت إما في قعة	each other for the triflest most dismissible of
	م يعت به، او التي عراف م بيف تها في رفعه - التددن،	wants, or for the weakest of superstitions that has
	ונדאבט:	no refuge in civilization's domain.

- Three professional linguists
- Guidelines:
 - 1) Check the readability level assigned to each word
 - 2) Correct any readability leveling errors

	وليس جميع الحوادث والأحوال تساوي الدم	Not all events and conditions are worth human
	الإنساني الذي لا به حد أثمن منه، و لا يحب	blood, for there is nothing more valuable than that
Original	مضارعة أرازان الشمير بالذبن بداد بينا	blood. Nor should one emulate those people who
Level 5	شينه الغادات وفتك دوضوه دوضاحا أقل	take to launching ₂ raids and slaughtering each
Level 5	الدير لا يعتر بدي أر أدني في افتر لا يدتر إدا ف	other for the triflest most dismissible4 of wants3,
	ا رب: 3 يعد4 به، او ادنی خرافه لا بيت نها في	other for the triflest most dismissible ⁴ of wants ³ , or for the weakest of superstitions that has no
	رفعه التمدن؛	refuge in civilization's domain.

- Three professional linguists
- Guidelines:
 - 1) Check the readability level assigned to each word
 - 2) Correct any readability leveling errors
 - 3) If the document is of:
 - a) Level 5, simplify to Level 4 then to Level 3
 - b) Level 4, simplify to Level 3
 - c) Level 3, no simplification is needed
 - 4) Simplification involves minimal replacements, deletions, and insertions to reduce the readability level of the document, while maintaining overall meaning and grammaticality

()		
Original Level 5	الإنساني الذي لا يوجد أثمن منه، ولا يجب مضارعة1 أولئك الشعوب الذين يبادرون إلى شن2 الغارات وفتك بعضهم بعضا على أقل أرب3 لا يعتد4 به، أو أدنى خرافة لا بيت لها في	Not all events and conditions are worth human blood, for there is nothing more valuable than that blood. Nor should one emulate1 those people who take to launching2 raids and slaughtering each other for the triflest most dismissible4 of wants3 , or for the weakest of superstitions that has no refuge in civilization's domain.
Level 4	وليس جميع الحوادث والأحوال تساوي الدم الإنساني الذي لا يوجد أثمن منه، ولا يجب أن تشبه، أولنك الشعوب الذين يبادرون، إلى بدءء الغارات، وفتك, بعضهم بعضا على أقل هدف. لا يؤخذ، به، أو أدنى خرافة لا بيت لها في	Not all events and conditions are worth human blood, for there is nothing more valuable than that blood. Nor should one imitate ₁ those people who take to ₅ starting ₂ raids ₆ and slaughtering ₇ each other for the triflest most negligible ₄ of aims ₃ , or for the weakest of superstitions that has no refuge in civilization's ₉ domain ₈ .
Level 3	وليس جميع الحوادث والأحوال تساوي الدم الإنساني الذي لا يوجد أثمن منه، ولا يجب أن نشبه أولنك الشعوب الذين يسر عون إلى بدء الهجمات، وقتل، بعضهم بعضا على أقل هدف لا يزخذ به، أو أدنى خرافة لا بيت لها في منطقة التعديد،	Not all events and conditions are worth human blood, for there is nothing more valuable than that blood. Nor should one imitate those people who rush tos starting attackss and killing ⁷ each other for the triflest most negligible of aims, or for the weakest of superstitions that has no refuge in the ranges of progresss.

Inter-Annotator Agreement:

- Selected ~1300 words from each book (~20K words) to be double annotated
- Qualitative Check:
 - ~6.8% word-level mismatches when simplifying the original text to Level 4
 - **~13.2%** word-level mismatches when simplifying Level 4 text to Level 3
 - Most mismatches are due to different lexical simplification choices
 - Annotation mistakes were ~10% of all mismatches

- The three versions of the texts are comparable in size in terms of the number of words
- Different distributions of readability levels:

	Original		L4		L3	
L0	2,631	1.7%	5,212	3.3%	5,246	3.3%
L1	83,772	52.6%	90,232	56.5%	95,898	59.9%
L2	23,103	14.5%	26,297	16.5%	30,015	18.7%
L3	22,517	14.1%	24,630	15.4%	28,990	18.1%
L4	14,965	9.4%	13,306	8.3%	0	0.0%
L5	9,463	5.9%	0	0.0%	0	0.0%
L6	2,814	1.8%	0	0.0%	0	0.0%
Total	159,265	100%	159,677	100%	160,149	100%

- The three versions of the texts are comparable in size in terms of the number of words
- Different distributions of readability levels:
 - 8.8% shift from higher to lower readability when simplifying the original text to level 4 (L4)

	Original		L4		L3	
L0	2,631	1.7%	5,212	3.3%	5,246	3.3%
L1	83,772	52.6%	90,232	56.5%	95,898	59.9%
L2	23,103	14.5%	26,297	16.5%	30,015	18.7%
L3	22,517	14.1%	24,630	15.4%	28,990	18.1%
L4	14,965	9.4%	13,306	8.3%	0	0.0%
L5	9,463	5.9%	0	0.0%	0	0.0%
L6	2,814	1.8%	0	0.0%	0	0.0%
Total	159,265	100%	159,677	100%	160,149	100%

- The three versions of the texts are comparable in size in terms of the number of words
- Different distributions of readability levels:
 - 8.8% shift from higher to lower readability when simplifying the original text to level 4 (L4)
 - 8.3% shift from higher to lower readability when simplifying level 4 (L4) text to level 3 (L3)

	Original		L4		L3	
L0	2,631	1.7%	5,212	3.3%	5,246	3.3%
L1	83,772	52.6%	90,232	56.5%	95,898	59.9%
L2	23,103	14.5%	26,297	16.5%	30,015	18.7%
L3	22,517	14.1%	24,630	15.4%	28,990	18.1%
L4	14,965	9.4%	13,306	8.3%	0	0.0%
L5	9,463	5.9%	0	0.0%	0	0.0%
L6	2,814	1.8%	0	0.0%	0	0.0%
Total	159,265	100%	159,677	100%	160,149	100%

- The three versions of the texts are comparable in size in terms of the number of words
- Different distributions of readability levels:
 - 8.8% shift from higher to lower readability when simplifying the original text to level 4 (L4)

	Original		L4		L3	
L0	2,631	1.7%	5,212	3.3%	5,246	3.3%
L1	83,772	52.6%	90,232	56.5%	95,898	59.9%
L2	23,103	14.5%	26,297	16.5%	30,015	18.7%
L3	22,517	14.1%	24,630	15.4%	28,990	18.1%
L4	14,965	9.4%	13,306	8.3%	0	0.0%
L5	9,463	5.9%	0	0.0%	0	0.0%
L6	2,814	1.8%	0	0.0%	0	0.0%
Total	159,265	100%	159,677	100%	160,149	100%

- 8.3% shift from higher to lower readability when simplifying level 4 (L4) text to level 3 (L3)
- 17.1% in the overall readability levels from the original text to level 3 (L3)

Corpus Statistics (Transformations Statistics):

- Alignment between the original text and its Level 4 simplified version (Original-L4) and Level 4 and Level 3 (L4-L3)
- Edits in terms insertions, deletions, and replacements
- Each edit could have more than one word at a time

	Origina	I-L4	L4-L3		
No Change	152,214	95.5%	145,090	90.8%	
1-1	6,430	4.0%	13,508	8.5%	
1-m	337	0.2%	476	0.3%	
m-1	120	0.1%	235	0.1%	
Insert	209	0.1%	354	0.2%	
Delete	40	0.0%	122	0.1%	

Corpus Statistics (Transformations Statistics):

- Alignment between the original text and its Level 4 simplified version (Original-L4) and Level 4 and Level 3 (L4-L3)
- Edits in terms insertions, deletions, and replacements
- Each edit could have more than one word at a time

	Origina	I-L4	L4-L3		
No Change	152,214	95.5%	145,090	90.8%	
1-1	6,430	4.0%	13,508	8.5%	
1-m	337	0.2%	476	0.3%	
m-1	120	0.1%	235	0.1%	
Insert	209	0.1%	354	0.2%	
Delete	40	0.0%	122	0.1%	

Corpus Statistics (Transformations Statistics):

- Alignment between the original text and its Level 4 simplified version (Original-L4) and Level 4 and Level 3 (L4-L3)
- Edits in terms insertions, deletions, and replacements
- Each edit could have more than one word at a time

	Origina	I-L4	L4-L3		
No Change	152,214	95.5%	145,090	90.8%	
1-1	6,430	4.0%	13,508	8.5%	
1-m	337	0.2%	476	0.3%	
m-1	120	0.1%	235	0.1%	
Insert	209	0.1%	354	0.2%	
Delete	40	0.0%	122	0.1%	

Corpus Statistics (Fragments Statistics):

- Segmented the data based on punctuation
- Each fragment consisted of ~7.5 words on average

	n	n%	Example					
No Change	8,920	8,920		43.3%	42.20/	Original	ولكنها تعلمت في المدارس الفرنسية أيضا،	
No change					L4 & L3	But she studied in French school as well.		
		10 12.7%	Original	ونسمعها <mark>تندب وتنوح كالثكلي</mark> .				
Change in L4 only				And we hear her lamenting and wailing like a bereaved woman.				
Change in L4 only	2,010		12.770	12.770	,010 12.770	,010 12.770	L4 & L3	ونسمعها تبكي وتصرخ كفاقدة ابنها.
			L4 & L3	And we hear her crying and screaming like one who lost her son.				
Change in L3 only 6,		C	Original	يجب أن يترأس الجلسة،				
	6,369	30.9%	L4	He must preside over the session,				
Change in LS only	0,309	30.9%	L3	يجب أن يقود الجلسة،				
			LJ	He must lead the session,				
			Original	اُحدهما مس <i>ر</i> ج ملجم ،				
	Change in L4 & L3 2.704 13.1%		Original	One of them is saddled and bridled ₁ ,				
Change in L4 & L3		13.1%	L4	أحدهما مسرج مربوط،				
Change in L4 & L3	2,704	13.170		One of them is saddled ₂ and tied ₁ ,				
			L3	أحدهما معد مربوط،				
				One of them is readied ² and tied,				
Total	20,603	100.0%						

Corpus Statistics (Fragments Statistics):

- Segmented the data based on punctuation
- Each fragment consisted of ~7.5 words on average

	n	n%	Example														
No Change	8.920	43.3%	Original	ولكنها تعلمت في المدارس الفرنسية أيضا،													
No change	0,920		L4 & L3	But she studied in French school as well.													
	2 610	610 12.7%	Original	ونسمعها <mark>تندب وتنوح كالثكلي</mark> .													
Change in L4 only			12.7%	s10 12 7%	Onginai	And we hear her lamenting and wailing like a bereaved woman.											
Change in L4 only	2,010			L4 & L3	ونسمعها تبكي وتصرخ كفاقدة ابنها.												
			L4 α L3	And we hear her crying and screaming like one who lost her son.													
Change in L2 only 6 260		Orig	Original	يجب أن يترأس الجلسة،													
	6,369	30.9%	L4	He must preside over the session,													
Change in L5 only	Change in L3 only 6,369	30.9%	50.9%	يجب أن يقود الجلسة،													
																	L3
	Change in L4 & L3 2.704 13.1%	Or	Original	أحدهما مسرج <mark>ملجم</mark> ،													
			Original	One of them is saddled and bridled ₁ ,													
Change in L4 & L3		13.1%	L4	أحدهما مسرج مربوط،													
Change III L4 & L3 2,704 13.	2,704	13.170	L4	One of them is saddled ₂ and tied ₁ ,													
			L3	أحدهما معد مربوط،													
		L3	One of them is readied ² and tied,														
Total	20,603	100.0%															

Roadmap

- Motivation
- Arabic Linguistic Facts
- The SAMER Project
- The SAMER Simplification Corpus
- Conclusion

Conclusion & Future Work

- Presented the first manually annotated publicly available Arabic corpus for lexical simplification
- Our corpus includes readability level annotations at both the document and word levels
- Two simplified parallels for each text targeting learners are different readability levels
- Future work:
 - Extensions to other genres and domains
 - Models for readability assessment and text simplification

Thank you! Q&A



The SAMER Arabic Text Simplification Corpus

Bashar Alhafni, Reem Hazim, Juan Piñeros Liberato, Muhamed Al Khalil, Nizar Habash



