

arTenTen: Arabic Corpus and Word Sketches

Tressy Arts

Lexicographer
tressy.arts@gmail.com

Yonatan Belinkov

MIT, USA
belinkov@mit.edu

Nizar Habash

New York University Abu Dhabi, UAE
nizar.habash@nyu.edu

Adam Kilgarriff

Lexical Computing Ltd, UK
adam@lexmasterclass.com

Vit Suchomel

Masaryk Univ., Cz.,
Lexical Computing Ltd, UK
xsuchom2@fi.muni.cz

Abstract

We present arTenTen, a web-crawled corpus of Arabic, gathered in 2012. arTenTen consists of 5.8-billion words. A chunk of it has been lemmatized and part-of-speech (POS) tagged with the MADA tool and subsequently loaded into Sketch Engine, a leading corpus query tool, where it is open for all to use. We have also created 'word sketches': one-page, automatic, corpus-derived summaries of a word's grammatical and collocational behavior. We use examples to demonstrate what the corpus can show us regarding Arabic words and phrases and how this can support lexicography and inform linguistic research.

The article also presents the 'sketch grammar' (the basis for the word sketches) in detail, describes the process of building and processing the corpus, and considers the role of the corpus in additional research on Arabic.

1 Introduction

Without data, nothing. Corpora are critical resources for many types of language research, particularly at the grammatical and lexical levels. In this article, we present arTenTen, a web-crawled corpus of Arabic, gathered in 2012, and a member of the TenTen Corpus Family (Jakubíček *et al.* 2013). arTenTen comprises 5.8-billion words. Since 2003, the key resource for Arabic has been Arabic Gigaword.¹ It contains exclusively newswire text. arTenTen improves on Gigaword, for dictionary-editing and related purposes, by covering many more types of text. A 115-million word chunk has been tokenized, lemmatized and part-of-speech tagged with the leading Arabic processing toolset, MADA (Habash and Rambow 2005; Habash *et al.* 2009), and installed in the Sketch Engine (Kilgarriff *et al.* 2004), a leading corpus query tool, where it is available for all to investigate.² There have been other important efforts in creating large collections of Modern Standard Arabic text, such as the Corpus of Contemporary Arabic (al-Sulaiti and Atwell 2006), International Corpus of Arabic (Alansary *et al.* 2007) and the Leipzig University Arabic collection (Eckart *et al.* 2014). Zaghouani (2014) has also presented a survey of several freely available corpora. These various corpora come in a range of sizes, but all of them are smaller than arTenTen.

One feature of interest in the Sketch Engine is the 'word sketch', a one-page, automatically derived summary of a word's grammatical and collocational behavior. Word sketches have been in use for English lexicography since 1999 (Rundell and Kilgarriff 2002) and are now available for twenty languages. In section 2, we describe how word sketches (and two related reports; thesaurus and 'sketch diff') can be used to give a better understanding of the behavior of Arabic words and phrases.³

To provide word sketches, we must parse the corpus either with an external parser or with the Sketch

¹ Arabic Gigaword is created and distributed by the Linguistic Data Consortium (Graff 2003). It is regularly updated and is now in its fifth edition.

² <http://www.sketchengine.co.uk>

³ The methods and approach described here are similar to those used in the creation of the Oxford Arabic Dictionary (Arts *et al.* 2014)

Engine's built-in shallow parser, as here. For this process, we need a 'sketch grammar' for Arabic, which is presented in a tutorial-style introduction in section 3. Section 4 describes how arTenTen was created and prepared for the Sketch Engine. In section 5, we conclude with a summary and a brief discussion of future work.

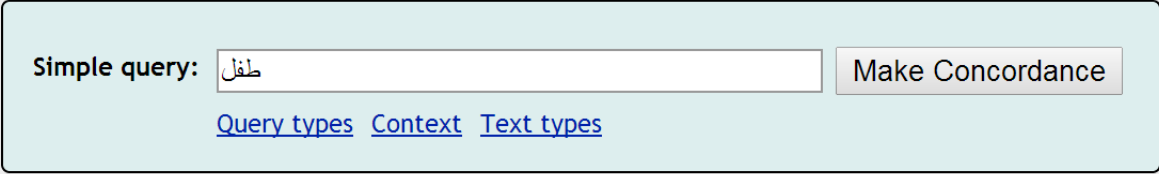
2 Using arTenTen in the Sketch Engine for language research

The Sketch Engine is in use for lexicography at four of the five UK dictionary publishers (Oxford University Press, Cambridge University Press, Collins, and Macmillan), at national institutes for Bulgarian, Czech, Dutch,⁴ Estonian, Irish,⁵ and Slovak, and for a range of teaching and research purposes at over 200 universities worldwide.

Before discussing the details of how we built the arTenTen corpus and annotated it, we provide several examples of its utility in the context of language research, e.g., for lexicography. This section is organized around the different functions available to the linguist using the Sketch Engine to study Arabic words in their context.

2.1 The Simple Concordance Query function

A Simple concordance query shows the word as it is used in different texts in the corpus. Figure 1 shows the query box, while Figure 2 shows its output. A simple search query for a word such as **طفل** (child) searches for the lemma as well as the string; so, the strings **الطفل** (the+child), **طفلهما** (child+their), **كالأطفال** (like+the+children), etc., are all retrieved.



Simple query:

[Query types](#) [Context](#) [Text types](#)

Figure 1: Simple concordance query

⁴ Dutch is an official language in both the Netherlands and Belgium (where it is also called Flemish), and the institute in question (INL) is a joint one from both countries.

⁵ Much of the development work for the Sketch Engine was undertaken under a contract from Foras na Gaeilge (the official body for the Irish language) in preparation for the creation of a new English-Irish dictionary (<http://www.focloir.ie>). Irish is spoken in both the Irish Republic and Northern Ireland (which is part of the UK), and Foras na Gaeilge is a joint institute of both countries.



Figure 2: The resulting concordance lines

2.2 The Frequency functions

The Sketch Engine interface provides easy access to tools for visualizing different aspects of the word frequency (see Figures 3 and 4). The Frequency Node⁶ forms function on the left hand menu (Figure 3) shows which of the returned forms are most frequent.

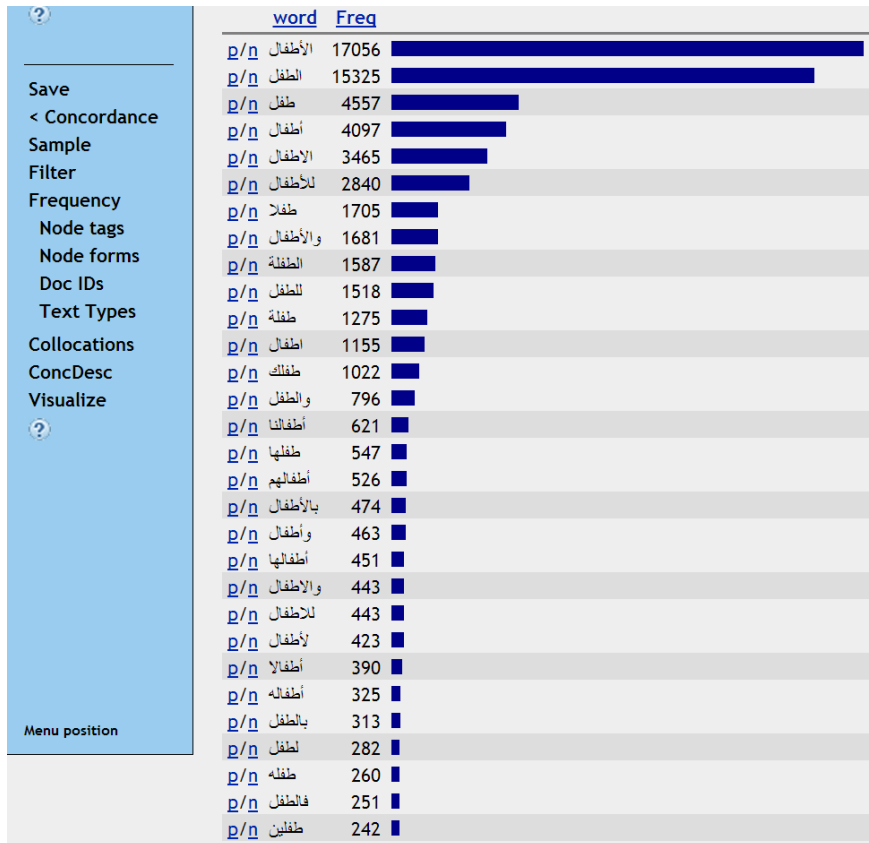


Figure 3: Frequency of node forms of طفل

The **p/n** links are for positive and negative examples. Clicking on **p** gives a concordance for the word form, while clicking on **n** gives the whole concordance *except* for the word form.

The Frequency Text Types function shows which top-level domain is most frequent (Figure 4).

Top level domain	Freq	Rel [%]
p/n com	38068	95.1
p/n net	14110	103.8
p/n org	10128	118.0
p/n ps	947	130.5
p/n sa	901	160.7
p/n info	744	62.1
p/n sy	435	126.5
p/n ae	357	138.6
p/n ws	338	54.4
p/n edu	305	612.1
p/n uk	284	87.4
p/n jo	271	103.8
p/n ma	267	132.0
p/n eg	256	71.8
p/n sd	208	75.5

Figure 4: Frequency list of domain extensions of sites that contain forms of طفل

Both hit counts and normalized figures are presented to account for the different quantities of material from different domains. If the word was equally frequent (per million words) in all of the domains, the figures in the fourth column would all be 100%. The bars are based on the normalized figures (with the height of the bar corresponding to the quantity of data). We see that طفل is frequent on .edu sites.

This utility is useful when researching regional differences. For example, making a frequency list for حَوْصَصَة (privatization), we see (Figure 5) that it is used almost exclusively in Moroccan and Algerian newspapers.

p/n	Second level domain	Freq	Rel [%]
p/n	sawt-alahrar.net	5	516.7
p/n	assif.info	4	1014.8
p/n	annahjaddimocrati.org	2	18471.3
p/n	wordpress.com	1	12.0
p/n	voltairenet.org	1	133.4
p/n	rifftoday.com	1	1759.3
p/n	odabasham.net	1	136.9
p/n	marxy.com	1	1279.2
p/n	kassioun.org	1	191.0
p/n	justgoo.com	1	544.6
p/n	essaha.info	1	1469.4
p/n	educpress.com	1	834.9
p/n	echoroukonline.com	1	867.2
p/n	djazairress.com	1	98.0

Figure 5: Frequency list of sites containing forms of *خوصصة*

2.3 The Word List function

The Word List function allows the user to make frequency lists of many varieties. Figures 6, 7 and 8 show the tops of frequency lists for word forms, lemmas and diacritized⁷ lemmas for the corpus.

Word list		Word list		Word list	
Corpus: arTenTen12 [sample 115M]		Corpus: arTenTen12 [sample 115M]		Corpus: arTenTen12 [sample 115M]	
Page	1	Go	Next >	Page	1
				Go	Next >
word	Freq	lemma	Freq	lemma_voc	Freq
في	3242280	في	3962066	فِي	3962066
من	2914934	من	3500214	مِنْ	3413373
على	1593477	على	2285678	عَلَى	2283548
أن	1184760	أن	2184612	أَنَّ	1332439
إلى	754664	الذي	1310294	الَّذِي	1310294
عن	738288	هذا	1245137	هَذَا	1245137
لا	659851	إلى	1231294	إِلَى	1231294
و	637527	كان	1102480	كَانَ	1102480
الله	629086	ما	1009041	مَا	1009041
ما	610949	لا	984894	لَا	984894
التي	585503	عن	927882	عَنْ	927455
هذا	518842	إن	899877	أَنَّ	850948
أو	453099	قال	746755	قَالَ	746755
الذي	416753	الله	722219	اللَّهِ	722219
ان	413353	ذلك	640003	ذَلِكَ	640003
مع	402313	و	638513	وَ	638513
هذه	402083	ل-	588106	ل-	588106
كان	361499	أو	545702	أَوْ	557135

Figures 6, 7, and 8: Frequency list of the whole corpus for word forms, lemmas and diacritized lemmas

2.4 The Word Sketch and Collocation Concordance functions

The Word Sketch function is invaluable for finding collocations. The word sketch for أخضر (green, Figure 9) shows expected collocates such as وأصفر (and yellow) and لون (color) but also the idiomatic الأخضر واليابس (literally "the green and the dry"). Clicking on the number after the collocate gives a concordance of the combination (Figure 10).

⁷ Diacritics and diacritization are often referred to as vowels and vocalization because the most common use of Arabic diacritics is to indicate short vowels. We use the more general term here to account for non-vowel diacritical marks, such as the consonant gemination marker, the shadda.

أخضر		arTenTen12 [sample 115M] freq =	
and/or	472 0.7	adjective-of	4865 4.6
ياهن	174 12.6	ضوء	456 10.48
أصفر	52 9.52	شاي	178 10.14
برتقالي	2 8.45	نون	300 9.42
أزرق	13 8.31	خط	355 8.94
بنفسجي	5 8.05	مسطح	61 8.54
وردي	6 7.88	منطقة	637 8.41
أحمر	77 7.87	جبل	119 8.35
رمادي	4 7.19	حزام	57 8.23
أبيض	42 6.99	عشب	51 8.1
خاف	5 6.78	مسيرة	83 7.76
أسود	12 6.36	مساح	67 7.65
مستقل	6 5.72	نبات	38 7.25
نظري	4 5.13	عنف	26 7.24
مفتوح	4 4.48	شجرة	48 7.23
		فلفل	22 7.12
		بصل	21 6.95
		جزيرة	57 6.94
		رقعة	22 6.87
		وادي	21 6.86
		قبة	23 6.83
		راية	25 6.8
		ورق	54 6.7

Page 1 of 9 Go Next | Last

الأراضي البنائية ذ 5 كانون الأول : حرائق هائلة تقضي على الأخضر واليابس في كل لبنان والحكمة الدولية تنفي تقريراً
العولمة ، الإمبراطور الأخير للحادثة المسيحية ، ستأتي على الأخضر واليابس في تربتنا وفي ترتيبهم ، عدا في حالة واحدة
والمجتمع والجماعات والأفراد في تصاعد وتيرته حتى أكل الأخضر واليابس ؟ ولا شك أن الظاهرة الاحتجاجية تندفعنا إلى
مجتمعية مسكنة اتقاء للهزات والعواصف التي تأتي على الأخضر واليابس ، كما حدث فعلا في بعض الأقطار العربية في
ان هذه القلعة هي المصدر الوحيد للرزق لكنهم قضاوا على الأخضر واليابس فقد بيعت هذه الشركة بمبلغ 1.3 مليار
وجاهليته وحارته إذا كان الهدف فتح حرب الطوائف لبحرق الأخضر واليابس . لكننا لا نعرف بالتحديد من يقف وراء العملية
" الوهم الطويل مندو ومرعب حتى لأعدائه ، هو الذي أدمن الأخضر واليابسة " في مرحلة صعوده ، يأتي إلا أن يطالهما وهو
القرى والمدن ، وجميعنا يعلم أن ظلم نظام بن علي أتى على الأخضر واليابس من بنزرت لبرج الخضراء ، ولكن أن تصبح شاشاتنا
سكان عدد كبير من القرى إذ قضت عليها تماما وأنت على الأخضر هل المصائب التي <p></p> . ؟ واليابس معا في ميامنار
والإتفاق تتراقب مع دعوات الي المواجعة والحرب وحرق الأخضر واليابس وترويح الناس ؟ وأي نظام أمني رسمي يمكن أن
حاق بفرص عمل ثمينة فارتموا في احضان مصارف قضت على الأخضر ويكثر من الاعتراض تقبل <p></p> . واليابس في جيوبهم
ومن يتبعه من المتعودين على الولايم والبنزس بدرس وطن الأخضر واليابس الذي لا يمت بصلة لشاريعهم بالاقتبال وتقسيم
اليوم فلا حيلة لي سوى المصمت فيواخرهم أصبحت تأتي على الأخضر واليابس وانا رغم ماقدمت لإزلت انا لا مرفاق ولا
سبعة أعوام عجاف أكل المستعمر وعسلاته <p></p> في عرافة الجريح ودمر جميع مؤسسات الدولة
حقيقي لها ، بل هي مسمى فقط ، وإن هي الا فرضي أنت على الأخضر <p></p> . واليابس أضاعت وحدة البلاد واستقرارها وحريتها
القوى المختلفة ، بل قد تشهد مصر حربا أهلية تأتي على الأخضر واليابس ، وفي ظل هذه المعطيات فإن من يدعو إلى إسقاط
وإن نواد هذه الفتن مع مهدفا قبل ان تنقش وتقتضى على الأخضر .. واليابس وتكون العواقب وخيمة ويكون قد فات الوقت
استفاقت مدينة ليون الفرنسية على حريق أتى على <p></p> الأخضر واليابس في مستودع كبير للحافلات مملوك لشركة كيبوليس
البدري ولعبه بطريقة لعب لاتناسب الفريق كادت ان تأتي على الأخضر واليابس في الفريق وتسببت بالفعل في خروج الفريق من
واحتدمت الأمور وانفجرت الحرب الأهلية التي أنت على الأخضر يغادر أعضاء الفريق البرتقالي مدينة <p></p> . واليابس

Page 1 of 9 Go Next | Last

Figure 9: Word Sketch results for أخضر (left)

Figure 10: Concordance lines for أخضر in combination with its collocate يابس (right)

In this concordance, we see that this combination usually occurs with 10) أتى على of the 20 lines in Figure 10) or verbs denoting destruction, such as قضى على (to destroy) for lines 1, 5, 11, and 17; and حرق (to burn) for line 10. Therefore, looking at the context, we can deduce the meaning “everything” for الأخضر واليابس and the idiom أتى على الأخضر واليابس (to destroy everything).

Additionally, in the Word Sketch, we see that a top collocate noun for the adjective أخضر is ضوء (light). Green light is not such a common phenomenon that it would account for this, so again, we look at the concordance (Figure 11).

سئون عسكري : الحكومة الإسرائيلية أعطت الضوء الأخضر كشف مسئول <p></p> writer أ <p></p> ل
الكوميكس الشهيرة جدا ، وهو فيلم الحركة البطولي الضوء الأخضر المنتظر عرضه خلال شهر يونيو القادم ، Lantern
ور قصة فيلم الضوء <p></p> . بروس الأميركيتين في مصر الأخضر حول هال جوردان ، وهو طيار في القوات الجوية الأمير
البريطاني ان " جمعية خيرية بريطانية حصلت على الضوء الأخضر لاطلاق لعبة بتصويب والجائزة عبارة عن علاجات للذ
وهناك أقول كثير مكتوبة في دينهم يشجعهم ويعطيهم الضوء الأخضر في التعامل مع الآخر بكل وحشية ودموية ، وهذا باطبع
اغتيال قيادات الحركة الشعبية ولكن بعد ما وجد الضوء الأخضر من الرئيس عمر البشير . في الاجتماع الذي تم بين الذك
أمريكية ، ومحمية بغيثو أمريكي يعطيها أتى شاءت الضوء الأخضر لتواصل جرائمها على مرأى من العالم وسماع ، فهي لا
ينتظر أن يوافق مسؤولو " العميد " على إعطائه الضوء الأخضر لمسح الديون من عائدات الفريق من حقوق البث التلفزيوني
تومي ، في الدقيقة الأولى من اليوم 5 جويلية ، الضوء الأخضر لانطلاق فعاليات المهرجان الثقافي الإفريقي ، بقصر ،
الإثنين الفارط وأشعره بقراره الأخير القاضي بإعطاء الضوء الأخضر للتشكيله البلدية للعودة إلى ملعب تشاكر أوراك المدرب
نفت مصادر سياسية رفيعة المستوى ان يكون هناك اي ضوء اخضر <p></p> . يتعلق برئاسة مجلس النواب ، وتفضيل مرشح
من لبنان وتوطين الفلسطينيين في لبنان وإعطاء الضوء الأخضر لإسرائيل بضم الضفة الغربية رسميا ، أو عمليا على الأ
جان بفلت من العقاب نتيجة لتخاذل السلطات بشكل ضوء أخضر بأن هذه السلطات لن تبالي بمحنة ضحايا العنف الجنسي
ج كيفن راد <p></p> [...] سوف تستضيف المؤتمر الأخضر من مراكز القوى لتحدي جوليا غيلارد على قيادة حزب
اللجنة العليا الصعود لتنفيذ التمرين بدون إشارة أو ضوء أخضر الدرجة النهائية = صفر مخالفات الفريق تنافس الجبازي
أصحابها عن أن السيد رئيس بلدية بنكبير قد أعطى الضوء الأخضر للمواطنين بالبناء والإصلاح دونما الحصول على رخص
سلاحه إلى أخيه سواء كان عنصرًا أو قياديا يعطي الضوء الأخضر للقتل (لان كلاهما قاتل) .. حتى يسفك الدم الذي تباكم
ن رب العمل قد أعطاني الضوء <p></p> السؤال <p></p> الأخضر لأخذ ما يكفيني من أرباح أمواله الذي هو خاصته ، لكن
بخفة ، بنادي لبيع بضاعته ، وحين تضيء الإشارة بال ضوء الأخضر بهم بالابتعاد خوفا من دهب بغير حساب ، وكثيرا ما ،
والتفاز والصحف . وقد أعطت إدارة بوش رايلي الضوء الأخضر لتشغيل إذاعة العراق الحر . رايلي يرتبط بخطة إدارة

Figure 11: Concordance lines for أخضر in combination with ضوء

In these lines, we can see that الضوء الأخضر (the green light) is used in much the same way as the English, in “to give/get the green light”, meaning to be allowed to go forward.

2.5 The Bilingual Word Sketch function

A new function of the Word Sketch is the bilingual word sketch, which allows the user to see word sketches for two words side-by-side. Figure 12 shows a comparison between أحمر and red.

Some of the same things are أحمر/red in Arabic and English; thus, we find the matched pairs لحم/meat, سجاد/carpet, and فلفل/pepper. All three are to an extent idiomatic, with the same idiomatic meaning in both languages. The Red Cross and Red Crescent are discussed more in Arabic media than in English, reflecting the unfortunate reality of several Arabic-speaking countries today. In contrast, wine is high in the English list but absent in the Arabic one.

adjective-of	9782	4.8	modifies	778081	0.4
صنّيب	1132	11.52	flag	27161	9.32
هلال	1101	11.38	carpet	21070	9.04
بحر	1498	11.0	wine	34956	8.8
خط	845	9.92	tape	15035	8.44
نون	507	9.75	meat	17594	8.34
قلعة	214	9.19	pepper	10687	8.26
بطاقة	249	9.08	herring	6067	7.93
نحم	208	8.87	light	25976	7.58
هندي	157	8.6	onion	5258	7.33
ساقية	93	8.24	rose	4714	7.3
دم	212	7.95	cell	16369	7.22
كربية	75	7.93	lipstick	3047	6.86
سجاد	77	7.92	ink	3605	6.73
شمع	67	7.76	bump	2858	6.59
شيطان	87	7.7	grape	2877	6.51
زاوية	69	7.54	stripe	2462	6.43
بافوت	47	7.23	ribbon	2586	6.41
فلفل	46	7.22	sole	2300	6.39
ورد	58	7.07	lip	3592	6.35
عظوب	34	6.78	shirt	4524	6.33
زاعوق	32	6.72	berry	2563	6.31
كبريت	32	6.71	hair	9863	6.29
خمبر	32	6.71	dress	6136	6.24
درب	34	6.64	snapper	1856	6.24
بقعة	33	6.62	arrow	2207	6.19

Figure 12: Adjective results of a bilingual word sketch for Arabic أحمر and English red

2.6 The Distributional Thesaurus function

The Sketch Engine also offers a distributional thesaurus, where, for the input word, the words 'sharing' the most collocates are presented. Figure 13 shows the top entries in similarity to تصدير (export). The top result is استيراد (import). Clicking on this word takes us to a 'sketch diff', which is a report that shows the similarities and differences between the two words in Figure 14.

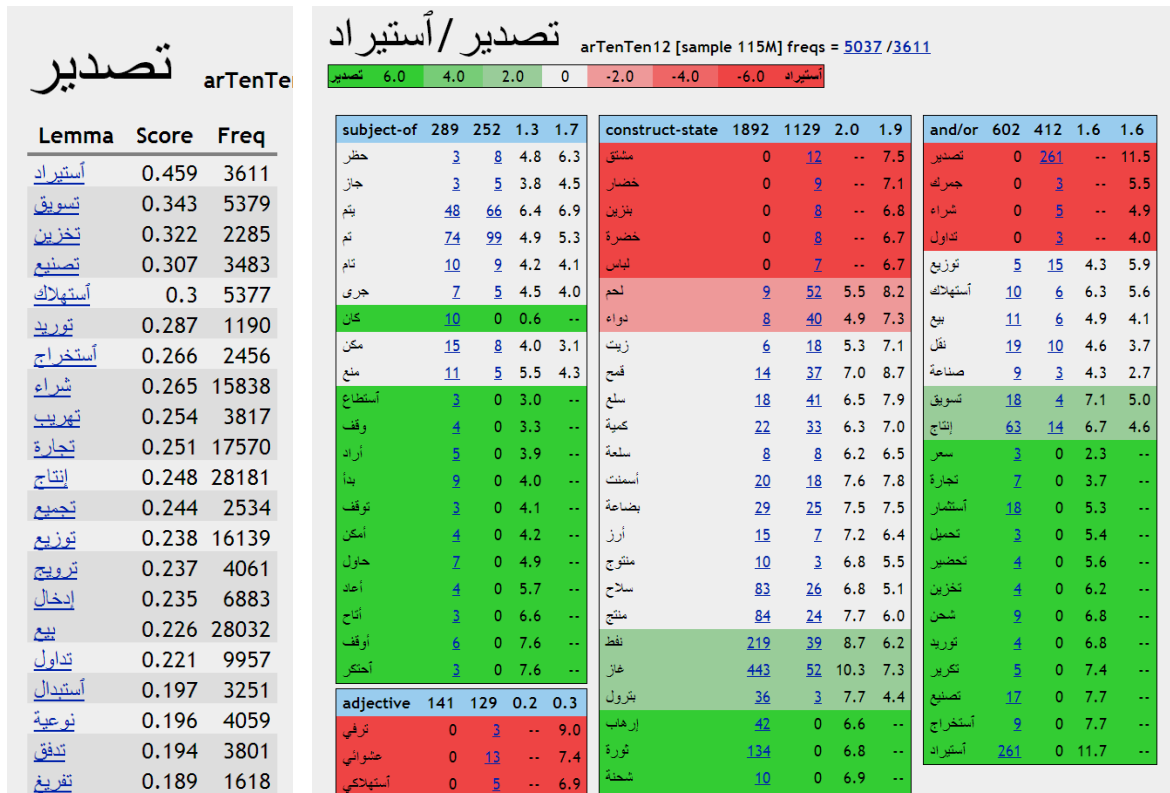


Figure 13: Thesaurus search showing entries similar to تصدير (export) (left)
Figure 14: Sketch Diff comparing collocates of تصدير and استيراد (export and import) (right)

The first number following the collocates shows the number of occurrences of this collocate with تصدير, the second number shows the number of occurrences with استيراد. A color scale from green to red visualizes the distribution.

2.7 Collocations and lexicographic research: Two case studies

The information in the Sketch Engine reports is particularly useful for lexicographers. It presents collocations, idioms, prepositions commonly occurring with verbs, and so forth.

It also gives insight into the use of words, often assisting the lexicographer in finding definitions for new words, for example, for توحدي (autistic), as shown in Figure 15. The immediate context of *child* and *patient* indicate that the word might be an adjective for an ailment.

وسيل , illuminative , الروح الى الله من خلال المسهل **توحدي** ويوجه الرسالة البابوية لآوون الثالث عشر ء
</p><p> : طفل " اعاني من التوحد وليست هي صفتي الوحيدة **توحدي** هو جزء بسيط من شخصيتي ولا يحددني كش
</p><p> صبراً .. صبراً .. صبراً </p><p> ... النهاية **توحدي** علي انه قدرة مختلفة وليست اعاقه , حاول ا
للتوحد يتوصل إلى أسبابه ... د . العياضي : 120 ألف طفل **توحدي** يوصل الفريق البحثي التابع للمركز </p><p> توحده . وقد نظم برنامج وسحور خيرى للرجال حيث قراء طفل **توحدي** آيات من الذكر الحكيم ثم عرض فيلم عن التو
على تخليص أبنائهم منه , مشيرة إلى تدهور صحة 400 مريض **توحدي** يتلقون العلاج عن طريق الأكسجين في أحد ا
أنا هنا الآن , في هذه النقطة الساحرة عادت إلي صبيحات **توحدي** مع آخر سأسميه (حياة) ... ذكراه تنهال ا!
حالتان : الحالة الأولى هي التوحد الكلي وانسلاخي أثناء **توحدي** عن هذا العالم ولحظة أخرى لا أستطيع الآن
لتوظيفه لخدمة توجه خاص لفئة معينة أو هدف أحادي أو فكر **توحدي** لأن هذا التحايل الرخيص ينم عن ضعف الحد
إذا ثبت أن الطفل - </p><p> . والمقاييس الخاصة بالتوحد **توحدي** فنوصي بإلحاقه بمركز خاص لذلك حتى يكو
حالته , ومنوهة بأهمية أن يكون لديه بطاقة تثبت أنه **توحدي** وطالبت السديري بتفعيل </p><p> . في .
بين ثلاثة أطفال , أما ناديه فإنها تعاني من مرض ذهاني **توحدي** أي لا تستطيع الكلام , و تعيش في عالمها الا

Figure 15: Concordance for توحدي

It also occasionally reveals new senses of words. For example, the word نسق is traditionally known to mean “order/manner”, as illustrated in Figure 16.

نَسَق *nasaq* order, array, layout, arrangement, disposition; connection, succession, sequence; manner, mode, system, method; symmetry; *nasaqan* in regular order, in rows | على نسق in the manner of; على نسق واحد in the same manner, equally, evenly, uniformly; see *nasaq* حروف النسق

النَّسَقُ - نَسَقٌ :
النَّسَقُ : ما كان على نظام واحد من كل شيء .
يقال : جاءَ القومُ نَسَقًا ، وَزَعَتِ الأشجارُ نَسَقًا .
ويقال : شَعَرَ نَسَقًا : مستويي البينة حسنَ التركيب ، ودرَّ نَسَقًا : منظمًا .
و النَّسَقُ المنسوق .
يقال : كلامٌ نَسَقٌ : متلائم على نظام واحد .
و (حروفُ النَّسَقِ) : حروفُ العطف .
المعجم: المعجم الوسيط

Figure 16: Dictionary entries for نسق from Wehr's *Dictionary of Modern Written Arabic* 4th ed. 1979, and *al-mu'jam al-wasit* (Academy of the Arabic Language in Cairo)⁸

However, looking at the concordance for the top adjective collocates تصاعبي (increasing, Figure 17), we see that these sentences do not seem to refer to “increasing order” but to an “increasing pace”

للتجار والقطاعات . فقد عرف حجم المشاركة من قبل التجار **نسقا** /نسق **تصاعديا** حيث تجاوز السنة الفارطة الألف تاجر واتسعت
قبل عشيرة وأقارب محمد البوعزيزي لتأخذ هذه المظاهرات **نسقا** /نسق **تصاعديا** مع انتحار شاب آخر بصعقة كهربائية في 22 ديسمبر
الشغالات (النحلة العاملة) يبدأ بالتناقص تدريجيا **وينسق** /نسق **تصاعدي** وسريع ينتهي بضعف ملفت للخلية ثم إتلافها نهائيا
الكبير في تنوعه . وكذلك الإستعمال المكثف والعشوائي **وينسق** /نسق **تصاعدي** للمبيدات الحشرية الفتاكة في النشاط الفلاحي
فرنسا وألمانيا وكندا بمكاتب تسجيل متنقلة مشيرا إلى **النسق** /نسق **التصاعدي** الملحوظ في عمليات تسجيل التونسيين بالخارج
كاكا " أحلى الأوقات " مع ريال مدريد هذه الفترة , بعد **النسق** /نسق **التصاعدي** الذي طرأ على أدائه منذ بداية الموسم الجاري

Figure 17: Concordance for نسق with تصاعبي

⁸ Entry as found at almaany.com, February 2014.

Investigating the word further, we find that “pace” is a common contemporary meaning of the word `نسق`.

Having shown the functions of the Sketch Engine and its functionality for Arabic, we will now go into more detail on developing the corpus and deploying it in the Sketch Engine.

3 A Sketch Grammar for Arabic

A sketch grammar is a grammar for the language based on regular expressions over part-of-speech tags (see Kilgarriff *et al.* 2004). It underlies the word sketches and is written in the corpus query language (CQL). A sketch grammar is designed particularly to identify head-and-dependent pairs of words (e. g., `نسق, تصاعبي`) in specified grammatical relations (here, adjective-modifier) so that the dependent can be entered into the head's word sketch and vice versa. Prior to the work described here, there has only been one sketch grammar for Arabic, developed at Oxford University Press (OUP) as part of the development phase for the Oxford Arabic Dictionary (Arts *et al.* 2014). It (and the word sketches resulting from it) is accessible only on arrangement with OUP.

The sketch grammar is one of the two components needed to build word sketches. The grammar is run over the corpus to identify all of the `<word1, grammatical-relation, word2>` triples in the corpus. The other component is a statistic. For each lemma occurring in the `word1` slot (the node word) and for each grammatical relation, we count the number of times each different lemma occurs in the `word2`, or 'collocate', slot. We use these numbers to calculate an association score⁹ between the node word and the collocate. The collocates with the highest association scores go into the word sketch.

A sketch grammar contains a set of definitions for grammatical relations. A simple grammatical relation definition is just:

```
=adjective  
1:"noun" 2:"adj"
```

This definition says that if we have a word with part-of-speech tag `noun` followed by one with part-of-speech tag `adj`, the grammatical relation `adjective` holds between the node word (the noun) and the collocate (the adjective). The 1: identifies the noun as the first argument of the grammatical relation, and the 2: identifies the adjective as the second argument.

We would also like to identify the noun as a collocate, when the adjective is the node word. To do that, we tell the system that the relation is `dual` and give a name for the inverse relation: here, `adjective-of`, as follows.

```
*DUAL  
=adjective/adjective-of  
1:"noun" 2:"adj"
```

There is some shorthand here. There may be many different fields of information associated with a word, of which the part-of-speech tag is just one field. In the case of arTenTen, there are many fields, including the word form itself, the lemma (with and without diacritics), the case and the state.¹⁰ The part-of-speech tag is called simply `tag` and in the formulation above, this has been set as the default. A non-shorthand version is

```
*DUAL  
=adjective/adjective-of  
1:[tag="noun"] 2:[tag="adj"]
```

All of the constraints on a word (or, technically, a *token*: tokens are usually either words or punctuation) are placed within square brackets, and each square-bracketed item relates to one token in a sequence.

9 The association score currently in use is a variant of the Dice coefficient; see Rychlý (2008) for full details.

10 See also section 4.2.

Now, the linguist will immediately note that there are many cases where adjectives happen to follow nouns but are not their modifiers. The definition above is insufficiently constrained and will give rise to many false positives. One constraint we want to add is that the adjective and noun agree, in case and in state. This is enforced in the next version.

**DUAL*

=adjective/adjective-of

1: [tag="noun"] **2:** [tag="adj"] & 1.state = 2.state & 1.case = 2.case

Now, an adjective followed by a noun only matches if the *state* value of the token indexed by 1: is the same as the *state* value of the token indexed by 2:, and likewise for *case*.¹¹

This is better and will not include many false positives. However, we should also be alert to valid cases of adjectives modifying nouns, which the definition above misses. One case is where two adjectives in succession modify a noun, e.g., السعودية المملكة العربية (lit: the Saudi-Arabian Kingdom). Only the adjective closest to the noun is captured by the clause above. To capture the other adjective, we add another clause to the definition:

1: [tag="noun"] [tag="adj"] **2:** [tag="adj" & pref1tag!="prep"] & 1.state = 2.state & 1.case = 2.case

This version allows an intervening adjective between the noun and its collocate adjective, which must not have a prefixed preposition.

The process of developing a sketch grammar is supported by the Sketch Engine because the CQL queries can be posed directly to the corpus, using the 'CQL' option in the concordance form. Thus, the strings above can be cut and pasted into the CQL box (Figure 18), and the developer can immediately see all of the hits (Figure 19).

Simple query:

[Query types](#) [Context](#) [Text types](#)

Query type simple lemma phrase word character CQL

Lemma:

Phrase:

Word Form: match case

Character:

CQL: Default attribute: [Tagset summary](#)

Figure 18: Using CQL in the concordance search form (with *tag* as default attribute)

11 Gender and number may seem to be good candidate features for this sketch grammar. However, since MADA uses what Habash (2010) terms *form-based* gender and number, and given the prevalence of deflected agreement (irrational plural nouns take feminine singular adjectives), these features are not good indicators of noun-adjective agreement. For more on issues of Arabic agreement, see Alkuhlani and Habash (2011).

Query noun, adj 25 > Random sample 25 (0.2 per million)		
Page 1	of 2	Go Next Last
307351	في إيران . وفي الليلة التي سبقتها ورغم	الأجواء القمعية المشددة
417451	وواشنطن ستتخذ كل الخطوات الضرورية لوقف	المشروع الذري الإيراني
544001	تسمع عن الفساد والفوضى والحرب في الصعده	والاختلالات السياسية والامنية
666351	وتطويرها تحت شعار (نحو منهج تربوي متطور لبناء	الانسان العراقي الجديد
862751	" , عبدالرحمن محمد </p><p>!! وهذا واحد منهم	المملكة العربية السعودية
1121401	رجال قالوا انهم تلقوا أموالا وأسلحة من	نائب لبياني مناهض
2150051	قسرا بما تقول فما الحال إذن في عقود منح	المصانع المحلية التابعة
2821501	هالفكرة عفا عليها الزمن علينا مراجعتها	مراجعة دقيقة وصرحة
3232751	لرصد العلاقة بين ارهاسات التفكير بجدارة	للإجتماع العربي الاسلامي
3559851	تحتوي على نسبة تفوق الحد المسموح به من قبل	المواصفات السعودية والأمريكية
4803851	فكم من عارض صحي يبدو تافها يكون مؤثرا ,	لأزمة صحية حقيقية
4816301	الموجبة إلى العراق في شهر رمضان على أنها	حرب صليبية جديدة
5737151	- قانون الاعمار وهو قانون 10 </p><p>. الانتخابات	برؤية علمية وعصرية
6352651	نبيه بري اعتقاده أن " الأمور ستتبلور خلال	الأيام القليلة المقبلة
6510601	فألذين في هذا العمل , ولا أي شكل من أشكال	البحث المنهجي الملائم
8194451	ليش كل الهجوم على معاليه مع انه للامانة	انسان متواضع ونظيف
8244451	وأوضح أن </p><p>. بما فيها الرغبات الشاذة	السبب المساند والمؤثر

Figure 19: Resulting concordance with noun-adj-adj sequences

Typically, this will include false positives, and the developer can then add constraints to rule them out. They should also think about the cases they are missing (in this example, the two-adjective case) and need to aim for as large a population of hits as possible, without too many false positives. In the terminology of information theory, they need to attend to recall - missing items that should be found - as well as precision - avoiding false positives. Recall tends to be a harder problem because a tool cannot show the items that are not found.

The Arabic sketch grammar aims at identifying the main grammatical relations while ensuring high-quality results. The grammatical patterns it covers are:

- **subject, subject-of:** These relations capture the relationship between verbs and their subjects. The noun is required to appear in the nominative case and may not have a prefixed preposition or conjunction.

The phrase نزل المطر (the rain fell) produces two grammatical relations. When نزل (fell) is the node word, the grammatical relation *subject* holds between it and its collocate المطر (rain). Conversely, if المطر is the node word, then it stands in the grammatical relation *subject-of* with نزل.

- **adjective, adjective-of:** These two relations capture noun-adjective pairs. We enforce agreement in state (definite/indefinite) and case. Enforcing agreement in gender and number is not trivial and left for future versions.

In the phrase بحث علمي (scientific research), the noun بحث takes the *adjective* علمي, which itself is *adjective-of* for بحث.

- **construct-state:** Captures construct state (idafa) constructions between two nouns. The first noun is required to be in the construct state and the second noun is required to be in the genitive case with no prefixed preposition or conjunction.

In the phrase مدير المدرسة (the school principal), the grammatical relation *construct-state* holds between the node word مدير (principal) and the collocate المدرسة (the school).

- **and/or**: This relation captures conjunctive constructions of pairs of nouns, adjectives, and verbs. We enforce agreement in certain grammatical features between the two words: for nouns and adjectives, we enforce agreement in case and state; for verbs, in aspect. This relation is declared as *symmetric*, which tells the system that both words can be the head node in turn.

Examples for pairs of adjectives include: كبير وصغير (large and small) and كبير أو صغير (large or small). In these examples, the word كبير (large) stands in grammatical relation of *and/or* with صغير (small) and vice versa. Similarly, we obtain pairs of nouns (e.g., النساء والرجال, "women and men") and verbs (e.g., يضحك أو يبكي, "laughs or cries").

The grammar focuses on the highest-confidence patterns for each grammatical relation. There are many constructions it does not yet cover. The quality of the identification of the different relations depends on the correctness of the automatic disambiguation component. Since the accuracy of automatic prediction of case is somewhere in the mid 80%, we can expect a fair amount of failed matches, e.g., verb-object pairs analyzed as verb-subject pairs. Future versions will increase coverage for current relations and add additional relations such as **verb-preposition** and **direct-object**. See Appendix 1 for the full grammar and the Sketch Engine documentation¹² for a full account of the formalism.

4 Creating and preparing the corpus

4.1 Crawling and text preparation

The following describes the processing chain for creating the corpus.

- We use texts from Arabic Wikipedia and other Arabic web pages to build the language-specific models that we need: (a) a character trigram model for language identification, (b) a byte trigram model for character encoding detection, (c) the most common Arabic words for seeding the crawl and for distinguishing sentences from lists and headers, and (d) parameters for the boilerplate cleaning utility.
- We crawl the Arabic web with SpiderLing¹³ (Pomikalek and Suchomel 2012), a crawler designed specifically for preparing linguistic corpora. The seeds for the crawl were generated by taking the top 1000 words from Arabic Wikipedia, randomly combining them into triples, and using the triples as Yahoo queries. The Yahoo search hits gave 4583 URLs, which were used as starting points for the crawl.
- We remove the non-textual material and boilerplate with jusText (Pomikalek 2011). JusText uses the working definition that we want only 'text in sentences' (excluding e.g., headers and footers). The algorithm is linguistically informed, rejecting material that does not have a high proportion of tokens that are the grammar words of the language; therefore, in the course of data cleaning, most material, which is not in the desired language, is removed.
- We de-duplicate with Onion (Pomikalek 2011) to remove near-duplicate paragraphs. We de-duplicate at the paragraph level because for many linguistic purposes, a sentence is too small a unit, but a whole web page (which may contain large chunks of quoted material) is too large.

These tools are designed for speed and are installed on a cluster of servers. For a language where there is plenty of material available, we can gather, clean and de-duplicate a billion words a day. ArTenTen was collected in 14 days. Table 1 presents the various statistics from arTenTen.

Data statistics	Documents (web pages; millions)	Sentences (millions)	Words (millions)	Data size
HTTP requests issued	87.8	–	–	–
Web pages received	58.8	–	–	2,015 GB
Cleaned text without exact duplicates	21.5	463	17,500	152 GB

¹² <http://www.sketchengine.co.uk/documentation>

¹³ <http://nlp.fi.muni.cz/trac/spiderling>

Final text without near duplicates	11.5	177	5,790	58.0 GB
Processed with MADA	0.23	4.5	115	1.32 GB ¹⁴

Table 1: Data sizes at the various stages of corpus preparation

4.2 Composition

The best-represented top level web domains in the corpus are .com, .net, .org, .info, .ps (Palestine), .sa (Saudi Arabia), .sy (Syria), .eg (Egypt), and .ae (United Arab Emirates), as shown in Table 2. There are 116,000 web domains represented by at least one document, and 43,000 represented by at least 10 (see Table 3), suggesting a heterogeneous corpus in contrast to corpora such as Arabic Gigaword or KSUCCA (Alrabiah *et al.* 2013), which are built from a small number of sources. The twenty domains that contributed the most documents are given in Table 4.

TLD	%	Note
.com	54.45	Generic commercial
.net	20.86	Generic network
.org	10.32	Generic organization
.info	1.69	Generic information
.ps	1.55	Palestine
.sa	1.41	Saudi Arabia
.sy	0.76	Syria
.eg	0.61	Egypt
.ae	0.60	United Arab Emirates
.cc	0.43	Cocos Islands/generic
.uk	0.41	UK
.cn	0.41	China
.jo	0.40	Jordan
.sd	0.38	Sudan
.ma	0.35	Morocco
.lb	0.30	Lebanon
.il	0.28	Israel
.biz	0.26	Generic business
.ws	0.26	Samoa/generic
.ir	0.25	Iran
Other	4.03	

Table 2: Document (web pages) by top-level domain (TLD)

>= 1 document	116,029 websites
>= 10 documents	43,282 websites
>= 100 documents	11,242 websites
>= 1,000 documents	2264 websites
>= 10,000 documents	112 websites

Table 3: Distribution of documents by website

aawsat.com	28,689
maghress.com	24,925
masress.com	23,818
sawt-alahrar.net	22,669
burnews.com	21,474
humum.net	21,084
chelseafarms.com	20,216

¹⁴ The size of the annotated corpus is 1.32 GB without morphological tags and 23.6 GB with full MADA morphological annotation.

nabanews.net	19,490
sarayanews.com	17,534
alghomhariah.net	17,090
anhri.net	16,718
tayyarcanda.org	16,315
arabic.xinhuanet.com	15,879
alsahafa.sd	15,774
m.islamweb.net	15,600
digital.ahram.org.eg	15,487
arabtimes.com	15,339
rosaonline.net	15,266
alwasatnews.com	15,210
elbiladonline.net	14,934

Table 4: Websites contributing the most documents

4.3 Processing with MADA

We chose to use the MADA tool for Arabic processing because of its state-of-the-art results on Arabic disambiguation, part-of-speech tagging and lemmatization and its holistic approach to modeling Arabic, predicting all of a word's morphological features in context. MADA has been successfully used by numerous Arabic NLP projects: in the NIST Open machine translation evaluation in 2012, nine out of twelve teams competing on Arabic-English translation used MADA. In a precursor to the work described in this article, Oxford University Press used MADA to prepare corpus materials used to create the Oxford Arabic Dictionary (Arts *et al.* 2014).

Within the framework of Arabic processing via MADA (Habash and Rambow 2005; Habash *et al.* 2009), we need to distinguish two concepts: **morphological analysis** and **morphological disambiguation**. **Morphological analysis** refers to the process that determines for a particular word all of its possible morphological analyses. The word, for MADA, is the orthographic word, defined as the sequence of letters delimited by spaces and punctuation. In Arabic, the word may include a variety of clitics, such as the definite article, prepositions, conjunctions and pronominals.

Each single analysis (out of many) includes a single choice or reading of the word with multiple dimensions of morphological information: the word's full diacritization, lemma, stem, part-of-speech (POS); the full Buckwalter Analyzer tag (Buckwalter 2002), values and POS tags for four possible proclitic slots; the values of eight inflection features -- person, aspect, voice, mood, gender, number, state and case; enclitic value and POS tag; English gloss; and whether the word had a spelling variation. Table 5 shows the MADA features for the example word *وبفكرة* *wbfrp* assuming a specific analysis corresponding to the English 'and with an idea'.

MADA Feature	Explanation of Feature
diac:wabifikorapK	Diacritization التشكيل
lex:fikorap_1	Lemma المفردة
stem:fikor	Stem الجذع
pos:noun	Part-of-speech قسم الكلام
BW:wa/CONJ+bi/PREP+fikor/NOUN+ap/NSUFF_FEM_SG+K/CASE_INDEF_GEN	Buckwalter POS tag قسم الكلام بنظام باكوالتر
prc3:null	Third proclitic position away from base word (typically, interrogative Hamza أداة \ سابقة لستفهام)
prc2:wa conj	Second proclitic position away from base word حرف \ سابقة عطف

prc1:bi_prep	First proclitic position away from base word حرف \ سابقة جر
prc0:0	Zeroth proclitic position away from base word (typically the determiner Al) ل \ سابقة التعريف
per:na	Person (not applicable here) الشخص
asp:na	Aspect (not applicable here) الزمن
vox:na	Voice (not applicable here) معلوم/مجهول \ البناء
mod:na	Mood (not applicable here) الصيغة
gen:f	Gender (feminine here) الجنس
num:s	Number (singular here) العدد
stt:i	State (indefinite here) التعريف
cas:g	Case (genitive here) الحالة لإعرابية
enc0:0	Only enclitic after the base word ضمير اللاحقة متصل
spvar:lex	Spelling Variant (none, exact lexicon match here) إملاء غير قبلي
gloss:idea;notion;concept	English gloss

Table 5: MADA analysis of *wbfrp* وبفكرة

Arabic words are highly ambiguous, primarily because diacritical marks are usually left out. A good analyzer produces the full set of choices for a particular word out of context. For example, the word *byn* can have many analyses, including:

Diacritization	Buckwalter POS tag	English Gloss
bay~an+a	PV+PVSUFF_SUBJ:3MS	He demonstrated
bay~an+~a	PV+PVSUFF_SUBJ:3FP	They demonstrated (f.p)
Biyn	NOUN_PROP	Ben
bay~in (dropping all case endings for simplicity)	ADJ	Clear
Bayn	PREP	Between, among

Morphological disambiguation refers to selecting the appropriate morphological analysis in context. Compare the following two sentences, which both contain *byn*. A good disambiguation model would select the proper noun reading for (1) and the preposition reading for (2):

(1) هل سينجح **بين** أفليك في دور باتمان?
Will **Ben** Affleck be a good Batman?

(2) كيري يحاول مجدداً اتقاذاً المفاوضات **بين** فلسطين وإسرائيل

Kerry tries again to save the negotiations **between** Palestine and Israel.

The task of morphological disambiguation for English is referred to as POS tagging because for English, a large part of the challenge is to determine what a noun, verb, or adjective is (for example, for base forms such as *promise*, s-forms such as *promises*, ing-forms such as *promising* and ed-forms such as *promised*). The standard English POS tag set, although only comprising 46 tags, completely disambiguates English morphologically. In Arabic, the corresponding tag set comprises thousands of tags, so the task is considerably harder. Reduced tag sets have been proposed for Arabic in which certain morphological differences are conflated, making the morphological disambiguation task easier. The term POS tagging is usually used for Arabic with respect to some of the smaller tag sets (Habash 2010).

MADA uses a morphological analyzer for MSA based on the standard Arabic morphological analyzer (SAMA) (Graff *et al.* 2009). It also uses a set of different classifiers that classify the values

of specific features from the analysis form in context, such as lemmas or gender. These features are trained on the Penn Arabic Treebank (Maamouri *et al.* 2004). The two sets of information (out-of-context analyses and in-context classified features) are combined to select the appropriate analysis in context (Habash and Rambow 2005; Roth *et al.* 2008).

A 115-million word subset of arTenTen was processed with MADA. The single preferred analysis for each word was output and used as the input to the next process. The work on MADA has been extended to handle Arabic dialects, specifically Egyptian Arabic (Habash *et al.* 2013). However, in this work, we only use MADA for MSA.

4.4 Into the Sketch Engine

Loading the arTenTen into the Sketch Engine required a conversion of MADA output into the format specified by the Sketch Engine. The Sketch Engine input format, often called “vertical” or “word-per-line”, is as defined at the University of Stuttgart in the 1990s and is widely used in the corpus linguistics community. Each token (e.g., word or punctuation mark) is on a separate line and where there are associated fields of information, such as lemma, POS-tag and morphological features, they are included in tab-separated fields. The conversion script extracts all of the MADA-generated features into fields and incorporates additional fields for ease of search in Sketch Engine, e.g., Arabic-script, diacritized and non-diacritized versions of the lemma (back-transliterated from the Buckwalter transliteration). Structural information, such as document beginnings and ends, sentence and paragraph mark-up, and any available metadata, are presented in XML-like form on separate lines. For web corpora, there is limited metadata available; date of collection and the URL from which the domain and top-level domain can be derived are useful. A sample of the vertical file is shown in Appendix 2.

In the Sketch Engine, each corpus has a corpus configuration file, which specifies the information fields that the corpus includes and various aspects on how they should be displayed. The next stage of the corpus preparation was to develop the arTenTen corpus configuration file. For instance, we needed to specify here that the word sketch attribute is the Arabic form of the lemma to facilitate searching by users in Arabic. This was problematic: it was not clear whether this should be the version of the lemma with diacritics or without. The no-diacritic option was desirable simply because it was the way that Arabic speakers usually write. If we did not permit no-diacritic input, beginner users would obtain no results and would be put off. However, if the diacritics are not written, the level of ambiguity is considerably higher, and it would not be possible to see a word sketch for صَافِر (to confiscate) without noise resulting from صَافِر (going out) because both are written as صافر when not diacritized. Thus, expert users would prefer that word sketches be computed on diacritized forms. The provisional solution is two versions of the corpus: one for users who know they need to use diacritized forms to obtain word sketches, the other for those who do not. We are currently building an interface option that allows users to use the undiacritized form while keeping the diacritized form as an option for advanced users.

We must note here that the quality of the output of the system depends heavily on the input, i.e., the quality of tagging and lemmatization. Errors in lemmatization and tagging will not go unnoticed and can lead to unexpected results for the lexicographer. There is generally a logical explanation, but it may require a closer view into the tagging and lemmatization to fully understand the output. One general difficulty is with proper nouns whose form is ambiguous with another word. For example, the name حَيِي (Huyay) is a common first name in religious texts. However, MADA usually tags it as an adjective meaning "modest", a mistake that stems from the fact that MADA is mostly built to process modern standard Arabic (MSA) texts, where this name is not a common one. It is also assigned the wrong lemma: حَيِي (Hayiy~) instead of حَيِي (Huyay~). Thus, when the lexicographer wants to search for words that may be read as proper nouns or adjectives, they must be aware of the ambiguity and either use the wrong lemma or search only with the simple string.

On the results page, the concordances are shown, by default, in a keyword-in-context (KWIC) view, as in Figure 2. With VIEW options, it is possible to change the concordance view to a number of alternative views. One is to view additional attributes such as POS tags or lemma alongside each word. This can be useful for finding out why an unexpected corpus line has matched a query, e.g., because of an incorrect POS-tag or lemma. By selecting fields in the references column, the user can decide what source of information should appear in blue at the left-hand end of the concordance line.

5 Summary and future plans

We have presented arTenTen, a very large web-crawled corpus of contemporary Arabic. We have also presented in some detail the subset of that corpus that has been processed by the MADA tool: how it has been set up and encoded and how we have produced word sketches for Arabic, with a full account of the sketch grammar that was used. We have discussed how this MADA-processed corpus can be used for dictionary-editing and related linguistic research, including how it can be used to find collocations, idioms, new words, new senses, and via the thesaurus, synonyms and related words. We have introduced the sketch diff, which shows how near-synonyms can be compared and contrasted.

We would of course like to apply MADA to the whole of arTenTen. To date, this has not been possible because of the speed of the program. This has recently been addressed with MADAMIRA (Pasha *et al.* 2014), a new and improved version of MADA combined with AMIRA (Diab 2009) that is orders of magnitude faster than MADA and has an output of comparable quality.

The method of compilation of arTenTen aims at a diverse corpus, including texts from many domains and genres. The nature of the Arabic language family also means that web texts are likely to appear in many language varieties: modern standard Arabic (MSA), classical Arabic, Quranic Arabic, and various dialects. Identifying the language variety of each text (or sub-text unit) is thus both a challenge and an opportunity: it is a non-trivial task, although standard language identification methods work quite well on identifying Arabic dialects (Zaidan and Callison-Burch 2013). The opportunity that lies in identifying the language varieties will facilitate lexicographic work on specific varieties and the comparative study of the dialects.

In preliminary experiments, we built a classifier to distinguish between MSA, classical Arabic, and Egyptian, Jordanian, and Saudi dialects. We trained a five-gram character level language model for each of these varieties based on published corpora and tested its performance on a small, manually selected subset of arTenTen texts in MSA, classical Arabic, and Egyptian Arabic, achieving 93% accuracy in this three-wise classification task. Then, we trained a combined dialectal model based on the Egyptian, Jordanian, and Saudi texts and processed a large number of arTenTen texts (40 k). We observed that the majority of the texts (~80%) are identified as MSA, and the rest are identified as classical or dialectal Arabic. This shows that a non-negligible portion of the texts is non-MSA. In future work, we intend to improve our language variety identification and increase its coverage to other dialects, using corpus-based approaches and resources, such as Buckwalter and Parkinson's Frequency Dictionary (2011) and the keywords method presented in Kilgarriff (2012). We will also consider the identification of sub-text units (Elfardy and Diab 2013), which is important for mixed texts.

arTenTen was gathered in 2012; so, it is already two years old. For each of the TenTen corpora, a program of re-crawling is planned, whereby material will regularly be added, both to keep the corpus current and so that empirical methods can be applied to the discovery of new words and meanings. We intend to gather newspaper feeds and blog feeds so that we have additional material with accurate time stamps.

We believe arTenTen, in combination with MADA/MADAMIRA and the Sketch Engine, possesses considerable promise for improved Arabic linguistic description and lexicography.

Acknowledgments

This work was partly supported by the Ministry of Education of the Czech Republic within the LINDAT-Clarin project LM2010013 and by the Ministry of the Interior of the Czech Republic within the project VF20102014003. Nizar Habash performed most of his work on this article while he was at the Center for Computational Learning Systems at Columbia University.

References

- Alansary, S., Nagi, M. and Adly, N. (2007) Building an International Corpus of Arabic (ICA): progress of compilation stage. In: *7th International Conference on Language Engineering*, Cairo, Egypt.
- Alkuhlani, S. and Habash, N. (2011) A Corpus for Modeling Morpho-Syntactic Agreement in Arabic: Gender, Number and Rationality. In: *Proceedings of the Association for Computational Linguistics (ACL'11)*, Portland, Oregon.
- Alrabiah, M., Al-Salman, A. and Atwell, E. (2013). The design and construction of the 50 million words KSUCCA King

- Saud University Corpus of Classical Arabic. In: *Second Workshop on Arabic Corpus Linguistics (WACL-2)*, Lancaster, UK 2013.
- Al-Sulaiti, L. and Atwell, E. (2006) The design of a corpus of Contemporary Arabic. *International Journal of Corpus Linguistics*, **11** (2).
- Arts, T. et al. (Forthcoming 2014) *Oxford Arabic Dictionary*. Oxford University Press.
- Buckwalter, T. *Buckwalter Arabic Morphological Analyzer v2.0*. LDC Catalog No.: LDC2004L02. Linguistic Data Consortium.
- Buckwalter, T. and Parkinson, D. (2011) *A Frequency Dictionary of Arabic*. Routledge. Frequency Dictionary Series.
- Diab, M. (2009) Second generation AMIRA tools for Arabic processing: Fast and robust tokenization, POS tagging, and base phrase chunking. In: *2nd International Conference on Arabic Language Resources and Tools*, Cairo, Egypt.
- Eckart, T., Quasthoff, U., Alshargi, F. and Goldhahn, D. (2014) Large Arabic Web Corpora of High Quality: The Dimensions Time and Origin. In: *Proceedings of the Workshop on Free/Open-Source Arabic Corpora and Corpora Processing Tools, LREC 2014*. Reykjavik, Iceland.
- Elfardy, H. and Diab, M. (2013) Sentence Level Dialect Identification in Arabic. In: *Proceedings of the Association for Computational Linguistics (ACL'13)*, Sofia, Bulgaria.
- Graff, D. (2003) Arabic Gigaword. LDC Catalog No.: LDC2003T12. Linguistic Data Consortium.
- Graff, D., Maamouri, M., Bouziri, B., Krouna, S., Kulick, S. and Buckwalter, T. (2009) Standard Arabic Morphological Analyzer (SAMA) Version 3.1. LDC Catalog No.: LDC2009E73. Linguistic Data Consortium.
- Habash, N. (2010) *Introduction to Arabic Natural Language Processing*. Morgan & Claypool Publishers. Synthesis Lectures on Human Language Technologies.
- Habash, N. and Rambow, O. (2005) Arabic Tokenization, Part-of-Speech Tagging and Morphological Disambiguation in One Fell Swoop. In: *Proceedings of the Association for Computational Linguistics (ACL'05)*, Ann Arbor, Michigan.
- Habash, N., Rambow, O. and Roth, R. (2009) MADA+TOKAN: A Toolkit for Arabic Tokenization, Diacritization, Morphological Disambiguation, POS Tagging, Stemming and Lemmatization. In: *Proceedings of the International Conference on Arabic Language Resources and Tools*, Cairo, Egypt.
- Habash, N., Roth, R., Rambow, O., Eskander, R. and Tomeh, N. (2013) Morphological Analysis and Disambiguation for Dialectal Arabic. In: *Proceedings of Conference of the North American Association for Computational Linguistics (NAACL)*, Atlanta, Georgia.
- Habash, N., Soudi, A. and Buckwalter, T. (2007) On Arabic Transliteration. In: A. van den Bosch and A. Soudi (eds.). *Arabic Computational Morphology: Knowledge-based and Empirical Methods*. Springer.
- Jakubíček, M., Kilgarriř, A., Kovář, V., Rychlý, P. and Suchomel, V. (2013) The TenTen Corpus Family. In: *International Conference on Corpus Linguistics*, Lancaster, UK.
- Kilgarriř, A. (2012) Getting to Know your Corpus. In: *Proceedings of Text, Speech, Dialogue Conference*, Brno, Czech Republic 2012. Springer.
- Kilgarriř, A., Rychly, P., Smrz, P. and Tugwell, D. (2004) The Sketch Engine. In: *Proceedings of EURALEX*, Lorient, France 2004. pp 105-116.
- Maamouri, M., Bies, A., Buckwalter, T. and Mekki, W. (2004) The Penn Arabic Treebank : Building a Large-Scale Annotated Arabic Corpus. In: *NEMLAR Conference on Arabic Language Resources and Tools*, Cairo, Egypt 2004. pp 102-109.
- Pasha, A., Al-Badrashiny, M., El Kholy, A., Eskander, M., Diab, M., Habash, N. Pooleery, M., Rambow, O. and Roth, R. (2014). MADAMIRA: A Fast, Comprehensive Tool for Morphological Analysis and Disambiguation of Arabic. In: *Proceedings of the 9th International Conference on Language Resources and Evaluation*, Reykjavik, Iceland.
- Pomikalek, J. (2011) *Removing Boilerplate and Duplicate Content from Web Corpora*. PhD, Masaryk University.
- Pomikalek, J. and Suchomel, V. (2012) Efficient web crawling for large text corpora. In: *Proceedings of the 7th Web as Corpus Workshop (WAC7)*, Lyon, France.
- Roth, R., Rambow, O., Habash, N., Diab, M. and Rudin, C. (2008) Arabic morphological tagging, diacritization, and lemmatization using lexeme models and feature ranking. In: *ACL 2008: The Conference of the Association for Computational Linguistics Companion Volume, Short Papers*, Columbus, Ohio.
- Rychlý, P. (2008). A lexicographer-friendly association score. In: *Proceedings of Recent Advances in Slavonic Natural Language Processing, RASLAN 2008*, Karlova Studanka, Czech Republic.
- Sharoff, S. (2006) Creating general-purpose corpora using automated search engine queries. In: M. Baroni & S. Bernardini (eds.). *WaCky! Working papers on the Web as Corpus*. Bologna, Italy.
- Wehr, H. (1979) *Dictionary of Modern Written Arabic* 4th ed., Spoken Language Services.
- Zaidan, O. and Callison-Burch, C. (2013) Arabic Dialect Identification. *Computational Linguistics*.
- Zaghouani, W. (2014) Critical Survey of the Freely Available Arabic Corpora. In: *Proceedings of the Workshop on Free/Open-Source Arabic Corpora and Corpora Processing Tools, LREC 2014*. Reykjavik, Iceland.

Appendix 1: Arabic Sketch Grammar

arTenTen Sketch Grammar, version 0.1 (7/20/2013)

*STRUCTLIMIT s

*DEFAULTATTR tag

*FIXORDER subject/subject-of adjective/adjective-of construct-state and/or

*DUAL

=subject/subject-of

1:"verb" 2:[tag="noun" & case="n" & pref1tag!="prep" & pref2tag!="conj"]

*DUAL

=adjective/adjective-of

1:"noun" 2:[tag="adj" & pref1tag!="prep" & pref2tag!="conj"] & 1.state = 2.state & 1.case = 2.case

1:"noun" [tag="adj" & pref1tag!="prep" & pref2tag!="conj"] 2:[tag="adj" & pref1tag!="prep"] & 1.state = 2.state & 1.case = 2.case

noun adjective pair; enforce agreement in state and case

=construct-state

1:[tag="noun" & state="c"] 2:[tag="noun" & case="g" & pref1tag!="prep" & pref2tag!="conj"]

simple annexation

#1:[tag="noun" & state="c"] [tag="noun" & case="g" & state="c" & pref1tag!="prep" & pref2tag!="conj"]+
[tag="noun" & case="g" & pref1tag!="prep" & pref2tag!="conj"]

more complex annexation

=and/or

*SYMMETRIC

1:"noun" [trans=">w"|trans=">m"|trans="w"] 2:"noun" & 1.state = 2.state & 1.case = 2.case

1:"noun" 2:[tag="noun" & pref2="wa"] & 1.state = 2.state & 1.case = 2.case

noun

1:"adj" [trans=">w"|trans=">m"|trans="w"] 2:"adj" & 1.state = 2.state & 1.case = 2.case

1:"adj" 2:[tag="adj" & pref2="wa"] & 1.state = 2.state & 1.case = 2.case

adjective

1:"verb" [trans=">w"|trans=">m"|trans="w"] 2:"verb" & 1.aspect = 2.aspect

1:"verb" 2:[tag="verb" & pref2="wa"] & 1.aspect = 2.aspect

verb

Appendix 2: Sample arTenTen XML 'vertical' format

With selected attributes of a morphological annotation by MADA. There are two paragraphs (<p>) each with one sentence (<s>) within one document (<doc>). The source of the document and other metadata is stored in attributes of structures (e.g. url="http://www.alsabar-mag.com/ar/article__419").

word	word latin	diac	lemma voc latin	lemm a voc latin	lemma latin	lemma stem	tag	bw	person	aspect	voice	mood	gender	number	state	case	gloss	lex/punc	
<doc id="301" length="6615" url="http://www.alsabar-mag.com/ar/article__419">																			
<p>																			
<s id="8135">																			
كلمات	klmAt	kalimAti	kalimap_1	كَلِمَة	klmp	كلمة	kalim	noun					f	p	c	a	words;remarks	lex	
للبحث	llbHv	lilbaHovi	baHov_1	بَحْث	bHv	بحث	baHov	noun					m	s	d	g	discussion	lex	
</s>																			
</p>																			
<p>																			
<s id="8136">																			
الناصرة	AlnASrp	Aln~ASirapi	nASir_2	نَاصِر	nASr	ناصر	nASir	adj					f	s	d	g	partisan;supporter	lex	
:	:	:	:_0	:	:	:		punc									:	punc	
انطباعات	AnTbAEAt	AinoTibAEAtN	{inoTibAE_1	أَنْطَبَاع	{nTbAE	أَنْطَبَاع	{inoTibAE	noun					f	p	i	n	impression	lex	
من	mn	min	min_1	مِنْ	Mn	من	min	prep									from	lex	
البرلمان	AlbrlmAn	AlbarolamAni	barolamAn_1	بَرْلَمَان	brlmAn	برلمان	barolamAn	noun					m	s	d	g	parliament	lex	
التي	Al*y	Al~a*iyy	Al~a*iyy_1	الَّتِي	Al*y	التي	Al~a*iyy	pron_rel					m	s	i	u	which;who;whom_ [masc.sg.]	lex	
عقد	Eqd	Euqida	Eaqad-i_1	عَقَدَ	Eqd	عقد	Euqid	verb			3		p	p	i	m	s	be_held;be_con- vened;be_con- cluded	lex
في	fy	fiy	fiy_1	فِي	Fy	في	fiy	prep									in	lex	
حديقة	Hdyqp	HadiyqapK	Hadiyqap_1	حَدِيقَة	Hdyqp	حديقة	Hadiyq	noun					f	s	i	g	garden	lex	
عامّة	EAm~p	EAm~apK	EAm~_1	عَامَ	EAm	عام	EAm~	adj					f	s	i	g	general;common;p ublic	lex	
</s>																			
</p>																			
More paragraphs follow..																			
</doc>																			