

جامعة نيويورك أبو ظبي



WIDH@NYCDH 2021



Introduction to Arabic Text Processing with Python and CAMEL Tools

Starting soon!

جامعة نيويورك أبو ظبي



WIDH@NYCDH 2021



Introduction to Arabic Text Processing with Python and CAMEL Tools

Salam Khalifa and Ossama Obeid

Computational Approaches to Modeling Language (CAMEL) Lab
New York University Abu Dhabi

Roadmap

خارطة الطريق

- **Arabic NLP Challenges**

- **CAMeL Lab Solutions:**

CAMeL Tools

- **Questions**

- **تحديات اللغة العربية**

- **حلول من مختبر كامل:**

أدوات كامل

- **أسئلة**

The Arabic Language

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

اللغة العربية

- فصحي التراث
 - لغة القرآن الكريم
 - نصوص تاريخية
- فصحي العصر
 - اللغة الرسمية في البلدان العربية
 - لغة الصحف والإعلام
 - لها إملاء قياسي وقواعد قياسية
 - لغة الأمة
- اللهجات العربية
 - لغة التحدث الرئيسية
 - بدون إملاء أو قواعد قياسية رسمية
 - لغة الأم
 - وجود متزايد على وسائل التواصل الاجتماعية

The Main Challenges for Arabic Language Processing

- Orthographic ambiguity
- Morphological richness
- Dialectal variation
- Orthographic inconsistency
- Resource poverty (data & tools)
- Limited research

التحديات الرئيسية للمعالجة الآلية للغة العربية

- الإبهام الإملائي
- الغنى الصرفي
- تعدد اللهجات
- الأخطاء الإملائية
- فقر موارد البيانات والأدوات
- البحث العلمي المحدود

Arabic Script

الخط العربي

العربية ←

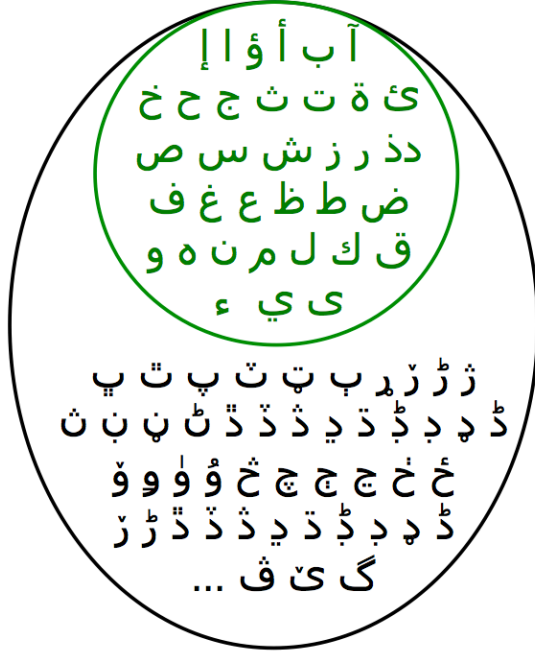
العربية ←

- Written right-to-left
- Letters have contextual variants
- Used to write many languages besides Arabic: Persian, Kurdish, Urdu, Pashto, etc.

- مكتوب من اليمين إلى اليسار
- حروف أشكالها متغيرة سياقيا
- يستخدم لكتابة العديد من اللغات إلى جانب العربية: الفارسية والكردية والأردية والباشتو، إلخ.

Arabic Script

الخط العربي



العربية ←

العربية ←

- Written right-to-left
- Letters have contextual variants
- Used to write many languages besides Arabic: Persian, Kurdish, Urdu, Pashto, etc.

- مكتوب من اليمين إلى اليسار
- حروف أشكالها متغيرة سياقيا
- يستخدم لكتابة العديد من اللغات إلى جانب العربية: الفارسية والكردية والأردية والباشتو، إلخ.

Orthographic Ambiguity

الإبهام الإملائي

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

- الإملاء بالخط العربي يستخدم التشكيل الاختياري
 - ١،٥٪ من كلمات نصوص الصحف تحتوي على علامات تشكيل
 - الكلمات العربية الفصحى ٦،٨ تشكيلات و٢،٧ مداخل معجمية بالمتوسط

وبعقدنا

وَبِعُقْدِنَا وَبِعِقْدِنَا وَبِعَقْدِنَا وَوَبِعَقْدُنَا

and with our (contract, necklace, psychoses) | and he stresses us out

Morphological Richness

الغنى الصرفي

- Arabic has a very rich inflectional system
 - For example, Arabic verbs have 5,400 inflected forms
 - Whereas English verbs have 6 and Chinese verbs have 1!

- اللغة العربية نظام صرفي غني جداً
ينتج عنه تراكيب كثيرة
 - فمثلاً، للفعل العربي حوالي ٥,٤٠٠ تصريف
 - بينما للفعل الإنجليزي ٦ تصاريف،
وللفعل الصيني تصريف واحد!

وسنقولها

/wasanaqūluhā/

و + س + ن + قول + ها

wa+sa+na+qūl+u+hā

and+will+we+say+it

And we will say it

قال، قالت، قالا، قالوا، قلت،

قلت، قلتما، قلتم، قلتن،

يقول، يقول، يقل، تقول، تقول،

تقل، تقولين، تقولي،

...فقال، فقالت، فقالا...

...وسأقولها، وسنقولها ... ،

Dialectal Variation

تعدد اللهجات

• Coarse Classification

- Gulf Arabic
- Levantine Arabic
- Egyptian Arabic
- Maghrebi Arabic

• Classification Problems

- Countries with multiple varieties
- Sub-dialectal varieties:
e.g., Urban, Rural, and Bedouin

• Diglossia and Multilinguality

• التصنيف العام

- اللهجة الخليجية
- اللهجة الشامية
- اللهجة المصرية
- اللهجة المغاربية

• مشاكل تصنيفية

- تعدد اللهجات في نفس البلد
- تعدد اللهجات الفرعية:
المدنية، والفلاحية والبدوية، مثلا

• الازدواجية والتعددية اللغوية

Phonological Variation

الاختلافات الصوتية

- First noticed differences
 - Not the most important

- الاختلافات الصوتية تُلاحظ فورياً
 - ليست الأكثر أهمية

MSA الفصحى		Dialects اللهجات	
ق	/q/	ق، ك، ء، گ، دج	/q/, /k/, /ʔ/, /g/, /dʒ/
ث	/θ/	ث، ت، س	/θ/, /t/, /s/
ذ	/ð/	ذ، د، ز	/ð/, /d/, /z/
ج	/dʒ/	ج، گ، دج	/dʒ/, /g/, /ʒ/

قلب → /qalb/, /galb/, /ʔalb/ (heart)

/fagr/ → فجر، فقر (poverty, dawn)

Morphological Variation

- Compared to Standard Arabic, dialectal morphology can be simpler or more complex
 - No case
 - New complex clitics

الاختلافات الصرفية

- مقارنة بصرف الفصحى، فإن صرف اللهجات أحياناً أبسط وأحياناً أكثر تعقيداً
 - اختفاء الحركات الإعرابية
 - تعدد التراكيب اللواصقية

كتابُ → كتابٌ، كتابٌ، كتاباً، كتابٍ، كتابٍ

(book)

وماكتبوها لوش → ولم تكتبوها له

(you all did not write it for him)

Morphological Variation

الاختلافات الصرفية

- Dialects vary systematically among themselves

- هناك العديد من الاختلافات المنهجية الصرفية بين اللهجات

	Past الماضي	Future المستقبل
فصحى Standard	كتب /kataba/	سيكتب /sajaktubu/
شامي Levantine	كتب /katab/	حيكتب /ħajiktob/
مصري Egyptian	كتب /katab/	هيكتب /hajiktib/
مغربي Moroccan	كتب /kteb/	غيكاتب /ʕajekteb/

Lexical Variation



*You say **to-MAY-to**,
I say **to-MAH-to**!*

الاختلافات المعجمية

براد؟



/b a r r a a d/



/b a r r a a d/



/b r a a d/

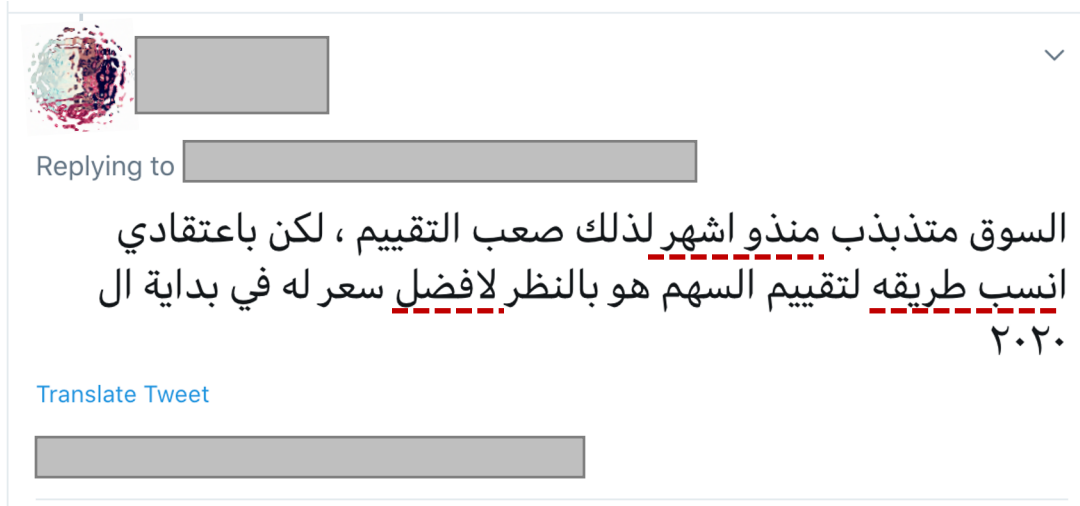
أنا أشهر براد بيت!
I'm the most
famous "brad"

Common Spelling Errors and Inconsistency

كثرة الأخطاء الإملائية والتضارب الإملائي

- Unedited standard Arabic text on news commentaries has an error rate of ~30% (QALB project)

- نسبة الخطأ في التعليقات الغير محررة المكتوبة بالعربية الفصحى على صفحات بعض مواقع الأخبار تصل إلى ٣٠٪ (مشروع «قلب»)



Common Spelling Errors and Inconsistency

- But in the dialects, there is no official orthography! So, there are many free forms including Arabizi.

كثرة الأخطاء الإملائية والتضارب الإملائي

- ولكن في اللهجات، لا يوجد معيار إملائي قياسي! لذا، هناك تنوع في التهجئات بما في ذلك الإملاء المرومن (العريزي).

وجدنا على الإنترنت ٢٧ تهجئة لهذه الكلمة باللهجة المصرية

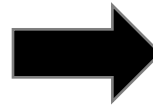
We found 27 spellings online of the Egyptian Arabic word

ما يقولهاش

/mabiʔulhāʃ/

“he does not say it”

(Habash et al., 2018)



Arabic Orthography	Arabic Transliteration	Frequency
مبيقولهاش	<i>mbyqwlhAš</i>	≈ 26,000
ما بيقولهاش	<i>mA byqwlhAš</i>	≈ 13,000
ما بقلهاش، مبقولهاش، مبقلهاش، ما بقلهاش، مابقولهاش	<i>mAbqlhAš, mbqwlhAš, mbqlhAš, mA bqlhAš, mAbbyqwlhAš</i>	≤ 10,000
ما بقلهاش، مابقولهاش، مبيقلهاش، ما بيقلهاش مابقلهاش	<i>mAbqwlhAš, mA bqwlhAš, mbyqlhAš, mA byqlhAš mAbbyqlhAš</i>	≤ 1,000
مبئلهاش، ما بيئولهاش، ما بيئولهاش، ما بيئولهاش	<i>mbÿlhAš, mAbÿwlhAš, mA byÿwlhAš, mAbÿwlhAš</i>	≤ 100
ما بيئلهاش، ما بيئلهاش، مبيئلهاش، ما بيئلهاش، ما بئلهاش، ما بئلهاش، مبيئلهاش، ما بئلهاش، مبيئلهاش، ما بئلهاش	<i>mA byÿlhAš, mAbÿlhAš, mbyÿwlhAš, mA byÿlhAš, mAbÿwlhAš, mA bÿlhAš, mA bÿlhAš, mbÿwlhAš, mbyÿwlhAš, mAbÿwlhAš, mbÿwlhAš</i>	≤ 10

بضعة تهجئات عريزية: Arabizi variants: mabi2ulhash, mabiquulhash, mabi'ulhash,...

The Main Challenges for Arabic Language Processing

التحديات الرئيسية للمعالجة الآلية للسغة العربية

- Orthographic ambiguity
- Morphological richness
- Dialectal variation
- Orthographic inconsistency
- **Resource poverty (data & tools)**
- Limited research

- الإبهام الإملائي
- الغنى الصرفي
- تعدد اللهجات
- الأخطاء الإملائية
- فقر موارد البيانات والأدوات
- البحث العلمي المحدود

Roadmap

- Arabic NLP Challenges
- CAMEL Lab Solutions:
CAMEL Tools
- Questions

خارطة الطريق

- تحديات اللغة العربية
- حلول من مختبر كامل:
أدوات كامل
- أسئلة

Computational Approaches to Modeling Language Lab

Researchers

Research

Teaching

Publications

Resources

News

Work With Us

Home / Research / Centers, Labs and Projects / Computational Approaches to Modeling Language Lab

Computational Approaches to Modeling Language Lab

The Computational Approaches to Modeling Language (CAMEL) Lab is a research lab at New York University Abu Dhabi established in September 2014. CAMEL's mission is research and education in artificial intelligence, specifically focusing on natural language processing, computational linguistics, and data science. The main lab research areas are Arabic natural language processing, machine translation, text analytics and dialogue systems.



مختبر كامل
CAMEL Lab

Featured Projects

TOIA

TOIA (time-offset interaction application) is a bilingual (Arabic-English) conversational agent, similar to a chat bot, except that the avatar is based on pre-recorded videos of an actual human being.



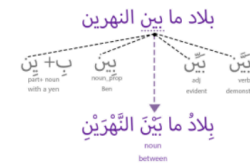
MADAR

Parallel Dialectal Corpus and Lexicon for 25+ cities



Gumar

A browsable 100-million word data set of Gulf Arabic



MADAMIRA

Arabic Artificial Intelligence

www.camel-lab.com

موارد كامل
CAMEL Resources

أدوات كامل
CAMEL Tools



موارد كامل

CAMeL Resources

Corpora

[\[+ \] Expand all](#)

+ ANERcorp - CAMeL Lab Train/Test Splits
+ Arab-Acquis: A parallel corpus of Arabic with all the languages of the European Parliament
+ The Arabic Parallel Gender Corpus
+ Gumar: The Gulf Arabic Corpus
+ Margarita Dialogue Corpus
+ Multi-Arabic Dialect Corpus (MADAR Corpus)
+ NYUAD Universal Dependency of Arabic (NUDAR) Treebank
+ Palm Treebank: A Multi-genre Arabic Dependency Treebank
+ Qatar Arabic Language Bank (QALB Corpus)

Lexicons

[\[+ \] Expand all](#)

+ Arabic Multidialectal Word Embeddings
+ ArabScribe
+ Multi-Arabic Dialect Lexicon (MADAR Lexicon)
+ SAMER Readability Lexicon

Tools

[\[+ \] Expand all](#)

+ ADIDA: Arabic Dialect Identification
+ BOTTA: An Arabic Dialect Chatbot
+ CALIMA-STAR
+ CAMeLParser
+ CAMeL Tools
+ DALILA: The Dialectal Arabic Linguistic Learning Assistant
+ Gulf Arabic Morphological Analyzer
+ MADAMIRA: Morphological Analyzer and Disambiguator for Arabic
+ MADARi: Web-base Morphological Annotation Tool
+ Palmyra
+ SIMMR
+ YAMAMA: Arabic Morphological Analyzer and Disambiguator

Guidelines

[\[+ \] Expand all](#)

+ CAMeL Arabic Phonetic Inventory (CAPHI)
+ CAMeL Part-of-speech tagset and guidelines (CAMEL POS)
+ CODA: Conventional Orthography for Dialectal Arabic

أدوات كامل CAMEL Tools

The Camel Tools Team

فريقي أدوات كامل



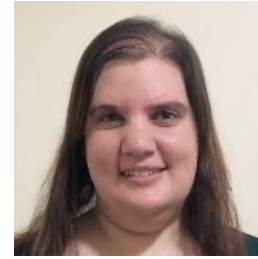
أسامة عبيد
Ossama Obeid



ناصر زلموط
Nasser Zalmout



سلام خليفة
Salam Khalifa



ديما التاجي
Dima Taji



مي عودة
Mai Oudah



بشار الحفني
Bashar Alhafni



جو إينوي
Go Inoue



فضل الإرياني
Fadhl Eryani



اليكساندر إردمان
Alex Erdmann



نزار حبش
Nizar Habash



أدوات كامل

CAMeL Tools

- An Open-source Python Toolkit for Arabic Natural Language Processing
 - Arabic Specific
 - Flexibility in use and reuse
 - Modularity
 - High Performance
 - Ease of use for beginner
- مجموعة أدوات مفتوحة المصدر بلغة بايثون لمعالجة اللغة العربية
 - التركيز على اللغة العربية
 - مرونة الاستخدام وإعادة الاستخدام
 - وحدات مستقلة سهلة التجميع
 - الأداء عالي الجودة
 - سهولة الاستخدام للمبتدئين

المقارنة مع رزم أخرى Comparison with other toolkits

	CAMEL Tools	MADAMIRA	Stanford CoreNLP	Farasa	
Target Language	Arabic خاص بالعربية	Arabic خاص بالعربية	Multilingual متعدد اللغات	Arabic خاص بالعربية	نوع الحزمة
Programming Language	Python بايثون	Java جافا	Java/Python جافا/بايثون	Java جافا	لغة البرمجة
CLI	✓	✓	✓	✓	واجهة سطر الأوامر
API	✓	✓	✓	✓	واجهة برمجة التطبيقات
Exposed Pre-processing	✓				المعالجة المسبقة للنصوص
Morphological Modeling	✓	✓			نمذجة التصريف
Disambiguation	✓	✓	✓	✓	حل الالتباس
POS Tagging	✓	✓	✓	✓	توسيم قسم الكلام
Diacritization	✓	✓		✓	التشكيل
Tokenization	✓	✓	✓	✓	التقطيع
Lemmatization	✓	✓		✓	المدخل المعجمي
Named Entity Recognition	✓	✓	✓	✓	التعرف على الكيانات المسماة
Sentiment Analysis	✓				تحليل المشاعر
Dialect ID	✓				تحديد اللهجة

- Other new toolkits include Adwat (Zerrouki, 2020) and ASAD (Mubarak et al., 2020)

- من بين الرزم الأخرى الحديثة، رزمة «أدوات» (زروقي، ٢٠٢٠) ورزمة «أسد» (مبارك وآخرون، ٢٠٢٠)

Roadmap

خارطة الطريق

- Arabic NLP Challenges
- CAMEL Lab Solutions:
CAMEL Tools
- Questions

- تحديات اللغة العربية
- حلول من مختبر كامل:
- أدوات كامل
- أسئلة

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

CAMeL-Lab / camel_tools

Unwatch 15

★ Unstar 113

Fork 27

<> Code

🔔 Issues 2

🔗 Pull requests

▶ Actions

📁 Projects

📖 Wiki

🛡 Security

📈 Insights

🔗 master ▾


🔗 8 branches

🔗 11 tags

Go to file

Add file ▾

📄 Code ▾

 owo Merge pull request #35 from slkh/docs ... 524fd52 on Dec 16, 2020 🕒 315 commits

📁 .github/ISSUE_TEMPLATE	Delete usage-and-support.md	5 months ago
📁 camel_tools	Bumped version to 1.1.0.	2 months ago
📁 docs	Added an example to dediac.	2 months ago
📁 tests	Updated copyright date.	5 months ago
📄 .gitignore	Finished CALIMA Star CLI plus extra documentation a...	2 years ago
📄 .readthedocs.yml	Fixed issues with Read the Docs builds.	6 months ago
📄 CONTRIBUTING.rst	Updated python version support mentions.	5 months ago
📄 LICENSE	Updated copyright date.	5 months ago
📄 MANIFEST.in	Remove Calima Star database from repo.	5 months ago
📄 README.rst	Updated installation instructions.	2 months ago
📄 camel_tools_logo.png	Updated CAMeL Tools Logo.	6 months ago
📄 setup.py	Implemented camel_data CLI utility.	2 months ago
📄 tox.ini	Updated python version support mentions.	5 months ago

About

A suite of Arabic natural language processing tools developed by the CAMeL Lab at New York University Abu Dhabi.

nlp sentiment-analysis
named-entity-recognition
nlp-apis arabic nlp-library
morphological-analysis
dialect-identification

📖 README

📄 MIT License

Releases 11

📄 camel-tools v1.1.0 Latest
on Nov 30, 2020

+ 10 releases

<https://camel-tools.readthedocs.io/>

camel_tools
latest

Search docs

CONTENTS:


- Overview
- Getting Started
- Command-line Tools
- Python API Reference
- Reference
- License

Private repos and priority support. Try Read the Docs for Business Today!

Sponsored · Ads served ethically

Docs » CAMEL Tools Documentation [Edit on GitHub](#)

CAMEL Tools Documentation



أدوات كامل CAMEL Tools

Contents:

- [Overview](#)
 - [About](#)
 - [License](#)
- [Getting Started](#)
 - [Installation](#)
 - [Datasets](#)
 - [Next Steps](#)
- [Command-line Tools](#)
 - [camel_transliterate](#)
 - [camel_arclean](#)
 - [camel_word_tokenize](#)
 - [camel_dediac](#)
 - [camel_diac](#)
 - [camel_morphology](#)
- [Python API Reference](#)

Read the Docs v: latest

<https://bit.ly/3a40q30>



CAMEL_Tools_Guided_Tour.ipynb ☆

File Edit View Insert Runtime Tools Help Last saved at 7:59 PM



Comment



Share



Table of contents



Introduction

Arabic NLP Refresher



Installation and Setup



Preprocessing

Simple Transliteration

Unicode Normalization

Orthographic Normalization

Diacritization

Word Tokenization

Morphology

Analysis

Generation

Reinflection

Disambiguation

Tagging

Tokenization

Dialect Identification

Sentiment Analysis

Named Entity Recognition

Section

+ Code + Text

Connect

Editing



Introduction

This notebook provides a quick overview of the functionalities provided by CAMEL Tools. While it doesn't cover every component in detail, it does provide a starting point for learning CAMEL Tools and navigating its APIs. For more detailed information please refer to the [CAMEL Tools documentation](#).

Arabic NLP Refresher

Before diving into CAMEL Tools, we recommend a quick refresher on Arabic NLP and the challenges it presents. [This webinar](#) presented by Dr. Nizar Habash provides a short introduction on Arabic NLP as well as short examples of how CAMEL Tools can be used to solve various problems. A similar presentation of this webinar in Arabic is also available [here](#).

Note: The examples in this talk may not be compatible with current or newer versions of CAMEL Tools.

Installation and Setup

The following steps are needed if you want to run the examples in this notebook on Google Colaboratory. If you want to run this notebook on your own machine, please follow the [installation instructions](#) instead.

First, we install the CAMEL Tools Python package.

```
[ ] %pip install camel-tools
```




شكراً على المتابعة! Thank You

www.camel-lab.com
{oobeid,salamkhalifa}@nyu.edu

diac	وَسَيَكْتُبُونَهَا	التشكيل
lex	كُتِبَ-u_1	المدخل المعجمي
root	ك.ت.ب	الجذر
pattern	وَسَيُ123ُونَهَا	الوزن
caphi	w_a_s_a_y_a_k_t_u_b_uu_n_a_h_aa	اللفظ
gloss	and+_will+_they_(people)+write+it;them;her	الترجمة
pos	verb	قسم الكلام
catib6	PRT+PRT+VRB+NOM	قسم الكلام/كاتب
ud	CONJ+AUX+VERB+PRON	قسم الكلام/يودي
prc2	wa_conj	سابقة ٢
prc1	sa_fut	سابقة ١
per	3	الشخص
gen	m	الجنس
num	p	العدد
asp	i	الزمن
vox	a	البناء
mod	i	الصيغة
enc0	3fs_dobj	لاحقة ١
d1tok	وَسَيَكْتُبُونَهَا	تقطيع د١
d2tok	وَسَيَكْتُبُونَهَا	تقطيع د٢
atbtok	وَسَيَكْتُبُونَهَا	تقطيع بنك بن الشجري
d3tok	وَسَيَكْتُبُونَهَا	تقطيع د٣
bwtok	وَسَيَكْتُبُونَهَا	تقطيع باكوالتر
bw	و/CONJ+س/FUT_PART+ي/IV3MP+كُتِبَ/IV+ون/و/DO:3FS/IVSUFF_ها	تحليل باكوالتر
pos_logprob	-1.023208	احتمالية قسم الكلام
lex_logprob	-3.648503	احتمالية المدخل المعجمي