

“ALKHALILDWS”: AN ARABIC DICTIONARY WRITING SYSTEM RICH IN LEXICAL RESOURCES

MOHAMMED REQQASS¹, ABDELHAK LAKHOUAJA¹ AND MOHAMAD BEBAH²

16 – 17 October 2019
ICALP 2019, Nancy France

Mohammed Reqqass¹, Abdelhak Lakhouaja¹ and Mohamad Bebah²

¹ Faculty of Sciences, Mohamed First University, Av Med VI BP 717, Oujda 60000, Morocco
reqqass.mohammed@gmail.com, abdel.lakh@gmail.com

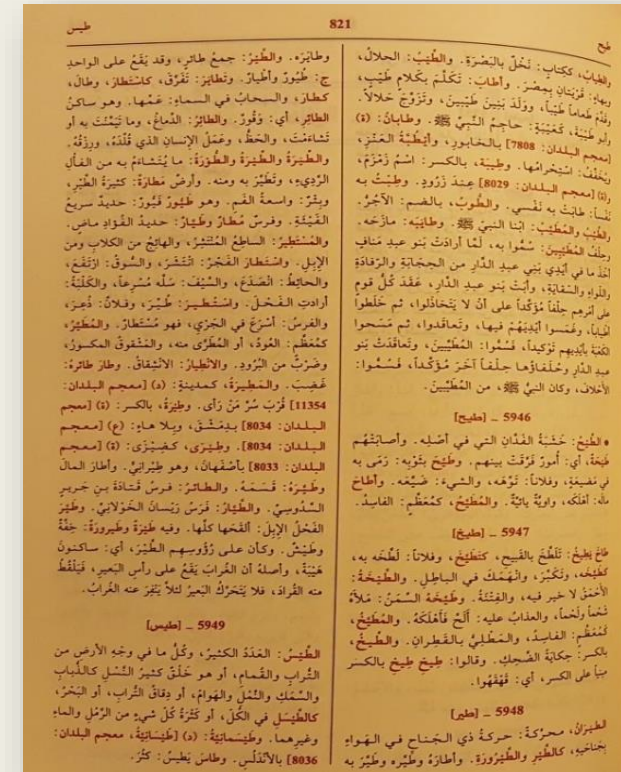
² Arab Center for Research and Policy Studies, Doha, Qatar
mohamad.bebah@dohainstitute.org

OUTLINE

- Introduction
- Dictionary Writing System
- “AlkhalilDWS”: An Arabic DWS rich in lexical resources
- Conclusion

Introduction

- A dictionary is a support that list the words of one or more specific languages, often listed alphabetically
- The process of making dictionary:
 - Identification of lexical entries
 - Editing lexical entries
 - Publication of lexical entries
- “AlkhalilDWS” assist the editor to make their Arabic dictionary project



Dictionary Writing System

Software for writing and producing a dictionary. It meet the following needs:

- write and edit dictionary entries
- control of the quality entries and ensuring their consistency
- exploitation of the available linguistic resources
- production of the dictionary in several forms: paper, electronic...
- exchange with other applications
- make collaborative group work

Dictionary Writing System

Dictionary Writing System	Description
TshwaneLex	a commercial writing system of monolingual dictionaries, bilingual dictionaries and multilingual dictionaries
Matapuna	an open source dictionary writing system, used to build a dictionary of Maori language
Glossword	is a system for writing and publishing multi-language dictionaries. It is a web application, open source.
ABBYY Lingvo Content	a dictionary writing system intended for compiling dictionaries, glossaries, encyclopedias, and other types of reference materials.

“AlkhalilDWS”: An Arabic Dictionary Writing System

- Nature of DWS
 - Based on a server– client architecture ;
- Data Base of DWS
 - Allows multi-user to work on the same dictionary-making project from any location ;
 - The data must be stored in a structured and extensible database;

“AlkhalilDWS”: An Arabic Dictionary Writing System

- Morphological analyzer
 - is analysis of structure of the word forms, it provides information about: word category (nouns, verb, adj, etc.), root, lemma, gender, etc.
 - find new lexical entries,
 - add grammatical information to lexical entries,
 - improve the search for a lexical entry in the dictionary.

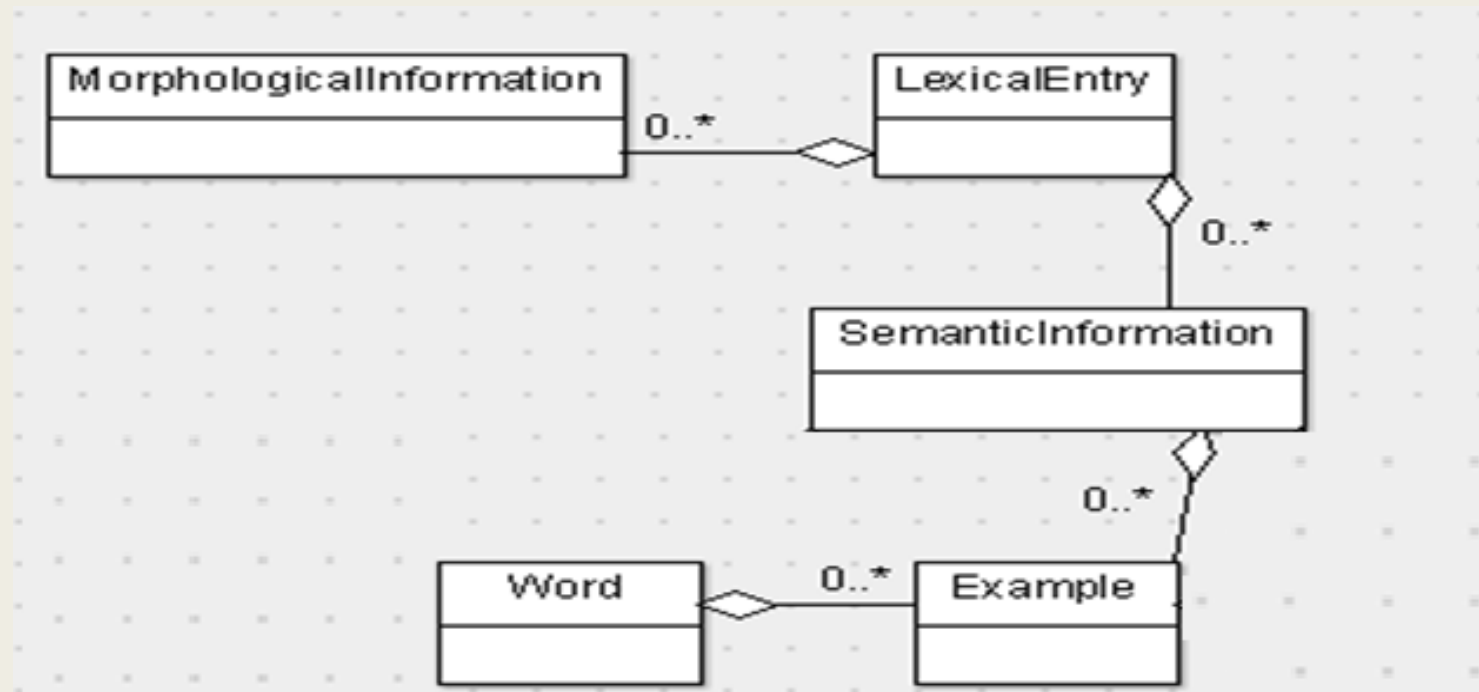
“AlkhalilDWS”: An Arabic Dictionary Writing System

- The spell checker
 - a software feature that checks for misspellings in a text;
 - errors categories : reading errors, hearing errors ,touch-typing errors, morphological errors, editing errors

“AlkhalilDWS”: An Arabic Dictionary Writing System

■ Lexical Markup Framework

- model that provides a common standardized framework for the construction of natural language processing lexicons



“AlkhalilDWS”: An Arabic Dictionary Writing System

- Dictionary Structure
- Lexical resources in “AlkhalilDWS”
- Operating “The Dictionary of the Modern Arabic Language”
- Profiles in “AlkhalilDWS”
- Identification of lexical entries
- Edition of lexical entries
- Publishing of lexical entries

“AlkhalilDWS”: An Arabic Dictionary Writing System

■ Dictionary Structure

- **macrostructure**: editor defines the global order of the lexical entries in the dictionary;
- **microstructure**: editor defines the internal structure of the lexical entry and the information it contains;

“AlkhalilDWS”: An Arabic Dictionary Writing System

- Lexical resources in “AlkhalilDWS”
 - “The interactive dictionary of the Arabic language”: an interactive open source web application based on “Alwassyt” dictionary.
 - “Almustakshif dictionary”: the first Arabic lexical encyclopedia of the dictionary “AlMoheet”
 - “The Dictionary of the Modern Arabic Language”: dictionary targeting to cover classical Arabic words still in usage, as well the new words used in different Arab countries.

“AlkhalilDWS”: An Arabic Dictionary Writing System

- Lexical resources in “AlkhalilDWS”
 - “The interactive dictionary of the Arabic language”: an interactive open source web application based on “Alwassyt” dictionary.
 - “Almustakshif dictionary”: the first Arabic lexical encyclopedia of the dictionary “AlMoheet”
 - **“The Dictionary of the Modern Arabic Language”**: dictionary targeting to cover classical Arabic words still in usage, as well the new words used in different Arab countries.

“AlkhalilDWS”: An Arabic Dictionary Writing System

- Operating “The Dictionary of the Modern Arabic Language”:
 - convert the document format of the dictionary to a text file,
 - keep only the lexical entries (remove introduction, index);
 - identify and build a mapping root, headword of lexical entries and the content of lexical entries;
 - The root is extracted from the line that start with a number followed by the character “-” and isolated Arabic characters. We keep the isolated Arabic characters only.
 - The headword is extracted from the line that does not start with a number or that start with the character “•”. We keep the word that appears before the character “:” and the character “]”.

“AlkhalilDWS”: An Arabic Dictionary Writing System

- Operating “The Dictionary of the Modern Arabic Language”:
 - The content of lexical entries is extracted from the line that does not contain root and from the line that contains a headword. If the line contains a headword, we keep only the words appearing before the first occurrence of the character ":";
 - extract meaning, semantic field and example from the content of lexical entry;
 - the semantic is the word that appears between the first occurrence of the character ")" and the first occurrence of the character "(";
 - the examples appear between the character “””, each example is separated by the character "-" or the character “° ”;
 - the meaning is the set of words that do not represent the semantics or the examples.
 - normalize the writing of the root by removing the space with the isolate Arabic character and replace the characters “ġ”, “ġ”, “ġ” by the character “ﺀ”;
 - normalize the writing of the lemma by analyzing them with the morphological analyzer Alkhalil2 and keeping only the result that corresponds to the lexical entries.

“AlkhalilDWS”: An Arabic Dictionary Writing System

- Operating “The Dictionary of the Modern Arabic Language”:

149 - أسد د

أَسَد [مفرد]: ج أساد وأُسُد وأُسُود، مؤ أَسَدَة:

1 - (حن) حيوان مفترس شديد الضراوة من فصيلة السِّنَّورِيَّات ورُتْبَة اللَّوَّاحِم، يشمل الذكر والأنثى ويطلق على الأنثى أَسَدَة ولَبُؤَة، وله في العَرَبِيَّة أسماء كثيرة أشهرها الليث والضيغم والغضنفر والضرغام "رأيت أسدًا- هذا السَّبَل من ذاك الأسد: يشبه الابن أباه في صفاته- *أَسَد عليّ وفي الحروب نعامَةٌ*" ° أسد الله: حمزة بن عبد المطلب رضي الله عنه- بين فكِّي الأسد: في خطر، في مأزق- حَصَّة الأسد: الجزء الأكبر.

“AlkhalilDWS”: An Arabic Dictionary Writing System

■ Operating “The Dictionary of the Modern Arabic Language”:

Element	Translation	Content
Root	A s d	أ س د
normalized root	a s d	ءسد
Lemma	Lion (assad)	أَسَد
normalized lemma	lion	أَسَد
Meaning	Very predatory animal. wild member of the felidae family. It includes males and females. The female is called lioness. it has many names in Arabic, including “leith”, “Ghadanfar”, “Dirgham” etc..	حيوان مفترس شديد الضراوة من فصيلة السِّتَوْرِيَّاتِ ورُثْبَةُ اللّٰوَحِمِ، يشمل الذكر والأنثى ويطلق على الأنثى أسدة ولبؤة، وله في العربية أسماء كثيرة أشهرها الليث والضيغم والغضنفر والضرغام
Semantic	Animal	حيوان
Examples	<ul style="list-style-type: none"> • I saw a lion • This cub is from that lion • A lion against me and soft in the war (proverb) • Lion of Allah: Hamza ibn Abdul-Muttalib (person) • Between lion’s jaws: in danger, in trouble • The lion's share: the major share of something 	<ul style="list-style-type: none"> • رأيت أسدًا • هذا السَّبَل من ذلك الأسد: يشبه الابن أباه في صفاته • *أَسَد عَلِيّ وفي الحروب نعامة* • أسد الله: حمزة بن عبد المطلب رضي الله عنه • بين فكيّ الأسد: في خطر، في مأزق • حصّة الأسد: الجزء الأكبر

“AlkhalilDWS”: An Arabic Dictionary Writing System

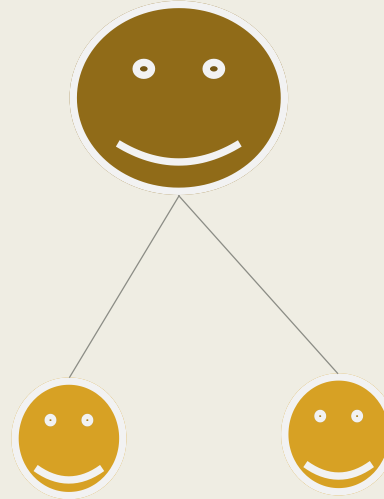
- Lexical resources in “AlkhalilDWS”

Dictionary	root	lemma	meaning	example	semantic field
the interactive dictionary of the Arabic language	7081	133591	133594	11643	247
the modern Arabic dictionary	5739	32258	91874	30909	35
almustakshif dictionary	5386	20112	26561	0	37

“AlkhalilDWS”: An Arabic Dictionary Writing System

■ Profiles in “AlkhalilDWS”:

chief editor: It allows the user to create his own working group, to identify lexical entries, to assign lexical entries to the members of his group, to approve the works submitted by editors, to export the dictionary in XML.



administrator: manage the accounts of chief editors.

editor: It allows the user to edit entries affected to him.

“AlkhalilDWS”: An Arabic Dictionary Writing System

- Identification of lexical entries
 - **manual entering**: the editor enters the words one by one. This is done by filling in the following information: lemma, type, and root.
 - **importation of lemmas**: the editor imports a list of words from a text file. Each line in the file contains the morphological information of word: lemma; type; root.
 - **automatic extraction of lemmas**: the editor can download a text file. “AlkhalilDWS” analyzes the words of the text and extract the new lemmas. Every lemma is linked with its morphological information.

“AlkhalilDWS”: An Arabic Dictionary Writing System

■ Identification of lexical entries

The screenshot displays the AlkhalilDWS web interface with a dark blue header and a light blue sidebar. The main content area is divided into three vertical panels:

- Left Panel: بناء مداخل معجمية من نص حر**
سيتتم تحليل النص المدخل بالحلل الصرفي الخليل 2، وبعد ذلك يتم استخراج المداخل المعجمية الجديدة.
تحدد الملف
No file chosen
- Middle Panel: تحميل ملف بالمداخل المرشحة**
(بالصيغة: الفرع:الوسم:الجذر)
تحدد الملف
No file chosen
- Right Panel: إضافة مدخل معجمي**
الفرع

الوسم

الجذر

The sidebar on the right contains the following menu items:

- الرئيسية
- المستخدمون
- التحليل الصرفي
- البحث في المعجم
- تسجيل الخروج
- موارد المعجم
- إضافة كتاب
- لائحة الكتب
- المداخل المعجمية
- مداخل معجمية للتحرير
- مداخل قيد التحرير
- مداخل محررة
- مداخل معتمدة
- ربط المعاني بالمداخل
- المعجم
- استخراج نسخة إكس إم إل
- استخراج نسخة ورقية

“AlkhalilDWS”: An Arabic Dictionary Writing System

■ Identification of lexical entries

تسجيل الخروج		المستخدمون		التحليل الصرفي		البحث في المعجم		الرئيسية	
« 4 3 2 1 »									
تجاهل	إحالة	اخيار المحرر	الرسم	الجذر	الفرع	1	أَعْلَمَ	علم	فعل
تجاهل	إحالة	اخيار المحرر	الرسم	الجذر	الفرع	2	أُعْلِمَ	علم	فعل
تجاهل	إحالة	اخيار المحرر	الرسم	الجذر	الفرع	3	إِعْلَامٌ	علم	اسم
تجاهل	إحالة	اخيار المحرر	الرسم	الجذر	الفرع	4	تُعَلِّمُ	علم	اسم
تجاهل	إحالة	اخيار المحرر	الرسم	الجذر	الفرع	5	عَالِمٌ	علم	اسم
تجاهل	إحالة	اخيار المحرر	الرسم	الجذر	الفرع	6	عَالِمٌ	علم	فعل

موارد المعجم

إضافة كتاب

لائحة الكتب

المدخل المعجمية

مدخل معجمية للتحريف

مدخل قيد التحريف

مدخل محررة

مدخل معتمدة

ربط المعاني بالمدخل

المعجم

استخراج نسخة إكس إم إل

استخراج نسخة ورقية

“AlkhalilDWS”: An Arabic Dictionary Writing System

■ Edition of lexical entries

The screenshot displays the AlkhalilDWS interface for editing a lexical entry. The top navigation bar includes 'الرئيسية', 'البحث في المعجم', 'التحليل الصرفي', 'المستخدمون', and 'تسجيل الخروج'. The main content area is titled 'المدخل المعجمي (الفرع: أدق، الرسم: فعل، المصدر: أدق)'. On the left, a sidebar shows 'معاني مقترحة' with a search bar and buttons for 'التكشاف' and 'نقل معنى'. The central area shows two entries for 'أدق' with their respective meanings and a 'حذف هذا المعنى' button. The right sidebar contains 'موارد المعجم' (including 'إضافة كتاب' and 'لائحة الكتب'), 'المدخل المعجمية' (including 'مداخل معجمية للتحرير' and 'مداخل قيد التحرير'), and 'المعجم' (including 'استخراج نسخة إكس إم إل' and 'استخراج نسخة ورقية').

“AlkhalilDWS”: An Arabic Dictionary Writing System

■ Publishing of lexical entries

the interactive version: "AlkhalilDWS" offers an interactive version of the dictionary. It allows the chief editor to search in his produced dictionary with the following options:

- **search with exact matching**: search the exact word in the lexical entries;
- **search with exact matching and ignore the diacritics**: search the lexical entries that correspond to the exact word ignoring the diacritics;
- **search after analyzing word**: search the lexical entries that correspond to the lemma of the search entry.

The screenshot displays the AlkhalilDWS web interface. At the top, there is a navigation bar with the following items from left to right: "تسجيل الخروج" (Logout), "المستخدمين" (Users), "التحليل الصرفي" (Morphological Analysis), "البحث في المعجم" (Search in the Dictionary), and "الرئيسية" (Home). The main content area is divided into two columns. The left column is titled "عرض المعاجم" (View Dictionaries) and contains two links: "تصفح الخفول الدلالية في معجم المستكشف" (Browse semantic anomalies in the Al-Khalil dictionary) and "تصفح معجم اللغة العربية المعاصرة" (Browse the Modern Arabic Dictionary). The right column is titled "البحث في المعجم" (Search in the Dictionary) and features a search input field labeled "الكلمة" (Word). Below the input field are several radio buttons for search options: "معجمي" (Lexical), "المستكشف" (Al-Khalil), "معجم اللغة العربية المعاصرة" (Modern Arabic Dictionary), "البحث مع تحليل الكلمة" (Search with word analysis), "البحث بدون تشكيل" (Search without morphology), and "البحث المطابق" (Exact search). A blue button labeled "عرض النتائج" (View Results) is positioned at the bottom of the search section.

“AlkhalilDWS”: An Arabic Dictionary Writing System

■ Publishing of lexical entries

the interactive version: "AlkhalilDWS" offers an interactive version of the dictionary. It allows the chief editor to search in his produced dictionary with the following options:

- **search with exact matching**: search the exact word in the lexical entries;
- **search with exact matching and ignore the diacritics**: search the lexical entries that correspond to the exact word ignoring the diacritics;
- **search after analyzing word**: search the lexical entries that correspond to the lemma of the search entry.

The screenshot displays the AlkhalilDWS dictionary interface. At the top, it says 'تصفح معجم اللغة العربية المعاصرة' (Browse the Modern Arabic Dictionary). Below that, there's a search bar with 'مرض الخلل الدلالي' (Semantic Defect Disease) entered. To the right, there's a dropdown menu showing 'الخلل الدلالي: الكيباء والصدبة' (Semantic Defect: Kibaa and Sadba). Below the search bar, there's a pagination control showing '« 4 3 2 1 » 1'. The main content area shows search results for 'أزوت' (Azot). The first result is 'أزوتية (أ ب ن و س، اسم)' (Azotiyya (A B N W S, Name)). The second result is 'مادة صلبة سوداء، غير موصلة للكهرباء، تنتج من خلط الكبريت بالمطاط القوي. (الخلل الدلالي: الكيباء والصدبة)' (A dark, solid, non-conductive material, produced from mixing sulfur with strong rubber. (Semantic Defect: Kibaa and Sadba)). The third result is 'أزوت (أ ز و ت، اسم)' (Azot (A Z W T, Name)). The fourth result is 'غاز شفاف لا لون له ولا رائحة ولا طعم يُعبر من أهم العناصر الطبيعية الحياتية وهو أكثر غازات الهواء مقداراً، يدخل في تركيب المواد البروتينية والأنسجة الحية الحيوانية والنباتية (النظر: أز و ت - أزوت). (الخلل الدلالي: الكيباء والصدبة)' (A colorless, odorless, and tasteless gas, one of the most important natural elements of life, and the most abundant gas in the atmosphere, it enters into the composition of protein materials and the living tissues of animals and plants (see: Az W T - Azot). (Semantic Defect: Kibaa and Sadba)). The fifth result is 'أزوتي (أ ز و ت، اسم)' (Azoti (A Z W T, Name)). The sixth result is 'حمض أزوتيّ سائل لا لون له، يُطلق أخرجه خائفة في الهواء، يُؤكسد الهيدروجين والكبريت وعصير الفحم وسوى ذلك، ويخرج بالصدوا والبوتاس فيتكوّن السمدان المعروفان بترات الصودا وترات البوتاس، وهو موصل جيد للكهرباء، ويُستعمل لصنع المفرقات ولإذابة الذهب والبلاتين (النظر: أز و ت - أزوتي). (الخلل الدلالي: الكيباء والصدبة)' (Azotiic acid, a colorless liquid, releases a sharp odor in the air, oxidizes hydrogen and sulfur and other things, and is extracted as soda ash and potassium ash, it is a good conductor of electricity, and is used for making insulators and for dissolving gold and platinum (see: Az W T - Azoti). (Semantic Defect: Kibaa and Sadba)). The seventh result is 'أزوتير (أ ن ر، اسم)' (Azotir (A N R, Name)). The eighth result is 'الأزوتير سائل عسوي طيار، لا لون له، يُستخدم مذيباً في الصناعة، ومُخدرًا في الطب. (الخلل الدلالي: الكيباء والصدبة)' (Azotir, a colorless, volatile, oily liquid, used as a solvent in industry, and as an anesthetic in medicine. (Semantic Defect: Kibaa and Sadba)).

“AlkhalilDWS”: An Arabic Dictionary Writing System

■ Publishing of lexical entries

XML format: The chief editor can export his produced dictionary in XML format. This format respects the standard LMF. Bellow an example of some entries edited with “AlkhalilDWS”

```
<?xml version="1.0" encoding="UTF-8" standalone="no"?>
<Dictionary>
<entry id="1887">
<feat att="partOfSpeech" val="اسم"/>
<root>
<feat att="writtenForm" val="آخر"/>
</root>
<lemme>
<feat att="writtenForm" val="آخِر"/>
</lemme>
<meaning>
<definition>
<feat att="text" val="مختلف، متباير أو بمعنى غيره"/>
</definition>
</meaning>
<meaning>
<definition>
<feat att="text" val="أحد شيئين يكونان من جنس واحد"/>
</definition>
</meaning>
</entry>
<entry id="2276">
<feat att="partOfSpeech" val="فعل"/>
<root>
<feat att="writtenForm" val="ون"/>
</root>
<lemme>
<feat att="writtenForm" val="آن"/>
</lemme>
<meaning>
<definition>
<feat att="text" val="أَنْ الماء: ضَبَّه"/>
</definition>
</meaning>
<meaning>
<definition>
<feat att="text" val="أنا خال هذا الفرس: صاحبها"/>
</definition>
</meaning>
</entry>
</Dictionary>
```

Conclusion

- "AlkhalilDWS" meets the requirements of Arabic lexicography. It offers many features and lexical resources that make the construction and updating of Arabic dictionaries easier.
- We have adopted the LMF standard to ensure the exchange and communication between produced dictionaries and any tool that complies with this standard
 - ➔ In the future, we plan to enrich "AlkhalilDWS" with new features as well as extracting good examples from the corpus, allowing the chief editor to define the structure of the lexical entry.

Thank you for your attention