



كلية دار العلوم

قسم علم اللغة والدراسات السامية والشرقية

# مُدَوَّنَةٌ مُعْجَمٍ عَرَبِيٍّ مُعَاَصِرٍ

## مُعَالَجَةٌ لُغَوِيَّةٌ حَاسُوبِيَّةٌ

أطروحة مُقَدِّمَةٌ لِلْحَصُولِ عَلَى دَرَجَةِ الْمَاجِسْتِيرِ فِي عِلْمِ اللُّغَةِ

إعداد

المُعْتَرِّضُ بِاللَّهِ السَّعِيدُ طه

المعيد بالقسم

1428هـ / 2007م

إشراف

الأستاذ الدكتور / محمد حسن عبد العزيز

أستاذ علم اللغة - بجامعة القاهرة

وعضو مجمع اللغة العربية بالقاهرة

وإشراف

الأستاذة الدكتورة / سلوى السيد حمادة

أستاذة الإلكترونيات والاتصالات بالمركز القومي للبحوث

١٩٨٨ هـ - ١٤١١ هـ

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



## الفهرس الموضوعي

أ	مقدمة
أ	موضوع الدراسة
أ	أهداف الدراسة
ب	منهج الدراسة
ب	الدراسات السابقة
ج	محاوِر الدراسة
1	الباب الأول (الإطار العام لمُدونة المعجم العربي المعاصر)
2	الفصل الأول (مباحث تمهيدية)
3	منهج دراسة المُدونة اللُّغوية
5	إرهاصات المنهج، ومراحل تطوُر دراسة المُدونات اللُّغوية
6	المرحلة الأولى: ما قبل ظهور الحاسب الآلي
12	المرحلة الثانية: المُدونات الإلكترونية في أطوارها الأولى
15	المرحلة الثالثة: مرحلة المشروعات اللُّغوية الكبرى
18	أنواع المُدونات الإلكترونية
20	المُدونات العربية المنجزة إلكترونياً
29	الفصل الثاني (الخطوات المنهجية لصناعة مُدونة معجم عربيٍّ معاصر)
30	تحديد مادة المُدونة اللُّغوية وتعيين مصادرها
42	إدخال مادة المُدونة اللُّغوية
42	أولاً: الإدخال الآلي
53	ثانياً: الإدخال اليدوي
54	ثالثاً: استخدام المادة المتاحة إلكترونياً
61	المراجعة الإملائية لنُصوص المُدونة اللُّغوية
63	تنسيق المادة المُدخلة تمهيداً لمعالجتها آلياً
64	ترميز نُصوص المُدونة اللُّغوية
66	لغة ترميز النُصوص XML

67	المكوّنات الأساسية لُغة الترميز XML .....
72	تصميم قاعدة بيانات المدوّنة اللُّغويّة .....
78	نماذج من نُصوص المدوّنة اللُّغويّة – موضوع الدّراسة .....
83	الفصل الثالث (المعالجة الآليّة لنُصوص المدوّنة اللُّغويّة) .....
84	مفهوم المعالجة الآليّة للُّغات الطّبيعيّة .....
85	مفهوم الخوارزمية، وعلاقتها بالمعالجة الآليّة للُّغة العربيّة .....
86	أدوات التّحليل الآليّ لنُصوص المدوّنات اللُّغويّة .....
86	أولاً: المُفهرس الآليّ .....
99	ثانياً: المُحلّل الصّرفيّ .....
116	ثالثاً: مُعنّون التّراكيب العربيّة .....
119	رابعاً: المعجم الإلكترونيّة المتّاحة على الشّبكة العنكبوتيّة .....
121	المعالجة الآليّة لمدوّنة المعجم المنشود .....
122	الباب الثاني (المعجم العربيّ المعاصر) .....
123	الفصل الأول (الخطوات التّنفيدية لصناعة معجم عربيّ معاصر) .....
125	تعيين المداخل المعجميّة، وتحديد الكلمات الرّأسيّة .....
128	تعيين معلومات التّكرار .....
131	تعيين المعاني المعجميّة لمدخل النّمودج المعجميّ .....
134	إدراج الأمثلة التّوضيحيّة لمدخل النّمودج المعجميّ .....
139	تعيين الجوانب البراجماتيّة .....
142	الكتابة الصّوتيّة لمدخل النّمودج المعجميّ .....
145	نمودج توضيحيّ يُوازن بين المُخرَج الحاسوبيّ والمُخرَج المعجميّ .....
149	الفصل الثاني (النّمودج المعجميّ، باب الباء) .....
301	الخاتمة .....
303	المراجع .....

## مقدمة

الحمد لله رب العالمين، والصلاة والسلام على محمد خاتم النبيين، وبعد...

فتشتمل مقدمة هذه الدراسة على العناصر التالية:

- موضوع الدراسة.
- أهداف الدراسة.
- منهج الدراسة.
- الدراسات السابقة.
- محاور الدراسة.

### أولاً: موضوع الدراسة:

موضوع هذه الدراسة (مدونة معجم عربي معاصر، معالجة لغوية حاسوبية)، وهو محاولة لتصميم مدونة لغوية Corpus Linguistics للعربية المعاصرة، ومعالجتها آلياً للإفادة منها في صناعة معجم عربي معاصر.

### ثانياً: أهداف الدراسة:

وتهدف الدراسة إلى:

1. التعريف بمنهج لغوي، حديث نسبياً، هو منهج دراسة المدونات اللغوية الإلكترونية Electronic Corpus Linguistics.
2. تصميم مدونة لغوية إلكترونية تعكس الواقع اللغوي للعربية في الوقت الراهن.
3. تقديم نموذج معجم عربي معاصر، يستمد مادته بأكملها من مدونة لغوية إلكترونية، ويفي بحاجة أرباب العربية ومتعلميها، في محاولة لسد بعض الفجوات المعجمية التي تعاني منها المعاجم العربية المعاصرة.

ثالثاً: منهج الدراسة:

صُمِّمَتِ المَدُونَةُ اللُّغَوِيَّةُ - مَوْضُوعُ الدَّرَاسَةِ - وَفَقَ المَنهَجُ الوَصْفِيُّ، وَيَتَنَوَّعُ مَجَالُ البَحْثِ بَيْنَ عِلْمِ اللُّغَةِ الحَاسُوبِيِّ Computational Linguistics الَّذِي يَقُومُ عَلَى إِخْضَاعِ الآلَةِ لِلتَّحْلِيلِ اللُّغَوِيِّ بِكَافَّةِ مُسْتَوِيَاتِهِ بِاسْتِخْدَامِ أَدَوَاتٍ حَاسُوبِيَّةٍ مُعَيَّنَةٍ، وَهَنْدَسَةِ اللُّغَةِ Language Engineering الَّتِي تُعْنَى بِتَصْمِيمِ هَذِهِ الأَدَوَاتِ وَتَطْوِيرِهَا وَفَقاً لِطَبِيعَةِ اللُّغَةِ مَوْضُوعِ الدَّرَاسَةِ. وَوُضِعَ النَّمُودَجُ المَعْجَمِيُّ المُلْحَقُ بِالدَّرَاسَةِ - كَذَلِكَ - وَفَقَ المَنهَجِ الوَصْفِيِّ، إِذْ اسْتَمَدَّ مَادَّتَهُ مِنْ مَدُونَةِ لُغَوِيَّةِ العَرَبِيَّةِ الفُصْحَى المَعَاصِرَةِ، تَضُمُّ نُصُوصاً تَنْتَمِي إِلَى الفَتْرَةِ مِنْ 1950 م، إِلَى 2007 م.

رابعاً: الدراسات السابقة:

لَمْ يَتِمَّ كُنَّ البَاحِثُ مِنَ العُثُورِ عَلَى دَرَاثٍ عَرَبِيَّةٍ عَنِ المَدُونَاتِ اللُّغَوِيَّةِ العَرَبِيَّةِ، وَرُبَّمَا يَرْجِعُ ذَلِكَ إِلَى جِدَّةِ المَنهَجِ نَسَبِيًّا؛ لَكِنَّ هُنَاكَ عِدَدًا مِنَ المَدُونَاتِ اللُّغَوِيَّةِ المَوْضُوعَةِ ضِمْنَ دَرَاثٍ غَيْرِ عَرَبِيَّةٍ لِبَاحِثِينَ عَرَبٍ، وَمِنْ هَذِهِ المَدُونَاتِ:

1. المَدُونَةُ مُتَعَدِّدَةُ اللُّغَاتِ Multilingual Corpus، صُمِّمَتِ ضِمْنَ الأَطْرُوحَةِ العِلْمِيَّةِ الَّتِي تَقَدَّمَ بِهَا البَاحِثُ سَتَّارُ عَزُوزِي نِي إِلَى جَامِعَةِ مَانَشِسْتَرِ لِلْحُصُولِ عَلَى دَرَجَةِ الدُّكْتُورَاهِ فِي (2003 م).
2. المَدُونَةُ العَرَبِيَّةُ العِلْمِيَّةُ العَامَّةُ General Scientific Arabic Corpus، صُمِّمَتِ ضِمْنَ الأَطْرُوحَةِ العِلْمِيَّةِ الَّتِي تَقَدَّمَ بِهَا البَاحِثُ أَمِينُ المُهَنِّي إِلَى جَامِعَةِ مَانَشِسْتَرِ لِلْحُصُولِ عَلَى دَرَجَةِ الدُّكْتُورَاهِ فِي (2003 م).
3. مَدُونَةُ العَرَبِيَّةِ الفُصْحَى Classical Arabic Corpus، صُمِّمَتِ ضِمْنَ الأَطْرُوحَةِ العِلْمِيَّةِ الَّتِي تَقَدَّمَ بِهَا البَاحِثُ عَبْدِ الحَامِدِ عَلِيَّةِ إِلَى جَامِعَةِ مَانَشِسْتَرِ لِلْحُصُولِ عَلَى دَرَجَةِ الدُّكْتُورَاهِ فِي (2004 م).

4. مَدَوْنَةُ الْعَرَبِيَّةِ الْمُعَاَصِرَةِ Contemporary Arabic Corpus، صُمِّمَتْ  
ضِمْنَ الْأَطْرُوحَةِ الْعِلْمِيَّةِ الَّتِي تَقَدَّمَتْ بِهَا الْبَاحِثَةُ الْقَطْرِيَّةُ لَطِيفَةُ السُّلَيْطِي -  
أَسْتَاذُ الدَّرَاسَاتِ اللُّغَوِيَّةِ بِجَامِعَةِ قَطْرَ - إِلَى جَامِعَةِ لِيدزِ لِلْحُصُولِ عَلَى دَرَجَةِ  
الْمَاجِسْتِيرِ فِي الْعُلُومِ فِي (2004 م).  
5. مَدَوْنَةُ تَطْوِيرِ تَقْنِيَةِ مُحَاكَاةِ الْأَصْوَاتِ Chatbot Corpus، صُمِّمَتْ ضِمْنَ  
الْأَطْرُوحَةِ الْعِلْمِيَّةِ الَّتِي تَقَدَّمَتْ بِهَا الْبَاحِثَةُ الْأُرْدُنِيَّةُ بِيَانِ عَارِفِ أَبُو شَاوَرِ إِلَى  
جَامِعَةِ لِيدزِ لِلْحُصُولِ عَلَى دَرَجَةِ الدُّكْتُورَاهِ فِي (2005 م).  
وَقَدْ تَعَرَّضَ الْبَاحِثُ لِهَذِهِ الْمَشْرُوعَاتِ الْفَرْدِيَّةِ وَغَيْرِهَا مِنْ الْمَشْرُوعَاتِ الْمُنْجَزَةِ بِوِاسِطَةِ  
مُؤَسَّسَاتٍ أَوْ هَيْئَاتٍ عِلْمِيَّةٍ، بِشَيْءٍ مِنَ التَّفْصِيلِ فِي الْبَابِ الْأَوَّلِ مِنْ هَذِهِ الدَّرَاسَةِ.

### خَامِسًا: مَحَاوِرُ الدَّرَاسَةِ:

وَسَعِيًّا وَرَاءَ بُلُوغِ الْمَهْدَفِ الْمَنْشُودِ، فَقَدْ قَسَمَ الْبَاحِثُ الدَّرَاسَةَ إِلَى بَابَيْنِ، تَعْقِبُهَا خَاتِمَةٌ،  
عَلَى النَّحْوِ التَّالِي:

- الْبَابُ الْأَوَّلُ: الْإِطَارُ الْعَامُّ لِمَدَوْنَةِ الْمُعْجَمِ الْعَرَبِيِّ الْمُعَاَصِرِ. وَيَشْتَمِلُ عَلَى ثَلَاثَةِ  
فُصُولٍ:

- الْفَصْلُ الْأَوَّلُ: مَبَاحِثُ تَمْهِيدِيَّةٍ. وَيُنَاقِشُ بَعْضَ الْقَضَايَا النَّظَرِيَّةِ الَّتِي تَتَعَلَّقُ  
بِمَوْضُوعِ الدَّرَاسَةِ، وَيُضْمُّ أَرْبَعَةَ مَحَاوِرَ أُسَاسِيَّةٍ:

1. مَنَهْجُ دِرَاسَةِ الْمَدَوْنَةِ اللُّغَوِيَّةِ (Corpus Linguistics).

2. إِرهَاصَاتُ الْمَنَهْجِ، وَمَرَاجِلُ تَطَوُّرِ دِرَاسَةِ الْمَدَوْنَاتِ اللُّغَوِيَّةِ.

- الْمَرْحَلَةُ الْأُولَى: مَا قَبْلَ الْحَاسِبِ الْآلِيِّ.

- الْمَرْحَلَةُ الثَّانِيَّةُ: الْمَدَوْنَاتُ الْإِلِكْتَرُونِيَّةُ فِي أَطْوَارِهَا الْأُولَى.

- الْمَرْحَلَةُ الثَّلَاثَةُ: مَرْحَلَةُ الْمَشْرُوعَاتِ اللُّغَوِيَّةِ الْكُبْرَى.

3. أنواع المُدَوَّنات اللُّغَوِيَّة الإلِكْترونيَّة.

4. المُدَوَّنات العرَبِيَّة المنجَزَة إلكْترونيًّا.

• الفصل الثَّاني: الخُطوات المنهجِيَّة لِصِناعة مُدَوَّنة مُعْجَمٍ عرَبِيٍّ مُعاصِرٍ، وَيُضْمُّ سِتَّةَ مَحاورٍ أَساسِيَّة:

1. تحديدها مادة المُدَوَّنة اللُّغَوِيَّة وتعيين مصادرها .

2. إدخال نصوص المُدَوَّنة اللُّغَوِيَّة.

- الإدخال الآلي للنصوص .

- الإدخال اليدوي للنصوص .

- النصوص المتاحة إلكترونيًّا.

3. المراجعة الإملائية لنصوص المُدَوَّنة اللُّغَوِيَّة.

4. تنسيق المادة المدخلة تمهيدًا لمعالجتها آليًّا.

5. ترميز النصوص المدخلة **Text Encoding**.

6. تصميم قاعدة بيانات المُدَوَّنة اللُّغَوِيَّة **Corpus DB**.

وتأتي تَتَمَّةُ الفصل لتضمَّ نماذج من المُدَوَّنة اللُّغَوِيَّة - موضوع الدراسة.

• الفصل الثالث: المُعالجة الآليَّة لنصوص المُدَوَّنة اللُّغَوِيَّة، وَيُضْمُّ أربعةَ مَحاورٍ أَساسِيَّة:

1. مفهوم المُعالجة الآليَّة للُّغات الطَّبِيعِيَّة.

2. مفهوم الخوارزمية وعلاقتها بالمُعالجة الآليَّة للُّغة العرَبِيَّة.

3. أدوات التَّحليل الآلي لنصوص المُدَوَّنة اللُّغَوِيَّة.

- المُفهرس الآلي للنصوص **Concordance**.

- المُحلِّل الصَّرفي العرَبِي **Arab Morpho**.

- مُعَنون التَّراكيب العرَبِيَّة **Arab Tagger**.

– المعاجم العربية المتاحة على الشبكة العنكبوتية.

#### 4. المعالجة الآلية لمُدونة المعجم المنشود.

- الباب الثاني: المعجم العربي المعاصر (الخطوات التنفيذية، والنموذج المعجمي). ويشتمل على فصلين:

- الفصل الأول: الخطوات التنفيذية لصناعة معجم عربي معاصر، ويضم سبعة محاور أساسية:

1. تعيين المداخل المعجمية، وتحديد الكلمات الرأسيّة.

2. تعيين معلومات التكرار **Frequency information**.

3. تعيين المعاني المعجمية لمداخل النموذج المعجمي.

4. إدراج الأمثلة التوضيحية لمداخل النموذج المعجمي.

5. تعيين الوصف النحوي بجانبه النحوي والتركيبي.

6. تعيين الجوانب البراجماتية لمداخل النموذج المعجمي.

7. الكتابة الصّوتية لمداخل النموذج المعجمي.

وتأتي تيمّة الفصل لتعرض نموذجاً يوازن بين المخرج الحاسوبّي والمخرج المعجمي.

- الفصل الثاني: النموذج المعجمي، وقد اشتمل على المداخل المعجمية لباب

(الباء) بكامل كلماته الرأسيّة الواردة في نصوص المدونة اللغوية موضوع

الدراسة، وقد وصل عدد المداخل المعجمية للنموذج المعجمي إلى مائتين

واثنين وخمسين (252) مدخلاً معجمياً، تبدأ بمدخل (ب)، وتنتهي بمدخل

(بي)؛ وتشتمل هذه المداخل على سبعمائة وأحد عشر (711) كلمة رأسيّة.

- الخاتمة، وتشتمل على أهمّ النتائج والتوصيات.

والله من وراء القصد وهو يهدي السبيل ...

رموز الكتابة الصوتية المستخدمة

الرمز الصوتي	الحرف	الرمز الصوتي	الحرف
ʕ	ع	ʔ	أ
ɣ	غ	b	ب
f	ف	t	ت
q	ق	ʔ	ث
k	ك	g	ج
l	ل	h	ح
m	م	ĥ	خ
n	ن	d	د
h	هـ	ɗ	ذ
w	و	r	ر
y	ي	z	ز
a	فتحة قصيرة	s	س
ā	فتحة طويلة	ʃ	ش
e	كسرة قصيرة	ʂ	ص
ē	كسرة طويلة	d'	ض
u	ضمّة قصيرة	ʦ	ط
ū	ضمّة طويلة	ʒ	ظ

الاختصارات المستخدمة

الاختصار	المصطلح
<b>BNC</b>	<b>British National Corpus</b>
<b>CAC1</b>	<b>Contemporary Arabic Corpus</b>
<b>CAC2</b>	<b>Classical Arabic Corpus</b>
<b>COBUILD</b>	<b>Collins Birmingham University Language Database</b>
<b>EAC</b>	<b>English – Arabic Corpus</b>
<b>ELRA</b>	<b>European Language Resources Association</b>
<b>GSAC</b>	<b>General Scientific Arabic Corpus</b>
<b>HTML</b>	<b>Hyper Text Markup Language</b>
<b>LDC</b>	<b>Linguistic Data Consortium</b>
<b>LOB</b>	<b>Lancaster-Oslo/ Bergen Corpus</b>
<b>MLC</b>	<b>Multilingual Corpus</b>
<b>OCR</b>	<b>Optical Character Readers</b>
<b>OED</b>	<b>Oxford English Dictionary</b>
<b>SDC</b>	<b>System Development Corporation</b>
<b>SSE</b>	<b>Survey Of Spoken English</b>
<b>SUE</b>	<b>Survey of English Usage</b>
<b>XML</b>	<b>Extensible Markup Language</b>
<b>CLARA</b>	<b>Corpus Linguae Arabicae</b>
<b>PDF</b>	<b>Portable Document Format</b>
<b>TXT</b>	<b>Text Document</b>
<b>RTF</b>	<b>Rich Text Format</b>
<b>DOC</b>	<b>Word Document</b>
<b>ASP</b>	<b>Active server pages</b>
<b>PHP</b>	<b>Personal Home Page</b>
<b>JSP</b>	<b>Java Server Pages</b>

# الباب الأول

## الإطار العام لمُدَوَّنَةِ الْمُعْجَمِ الْعَرَبِيِّ الْمُعَاْصِرِ

## الفصل الأول: مَبَاحِثُ تَمْهِيدِيَّة

1. منهج دراسة المدونة اللغوية (Corpus Linguistics).
2. إرهاصات المنهج، ومراحل تطوُّر دراسة المدونات اللغوية.
  - المرحلة الأولى: ما قبل الحاسب الآلي.
  - المرحلة الثانية: المدونات الإلكترونية في أطوارها الأولى.
  - المرحلة الثالثة: مرحلة المشروعات اللغوية الكبرى.
3. أنواع المدونات اللغوية الإلكترونية.
4. المدونات العربية المنجزة إلكترونياً.

## 1. منهج دراسة المدونة اللغوية

### Corpus Linguistics Study

يمكن تعريف المدونة اللغوية بأنها كتلة غير منتظمة من النصوص المكتوبة أو المنطوقة التي تُستخدم لدراسة جوانب اللغة، يمكن قراءتها والتعامل معها آلياً بعد إدخالها على الحاسب الآلي (1)، كما يمكن التحكم في بياناتها ومُدخلاتها، بالإضافة أو الحذف أو التعديل من خلال قواعد بيانات (Databases) صُممت خصيصاً للتعامل مع هذه النصوص. وتعتبر قاعدة البيانات الحاوية لنصوص المدونة اللغوية مخزناً كبيراً للغة، يُرجع إليه وقت الحاجة، ويتحمل أي قدر من النصوص التي تُضاف إلى المادة الأساسية مستقبلاً.

ومادة المدونة اللغوية ليست نصوصاً تقيديّة أو عشوائية؛ إنها كتلة غير منتظمة من النصوص التي تخضع لمجموعة من الأسس والمعايير، يُحددها الهدف المنشود من المدونة اللغوية؛ فالمدونة التي يُعتمد عليها في صناعة معجم لغويّ، تختلف مادتها عن تلك المستخدمة في حصر مجموعة من الأنماط التركيبية أو البنيوية للغة، كما تختلف مادة المدونة المستخدمة في صناعة معجم تكراريّ عن تلك التي يُعتمد عليها في صناعة المعاجم التاريخية. كذلك.. فإنّ المعالجة الآلية للنصوص تتفق وطبيعة المدونة؛ فالبرامج الحاسوبية المستخدمة، وطريقة معالجة النصوص، وطرائق إدارة قواعد البيانات، كلّ هذا يخضع لتلك الأسس والمعايير التي تُحددها طبيعة المدونة اللغوية.

(1) **Hartmann, R. R. K. and Stork, F. C.** (1972). Dictionary Of Language and Linguistics. London. P.55, **Kennedy, G.** (1998). An Introduction to Corpus Linguistics. Longman. P.1, and For More Information About Corpus Linguistics Definition, See: **McEnery, T. and Wilson, A.** (1997). Corpus Linguistics Edinburgh: Edinburgh University Press. P. 21. **Kübler, S.** (2005). Introduction to Corpus Linguistics. University of Tübingen. P.2. It is Available from: [http://www.sfs.uni-tuebingen.de/~kuebler/rocoli/intro\\_corp\\_ling.pdf](http://www.sfs.uni-tuebingen.de/~kuebler/rocoli/intro_corp_ling.pdf).