

بناء وتطوير ذخيرة لغوية مُحَرِّك بحث عربيّ مفتوح المصدر

دكتور المُعْتزّ بالله السَّعيد

ورقة عمل مُقدَّمة إلى

المنظمة العربية للتربية والثقافة والعلوم
ALECSO

اجتماع خبراء مُحَرِّكات البحث العربية
تونس، 26-27 أبريل 2017

1. توصيف الذخيرة المنجزة:

- استمدت الذخيرة *Linguistic Corpus* مادتها من موردين رئيسيين؛ على النحو الآتي:

| م | المورد | عدد الوثائق | عدد الكلمات | الكلمات الفريدة |
|---|------------------|-------------|-------------|-----------------|
| 1 | موسوعة ويكيبيديا | 306366 | 42894090 | 931790 |
| 2 | صحيفة الحياة | 26919 | 20072264 | 408681 |

- تمّت تهيئة الذخيرة في صورتين؛ إحداهما: الذخيرة بصورتها الخام *Raw Corpus*، والأخرى في قاعدة بيانات، بصيغة *XML*.
- تمّ استخلاص قوائم الكلمات الفريدة *Unique Words* لكلّ وثيقة على حدة؛ وتمّ التّعامل معها باعتبارها كلمات مفتاحية *Key Words*.
- تمّ استخلاص إحصاءات الذخيرة. واشتملت هذه الإحصاءات على:
 1. إحصاءات الكلمات.
 2. إحصاءات الكلمات الفريدة.
 3. إحصاءات تردّدات الكلمات الفريدة.
 4. إحصاءات وثائق الذخيرة، مع تصنيف المقالات.

2. نموذج من الدَّخيرة المنجزة:

| أيون الهيدروجين | | عنوان المقالة |
|--|-------------|---------------------------|
| موسوعة ويكيبيديا | | المورد |
| <p>أيون الهيدروجين. أوصى الاتحاد الدولي للكيمياء البحتة والتطبيقية IUPAC بمصطلح عام لأيونات الهيدروجين ونظائره بالاعتماد على شحنة الأيون.</p> <p>1- أيون موجب الشحنة. 2- أيون سالب الشحنة.</p> <p>وتحت ظروف مائية معينة كما توجد في الكيمياء الحيوية توجد أيونات الهيدروجين في صورة هيدرات مكونة جزيئات الهيدرونيوم H_3O^+، ويسمى غالباً أيون الهيدروجين أو بروتون في الكيمياء الحيوية.</p> <p>- أيون الهيدروجين: له شحنة أولية موجبة واحدة (H^+)</p> <p>- البروتون: H^+ هو أيون الهيدروجين، ويحمل شحنة أولية واحدة موجبة، يستخدم هذا التعبير مرادفاً لأيون الهيدروجين.</p> <p>- ديوترون: D^+، وهو أيون الهيدروجين الثقيل (نظير) حيث تتكون نواته ليس من بروتونا واحداً فقط وإنما يتكون من 1 بروتون و 1 نيوترون لهاذا يسمى الهيدروجين الثقيل، ويرمز إليه بالرمز D، وهو يحمل شحنة أولية واحدة موجبة مثل أيون الهيدروجين.</p> <p>...</p> | | نموذج التُّصوُّص |
| 179 | | إحصاء الكلمات |
| 98 | | إحصاء الكلمات الفريدة |
| التردد | الكلمة | الكلمات الفريدة المفتاحية |
| 16 | الهيدروجين | |
| 14 | أيون | |
| 5 | شحنة | |
| 3 | بروتون | |
| 3 | موجبة | |
| 2 | الحيوية | |
| 2 | الكيمياء | |
| 2 | الهيدرونيوم | |

3. موارد مُقترحة لتطوير الدَّخيرة:

- تقترحُ الورقة إثراءَ مادَّة الدَّخيرة بِنُصوص العربيَّة المُعاصرة المُستمدَّة من الموارد الآتية:

1. وثائق مُؤتمر (تيد *TED*):

- تأتي هذه المادَّة تمثيلاً للغة العربيَّة العلميَّة.
- يُتوقَّع ألا يقلَّ عددُ كلمات هذه المادَّة عن 2 مليون كلمة.
- تضمُّ التَّرجمة العربيَّة لمُحاضرات المُؤتمر الأكاديمي (تيد *TED*).
- يمتدُّ الإطارُ الزَّمنيُّ لهذه النُّصوص من عام (2012م) حتَّى عام (2017م).

2. وثائق مُنظمة الأمم المتَّحدة (*UN*):

- تأتي هذه المادَّة تمثيلاً للغة العربيَّة الرِّسميَّة.
- يُتوقَّع ألا يقلَّ عددُ كلمات هذه المادَّة عن 5 ملايين كلمة.
- تضمُّ النُّسخة العربيَّة لقرارات مُنظمة الأمم المتَّحدة (*UN*).
- يمتدُّ الإطارُ الزَّمنيُّ لهذه النُّصوص من عام (2000م) حتَّى عام (2010م).

3. وثائق الصَّحافة العربيَّة المُعاصرة:

- تأتي هذه المادَّة تمثيلاً للغة العربيَّة العامَّة.
- يُتوقَّع ألا يقلَّ عددُ كلمات هذه المادَّة عن 10 ملايين كلمة.
- تضمُّ نُصوص المقالات الصَّحفيَّة المُستمدَّة من الصُّحف العربيَّة.
- يمتدُّ الإطارُ الزَّمنيُّ لهذه النُّصوص من عام (2000م) حتَّى عام (2017م).

4. وثائق الموقع المفتوح (ويكي هاو *Wikihow*):

- تأتي هذه المادَّة تمثيلاً للغة العربيَّة العلميَّة.
- يُتوقَّع ألا يقلَّ عددُ كلمات هذه المادَّة عن 100 ألف كلمة.
- تضمُّ النُّسخة العربيَّة لمقالات موقع الويكي (*Wikihow*) المعرفي.
- يمتدُّ الإطارُ الزَّمنيُّ لهذه النُّصوص من عام (2008م) حتَّى عام (2017م).

5. وثائق الموسوعة الحرّة ويكيبيديا "مُستدرّكة":

| |
|--|
| - تأتي هذه المادّة تمثيلاً للغة العربيّة الموسوعيّة. |
| - يُتوقّع ألا يقلّ عددُ كلمات هذه المادّة عن 5 ملايين كلمة. |
| - تضمُّ النسخة العربيّة من الموسوعة الحرّة "ويكيبيديا Wikipedia" |
| - يمتدُّ الإطارُ الزمّنيُّ لهذه النُصوص من عام (2014م) حتّى عام (2017م). |

4. المهامّ والمخرجات والإطار الزمّني لبناء وتطوير الدّخيرة اللّغويّة:

| المرحلة | المهامّ | المُخرجات | الإطار الزمّنيّ | |
|------------------------------------|---|---|--|---|
| 1 | استيراد وتنقية الدّخيرة | الدّخيرة اللّغويّة الخام [XML]. وثائق الدّخيرة المُنقّاة من الكشائد والرّوائد. | 4 أشهر | |
| | 2 | العمليّات الإحصائيّة | تحليل <i>N-Gram</i> | الكلمات الأحاديّة (1) المُتلازمات الثنائيّة (2) المُتلازمات الثلاثيّة (3) |
| إحصاءات الكلمات والمُتلازمات | | | قائمة الكلمات الفريدة <i>Wordlist – 01 Gram</i> قائمة الثنائيّات الفريدة <i>Wordlist – 02 Gram</i> قائمة الثلاثيّات الفريدة <i>Wordlist – 03 Gram</i> | 4 أشهر |
| 3 | | | العمليّات الحاسوبية | الفهرسة الآليّة |
| | التّعريف الآليّ على عناوين الوثائق | قائمة عناوين وثائق الدّخيرة | | 4 أشهر |
| 4 | العمليّات اللّغويّة | الوسم التّركيبيّ | نماذج من الدّخيرة اللّغويّة مُوسّمة تركيبياً، توسيمًا آليًا جُزئيًا. | |
| | | التّجذيع <i>Stemming</i> | كلمات الدّخيرة بعد القطع الجُزئيّ للرّوائد | 4 أشهر |

| | | | |
|--|---|-----------|--|
| | نتائج تقييم عمل مُحرك البحث، استنادًا إلى مُخرجات الدَّخيرة | التَّقييم | |
|--|---|-----------|--|

5. فريق العمل المُقترح لبناء وتطوير الدَّخيرة:

| م | الباحث | التَّوصيف |
|---|--------------------------|-------------|
| 1 | المُعتمَر بالله السَّعيد | باحث رئيس |
| 2 | ريم الأَطر | باحث مُساعد |

6. توفيق المهام والمُخرجات مع مُتطلبات المُعالجة اللُّغويَّة وتقييم أداء مُحرك البحث:

| المُخرجات | إجراءات المُعالجة | م | المُعالجة |
|--|---|---|-----------------------|
| قائمة الكلمات الفريدة <i>Wordlist – 01 Gram</i> | توسيع مُعجم النِّظام | 1 | المُعالجة اللُّغويَّة |
| قوائم كلمات الدَّخيرة وتردُّداتها - لكلِّ وثيقة | مُطابقة مُفردات الدَّخيرة بمُفردات المُعجم | 2 | |
| قائمة عناوين وثنائى الدَّخيرة | تعيين أسماء الأعلام غير الموجودة في مُعجم النِّظام | 3 | |
| الكلمات الأحاديَّة (1) | تعيين المُصطلحات الدَّائرة غير الموجودة في مُعجم النِّظام | 4 | |
| المُتلازمات الثُّنائيَّة (2) | | | |
| المُتلازمات الثُّلاثيَّة (3) | | | |
| قائمة الكلمات الفريدة <i>Wordlist – 01 Gram</i> | إحصاءات المُفردات والمُصطلحات | 5 | |
| قائمة الثُّنائيَّات الفريدة <i>Wordlist – 02 Gram</i> | | | |
| قائمة الثُّلاثيَّات الفريدة <i>Wordlist – 03 Gram</i> | | | |
| كلمات الدَّخيرة بعد القطع الجُزئيِّ للزوائد | إحصاء مُفردات الدَّخيرة بعد القطع السَّطحيِّ للسَّوابق واللَّواحق | 6 | تقييم أداء مُحرك |
| كلمات الدَّخيرة بعد القطع الجُزئيِّ للزوائد | تحديد نسبة نجاح القطع السَّطحيِّ | 7 | |
| الدَّخيرة اللُّغويَّة الخام [XML] | بناء قاعدة بيانات خاصَّة بالنُّصوص العلميَّة المُعاصرة | 8 | |
| وثائق الدَّخيرة المُنقَّاة من الكشائد والزوائد | | | |

| | | | |
|--|---|----|--|
| الدَّخيرة اللُّغويَّة الخام [XML] | بناء استعلامات الدَّخيرة | 9 | |
| وثائق الدَّخيرة المُنقَّاة من الكشائد والزَّوائد | | | |
| الدَّخيرة اللُّغويَّة الخام [XML] | المُقارنة بين نتائج البحث الآليَّة والنتائج المُعدَّة يدويًّا | 10 | |
| قوائم كلمات الدَّخيرة وتردُّداتها - لكلِّ وثيقة | | | |

7. نموذج مُخرجات الدَّخيرة اللُّغويَّة لمُحرِّك البحث:
 - يُمَثِّلُ النَّموذجُ الآتي صُورة مُصَغَّرة من الدَّخيرة اللُّغويَّة المنشودة.
 أوَّلاً: إحصاءات الوثائق:

| م | عُنوان الوثيقة | عدد الكلمات | م | عُنوان الوثيقة | عدد الكلمات |
|----|------------------------|-------------|----|---------------------|-------------|
| 1 | آشور | 3940 | 26 | المحقق كونان | 2319 |
| 2 | آن بولين | 3760 | 27 | المشتري | 3948 |
| 3 | أغالبة | 3320 | 28 | الوطن العربي | 4473 |
| 4 | ألماس | 5990 | 29 | بابوية كاثوليكية | 2809 |
| 5 | ألومنيوم | 4607 | 30 | برج إيفل | 3717 |
| 6 | أورانوس | 4517 | 31 | بلاد الشام | 2750 |
| 7 | أور | 2881 | 32 | بنين | 2276 |
| 8 | إدوارد سعيد | 3240 | 33 | بيل غيتس | 3242 |
| 9 | إسحاق شامير | 2941 | 34 | تاريخ الإسكندرية | 5489 |
| 10 | إسحاق نيوتن | 3299 | 35 | تخلخل العظم | 2738 |
| 11 | إسرائيل | 3544 | 36 | تدخين | 3716 |
| 12 | إعصار قمعي | 3956 | 37 | تويوتا | 5498 |
| 13 | إلكترون | 4658 | 38 | جامعة الدول العربية | 3505 |
| 14 | الأرض | 6421 | 39 | جوجل | 4960 |
| 15 | الأندلس | 4073 | 40 | ج.ك. رولينج | 5140 |
| 16 | الإمبراطورية البيزنطية | 4246 | 41 | حديد | 4436 |

| | | | | | |
|------|----------------------|----|------|-------------------|----|
| 4651 | حرب بونيقية ثانية | 42 | 4956 | الدولة المرينية | 17 |
| 4194 | حشرة | 43 | 4156 | الزهرة | 18 |
| 2600 | ديمقراطية | 44 | 3099 | السكري | 19 |
| 2087 | ذكاء اصطناعي | 45 | 5249 | الشمس | 20 |
| 2791 | ربو | 46 | 3485 | الظاهر بيبرس | 21 |
| 3581 | روبوت | 47 | 4323 | القاهرة | 22 |
| 4240 | رومانيا | 48 | 5276 | القضية الفلسطينية | 23 |
| 3942 | زحل | 49 | 4037 | القمر | 24 |
| 4722 | زينون | 50 | 4597 | المجموعة الشمسية | 25 |

| عدد الكلمات | عنوان الوثيقة | م | عدد الكلمات | عنوان الوثيقة | م |
|-------------|-----------------------------|----|-------------|-----------------------|----|
| 4214 | كيمونو | 76 | 6116 | ستيف جوبز | 51 |
| 2754 | لاوس | 77 | 3592 | شركة فورد | 52 |
| 3480 | لغة سريانية | 78 | 4221 | صحافة | 53 |
| 2578 | لوحة أم | 79 | 2638 | صهيونية | 54 |
| 3501 | ليثيوم | 80 | 2992 | ضفدع | 55 |
| 2339 | ماري أنطوانيت | 81 | 3953 | طاقة شمسية | 56 |
| 3573 | ماري كوري | 82 | 2761 | عبد الحميد بن باديس | 57 |
| 4745 | مالكوم إكس | 83 | 5592 | عصر النهضة | 58 |
| 3175 | مايكروسوفت ويندوز | 84 | 3243 | عكا | 59 |
| 2731 | محرك بحث | 85 | 3987 | عيد الحب | 60 |
| 3806 | مدغشقر | 86 | 4177 | غابرييل غارثيا ماركيث | 61 |
| 3638 | مسيحية | 87 | 2240 | غابرييل فوري | 62 |
| 3067 | مندائية | 88 | 4861 | غيشا | 63 |
| 2646 | منظمة التحرير الفلسطينية | 89 | 3266 | فتوحات إسلامية | 64 |
| 12488 | مونستر | 90 | 1827 | فرانكلين روزفلت | 65 |
| 4258 | ميلانو | 91 | 2270 | فلاديمير لينين | 66 |
| 3852 | ميونخ | 92 | 2949 | فولتير | 67 |

| | | | | | |
|------|----------------------------|-----|------|-------------------|----|
| 2663 | نفط | 93 | 3642 | فيروس | 68 |
| 3933 | نيكولاي ريمسكي كورساكوف | 94 | 2679 | فيلكا | 69 |
| 2978 | هيدروجين | 95 | 3380 | فينسنت فان غوخ | 70 |
| 2813 | هيليوم | 96 | 3819 | كارلوس الثعلب | 71 |
| 5024 | ولاية برج بوعريج | 97 | 4128 | كازاخستان | 72 |
| 4608 | يهودية إصلاحية | 98 | 2691 | كهرباء | 73 |
| 2835 | يوفنتوس | 99 | 3263 | كوكب | 74 |
| 5367 | يوهانس برامس | 100 | 3309 | كولاجين | 75 |

ثانياً: إحصاءات الدخيرة:

Raw Corpus: Texts: 100 | Arabic Types: 59498 | Arabic Tokens: 383057

ثالثاً: من قائمة الكلمات وتردّداتها:

| التردد | الكلمة | م |
|--------|-----------|----|
| 8 | التخطيط | 1 |
| 1 | التخفيضات | 2 |
| 2 | التخفيف | 3 |
| 1 | التخفيفات | 4 |
| 10 | التخلص | 5 |
| 1 | التخلق | 6 |
| 1 | التخلل | 7 |
| 14 | التخلي | 8 |
| 7 | التخليق | 9 |
| 1 | التخمير | 10 |
| 1 | التخمين | 11 |
| 2 | التخوف | 12 |
| 1 | التخوم | 13 |
| 4 | التدابير | 14 |
| 2 | التداخل | 15 |
| 1 | التداخلات | 16 |

| | | |
|----|-----------|----|
| 11 | التدخل | 17 |
| 1 | التدخلات | 18 |
| 77 | التدخين | 19 |
| 1 | التدرب | 20 |
| 3 | التدرج | 21 |
| 1 | التدرن | 22 |
| 9 | التدريب | 23 |
| 1 | التدريبية | 24 |
| 9 | التدريجي | 25 |
| 6 | التدريس | 26 |
| 13 | التدفق | 27 |
| 6 | التدفئة | 28 |

رابعًا: قائمة المتلازمات الثنائية وتردُّداتها:

| م | الكلمة | التردد |
|----|--------------------|--------|
| 1 | الشريط البحري | 2 |
| 2 | الشريط الذي | 1 |
| 3 | الشريط الساحلي | 2 |
| 4 | الشريعة الإسلامية | 2 |
| 5 | الشريعة الموسوية | 1 |
| 6 | الشريعة لشمس | 1 |
| 7 | الشريف التلمساني | 1 |
| 8 | الشريف الخليلي | 1 |
| 9 | الشريف حسين | 3 |
| 10 | الشريك التجاري | 1 |
| 11 | الشريك الرئيسي | 1 |
| 12 | الشريكان الرئيسيان | 1 |
| 13 | الشط الغربي | 1 |
| 14 | الشط المالح | 1 |
| 15 | الشطر الثاني | 1 |
| 16 | الشطر الغربي | 1 |

| | | |
|---|-------------------|----|
| 1 | الشعار الحالي | 17 |
| 1 | الشعار باللون | 18 |
| 1 | الشعار ذو | 19 |
| 1 | الشعار هو | 20 |
| 1 | الشعارات المعادية | 21 |
| 1 | الشعاع الكهربائي | 22 |
| 1 | الشعاع تصرف | 23 |
| 1 | الشعائر الإصلاحية | 24 |
| 1 | الشعائر المرتبطة | 25 |
| 3 | الشعائر اليهودية | 26 |
| 1 | الشعب الإسرائيلي | 27 |
| 2 | الشعب الثوري | 28 |

خامساً: قائمة المتلازمات الثلاثية وتردُّداتها:

| م | الكلمة | التردد |
|----|----------------------------|--------|
| 1 | الاتحاد الصهيوني البريطاني | 1 |
| 2 | الاتحاد العالمي لليهودية | 2 |
| 3 | الاتحاد العبري المعهد | 1 |
| 4 | الاتحاد العبري عام | 1 |
| 5 | الاتحاد العبري فرعا | 1 |
| 6 | الاتحاد العربي لكرة | 1 |
| 7 | الاتحاد الفلكي الدولي | 13 |
| 8 | الاتحاد المسيحي الديمقراطي | 1 |
| 9 | الاتحاد اليهودي للتربية | 1 |
| 10 | الاتحاد اليهودي للموسيقى | 1 |
| 11 | الاتحاد اليهودي مركز | 1 |
| 12 | الاتحاد أو الجمهورية | 1 |
| 13 | الاتحاد بدون مشغل | 1 |
| 14 | الاتحاد بين الشركتين | 1 |
| 15 | الاتحاد بين اليمن | 2 |
| 16 | الاتحاد في سنة | 1 |

| | | |
|---|---------------------------|----|
| 1 | الاتحاد للإدارة العامة | 17 |
| 1 | الاتحاد مؤتمرا عالميا | 18 |
| 1 | الاتحاد مؤسسات في | 19 |
| 1 | الاتحاد وقالت بالطبع | 20 |
| 1 | الاتحادات القارية لكرة | 21 |
| 1 | الاتحادات القطرية العربية | 22 |
| 1 | الاتحادي الثالث وذلك | 23 |
| 1 | الاتحادي لتمثيل ألمانيا | 24 |
| 1 | الاتصال الأوروبي بدأ | 25 |
| 1 | الاتصال السريع مثل | 26 |
| 1 | الاتصال الشبكي الشهير | 27 |
| 1 | الاتصال الفكري بين | 28 |

سادسًا: نموذج التّوسيم التّركيبيّ الآليّ لنُصوص الدّخيرة:

[PO] إن [CN] النتائج [RP] التي [VI] تمكن [PO] أن [VI] يحصل [PRE] عليها [CN] الملك [CN] الآشوري، [PO] لا [CN] سيما [PRE] في [CN] التخلص [PRE] من [CN] النفوذ [CN] الميتماني [CO] و [CN] يتمكن [PRE] من اقتسام بلادهم، [PO] أن جعله [VI] يتوجه [AD] نحو توطيد أو اصر علاقاته [CN] السياسية [AD] مع [CN] القوى [CN] السياسية [CN] الفاعلة، [AD] حيث [VI] أقدم [PRE] على [CN] الزواج [PRE] من ابنة [CN] الملك [CN] الكاشي [RP] الذي [VP] كان [VI] يفرض نفوذه [PRE] على [CO] و [PO] قد [VI] حظيت مملكة آشور بملوك خلفوا آشور أو بطلت [CO] و [VP] كانوا [PRE] على مستوى [CN] المسؤولية [CO] و [VP] انتهجوا ذات [CN] الأسلوب [RP] الذي [VP] سار [PRE] عليه [PO] ل [VI] يثمر [PRE] عن [DE] ذلك [AD] خلال [CN] القرن [CNU] التاسع [AD] ق.م. بلوغ مستوى [CN] الإمبراطورية [CN] الآشورية [PRE] ب [CN] كل قوتها ونفوذها [CN] السياسي. [PRE] من [CN] الملوك [CN] الآشوريين [CN] البارزين شلمنصر [CNU] الأول [RP] الذي [VP] دام حكمه [CNU] 1266 - [CNU] 1243 [AD] ق.م.، [CO] و [VI] تطلع [PRE] إلى [CN] توجيه [CN] العديد [PRE] من [CN] الحملات [CN] العسكرية وعمل [PRE] على [CN] استبدال [CN] العاصمة آشور [PRE] ب [CN] مدينة نمرود. [PO] أما [CN] العمل [CN] الأبرز [CO] ف [VP] كان [PRE] على يد [CN] الملك توكلتي نورتا [CNU] 1243 - [CNU] 1221 [AD] ق.م.، [RP] الذي [VP] تمكن [PRE] من [CN] السيطرة [PRE] على بلاد بابل، وتوسيع سلطانه [PRE] في [CN] الجهات [CN] الشرقية [CO] و [CN] الغربية. [PO] لكن [AD] بعد وفاة [DE] هذا [CN] الملك [VP] دخلت آشور [PRE] في مرحلة [CN] الضعف [CN] السياسي، نتيجة لوصول ملوك ضعاف [CN] الشخصية، [EX] غير قادرين [PRE] على إدارة مقاليد [CN] الحكم، [CO] و [VP] استمرت [DE] هذه [CN] الفترة [AD] حوالي [CNU] مائة عام، [CO] حتى بلوغ [CN] الملك تجلات بلاسر [CNU] الأول [CNU] 1116 - [CNU] 1090 [AD] ق.م. [PRE] إلى سدة [CN] الحكم، [PO] ل [VI] تكون [DE] هذه [CN] الفترة مليئة [PRE] ب [CN] الإنجازات [CN] العسكرية [CN] الكبيرة، [AD] حيث [VP] تمكن [PRE] من تحقيق [CN] الانتصارات [CN] المتوالية [PRE] في [CN] الأصقاع [CN] البعيدة، [PRE] في [CN] البحر [CN] الأسود وسواحل آسيا [CN] الصغرى [CO] و [CN] المدن [CN] الفينيقية [PRE] على [CN] الساحل [CN] السوري.