

مقاربة جديدة للتشكيل الآلي باستخدام برنامج الخليل للتحليل الصرفي^(*)

محمد ولد عبد الله ولد بياه، عز الدين مزروعي^(**)،
عبد الحق خواجة، عبد الوافي مزيان، شنوفي أمين

ملخص :

نقدم في هذا البحث مقاربة جديدة للتشكيل الآلي لنصوص اللغة العربية الفصحى. وقد أطلقنا على برنامج التشكيل الآلي المنجز الذي يعتمد هذه المقاربة اسم برنامج الخليل للتشكيل الآلي. و تتكون عملية المعالجة التي تتم في هذا البرنامج من مستويين. يعتمد المستوى الأول على مخرجات برنامج عربي مفتوح المصدر هو برنامج الخليل للتحليل الصرفي، وتحتوي هذه المخرجات على التشكيلات الممكنة للكلمة خارج سياقها من النص. بينما يعتمد المستوى الثاني على منهج إحصائيٍّ باستخدام نماذج ماركوف الخفية (Hidden Markov Models)

(*) تم تقديم هذا البحث في صيغته الأولى مرفقاً ببرنامج التشكيل الآلي إلى جائزة المنظمة العربية للتربية والثقافة والعلوم للابداع والابتكار التقني للباحثين الشبان بالوطن العربي سنة 2012 ونال عليه الباحث محمد ولد عبد الله ولد بياه المركز الثاني في هذه المسابقة.

(**) فريق المعالجة الآلية للغة العربية، كلية العلوم - جامعة محمد الأول - وجدة/المغرب.

حيث قمنا باقتراح تصنیف جديد للكلمات العربية لتمثّل هذه الأصناف (classes) مشاهدات النموذج المارکوفي (observations) في حين تمثّل تسلسلاً علامات التشكيل الممكنة الحالات المخفية للنموذج (hidden states).

ولضبط نتائج المستوى الأول من المعالجة عمدنا إلى إنشاء معجم للكلمات العربية شائعة التكرار، قُمنا بدمجه في إطار قواعد معطيات المحلل الصافي. وقد تم استقاء هذا المعجم من ثمان مدونات عربية متاحة على الشابكة (الإنترنت) وتعطي مختلف المجالات. كما قُمنا بعملية تدريب البرنامج وتقييمه على نصوص متعددة تشمل نصوصاً معاصرة من مدونة نيملار (NEMLAR Written Arabic Corpus) ونصوصاً تراثية من مدونة "تشكيلة".

1. مقدمة

يتميز نظام الكتابة العربية، في أغلب النصوص، بغياب علامات التشكيل، وهي: الصوایت القصيرة (short vowels)؛ أي الفتحة والضمة والكسرة، بالإضافة إلى علامات التنوين والشدّة والسكون. ويؤدي غياب هذه العلامات إلى زيادة معتبرة في غموض النص العربيّ تصل بحسب (Debili and Achour 1998) إلى التباس أكثر من 90% من النص. وعلى الرغم من كون القارئ المتعلّم يتمكّن بسهولة من استرجاع علامات التشكيل الغائبة عن النص اعتماداً على سياق الكلمات وعلى مدى معرفته بالصرف والتحوّل العربيّين، إلا أن النص غير المشكول يمثّل عقبة بالنسبة لمتعلمي العربية من غير الناطقين بها والأشخاص ذوي الصعوبات في التعلم، كما يمثّل مصدراً أساسياً للغموض بالنسبة للحاسوب، وخصوصاً في حالة تطبيقات من قبيل تحويل النص المكتوب إلى كلام (Text to Speech) (Zitouni et al. 2006). فخلافاً للغات الأوروبية التي يسهل فيها نوعاً ما إيجاد التقابل بين المقاطع المكتوبة والфонيمات المنطقية، فإن النص العربيّ غير المشكول لا يحيد عن استرجاع علامات تشكيله لهذا الغرض (Vergyri and Kirchoff 2004). وتشير بعض البحوث من جهة

أخرى إلى أهمية استخدام النصوص المنشورة للرفع من كفاءة نظم التعرّف على الكلام (Messaoudi et al. 2004). كما تبرز أهمية التشكيل الآلي في تطبيقات أخرى من قبيل المدونات الشجرية (Treebanks) وال محللات النحوية .(Mammouri et al. 2006)

ونظراً لأهمية استرجاع علامات التشكيل فقد ظهرت محاولات من قبل جهات عدّة، طيلة العقود الماضية، من أجل تطوير مشكّل آلي ذي كفاءة مقبولة. ويمكن أن نقسم هذه المحاولات، على غرار أغلب تطبيقات المعالجة الآلية للغة، إلى نوعين : أولها محاولات أجزتها شركات تجارية بهدف إنشاء مشكّل آلي مستقل أو مشكّل آلي يمثل جزءاً من تطبيق آخر، كالناطق الآلي أو المدقق الإملائي. ومن أبرز هذه المحاولات يمكن أن نذكر شركة "RDI" ومشكّلها الآلي "ArabDiac" ، المشكّل الآلي لشركة صخر، والمشكّل الآلي لشركة "CIMOS" . كما أطلقت شركة غوغل منذ سنتين خدمة "Tashkeel" ، وهي خدمة مجانية للتشكيل الآلي للنصوص العربية بالنسبة لمتصفحـي الشبـكة (الإنـترنت) غير أن هذه الخـدمة لم تكتب لها الاستـمرارـية لأسبـاب نـجهـلـها. وبالرغم من الأهمـية البـالـغـة لـهـذـهـ الـمـحـاوـلـاتـ، إلاـ أنـ طـبـيعـتـهاـ التـجـارـيـةـ الـاحـتكـارـيـةـ تـحـولـ دونـ الـحـصـولـ عـلـىـ تـفـاصـيلـ كـافـيـةـ عـنـ آلـيـةـ عـمـلـهـاـ وـالمـصـادـرـ الـلـغـوـيـةـ الـتـيـ تـسـتـنـدـ عـلـيـهـاـ، فـضـلاـ عـنـ شـيـفـرـتـهاـ الـمـصـدـرـيـةـ مـاـ يـجـعـلـ مـاـ تـحـسـيـنـهـاـ وـدـمـجـهـاـ فـيـ تـطـبـيقـاتـ أـخـرىـ أـمـرـاـ بـعـيدـ المـنـالـ.

أما النوع الثاني من هذه المحاولات فيتمثل في جهود الباحثين ضمن مشاريع مراكز البحث الأكاديمية، وقد أسفرت هذه الجهود طيلة العقود الماضية عن ظهور محاولات متعددة في هذا المجال. ويساعد توفر أوراق بحثية منشورة حول هذه الأعمال في التعرف على المقاربـاتـ المتـهـجـةـ فـيـهـاـ، وهيـ مـقـارـبـاتـ إـحـصـائـيـةـ فـيـ أـغـلـبـ الـأـحـيـانـ، كـماـ يـسـاعـدـ عـلـىـ أـخـذـ فـكـرـةـ عـنـ كـفـاءـةـ الـبـرـجـيـاتـ المنـجزـةـ فـيـهـاـ تـحـتـ ظـرـوفـ الـاخـتـبـارـ الـتـيـ أـخـضـعـهـاـ مـطـورـوـهـاـ.

وستتناول بالتفصيل في الفقرة الموالية من هذا البحث أهم الأعمال السابقة المنجزة في مجال التشكيل الآلي للغة العربية مع التركيز على المقاربات المعتمدة في هذه الأعمال. وفي الفقرة الثالثة نقدم برنامج التشكيل الآلي الذي أنجزناه، عارضين مختلف مراحل المعالجة التي يقوم بها والنموذج الماركوفي الذي يعتمد عليه. ثم بعد ذلك نعرض تباعاً لمرحلة تدريب البرنامج ونتائج تقييمه والخلاصات التي خرجنا بها من هذا العمل.

2. الأعمال السابقة

يمكن تقسيم المقاربات المتبعة من خلال الأعمال السابقة في التشكيل الآلي للنصوص العربية إلى ثلاثة أقسام:

2.1. المقاربات اللغوية (rule-based approaches)

وفي هذا الإطار ظهرت بعض الأعمال التي تهدف إلى برمجة القواعد اللغوية الصوتية والصرفية والنحوية والإملائية بغية تشكيل الكلمات العربية. ومن بين بوادر هذه الأعمال يمكن أن نذكر الطريقة الواردة في (El-Sadany and Hashish 1988) التي تعتمد على استخدام القواعد الصرفية من أجل تشكيل نصف آلي (semi-automatic) للأفعال العربية. وفي (Debili and Achour 1998) يقدم المؤلفان دراسة عن التشكيل الآلي للنصوص العربية تناولاها من جانب غموض الكلمات المكتوبة ومدى إسهام كل من التحليل المعجمي والتحليل الصرفي والعنونة النحوية في كشف الغموض، وبالتالي في تشكيل كلمات النصوص العربية.

غير أن ارتفاع نسبة الغموض وكثرة القواعد الصرفية والنحوية وتشابكها وعدم توفر مُحَلّل نَحْوِي فعال، فضلاً على أن بعض حالات الالتباس تتطلب تحليلاً دلائياً، كل ذلك أسبابٌ تعيق إيجاد مُشَكّل آليًّا للنصوص العربية لا يعتمد سوى على القواعد اللغوية.

2.ب. المقاربـات الإحصـائية (statistical approaches)

بالنظر إلى النجاح الكبير الذي حققه الطرق الإحصائية في مجالات عدـة من المعالجة الآلية للغـة، كالتعرف الآلي على الكلام والترجمـة الآلـية واسترجـاع المعلومات وغيرها (Manning and Schütze 1999)، فقد انصبت أغلـب الجهـود في التشكـيل الآـلي للغـة العـربية على استـعمال هذه الـطرق. وقد تـنوعـت الأسـاليـب المتـبـعة في هذا الـباب، حيث اهـتم بعض البـاحـثـين بـنـماـذـج إـحـصـائـية تـعـالـج التـشـكـيل عـلـى مـسـتـوى الحـرـوف، وـرـكـز آخـرون عـلـى التـشـكـيل عـلـى مـسـتـوى الكلـمـات، بـينـما اتبـع فـرـيق ثـالـث أـسـلـوب المـزاـوـجـة بـيـنـ الطـرـيقـيـن السـابـقـيـن.

فقد سـجـل إـمام وـفـيـشـر (Emam and Fischer, 2005) بـراءـة اـخـتـرـاع مشـكـلـاً آـلـيـاً لـلـغـة العـربـية يـسـتـندـ فيـ فـكـرـته عـلـى طـرـيقـة التـرـجمـة الآـلـية المـعـتـمـدة عـلـى الـأـمـثـلة (exemplar-based machine translation)، حيث يـتـم الـبـحـث ضـمـن قـوـادـعـ معـطـيـات الـأـمـثـلة بـطـرـيقـة تـرـاتـبـية (hierarchical) عـنـ الجـمـلـة المـرـاد تـشـكـيلـها ثـم عـنـ أـجزـاءـ الجـمـلـة وـصـوـلاً إـلـىـ الـكـلـمـات لـاـخـتـيـارـ أـكـثـرـ التـشـكـيلـات رـجـحـانـاً. كـمـ يـتـم تـطـيـقـ نـمـوذـجـ إـحـصـائـيـ (n-gram model) عـلـى مـسـتـوىـ الحـرـوف لـتـشـكـيلـ الـكـلـمـاتـ التـيـ لـمـ يـتـمـ العـثـورـ عـلـيـهـاـ خـالـلـ عـمـلـيـةـ الـبـحـثـ التـرـاتـبـيـةـ السـابـقـةـ.

غـيرـ بـعـيدـ عـنـ هـذـهـ الفـكـرـةـ نـجـدـ فيـ (Schlippe et al, 2008) مـقارـبةـ تـعـتمـدـ عـلـىـ تقـنـيـةـ التـرـجمـةـ الآـلـيةـ الإـحـصـائـيةـ (statistical machine translation) انـطـلاـقاًـ مـنـ مـدوـنةـ مـتوـازـيـةـ (parallel corpus) مـكـوـنـةـ مـنـ النـصـوصـ المـشـكـولـةـ وـمـيـثـلـتـهاـ غـيرـ المـشـكـولـةـ.

كـمـ قـدـمـ (Gal 2002) مـقارـبةـ مـارـكـوـفـيةـ لـتـشـكـيلـ نـصـوصـ الـلـغـتـيـنـ العـربـيـةـ وـالـعـبرـيـةـ. وـقـدـ اـسـتـخـدـمـ المؤـلـفـ فيـ هـذـاـ الـعـمـلـ نـصـوصـ الـقـرـآنـ الـكـرـيمـ بـالـنـسـبـةـ لـلـعـربـيـةـ وـنـصـوصـ الـكـتـابـ الـمـقـدـسـ بـالـنـسـبـةـ لـلـغـةـ الـعـربـيـةـ. وـتـمـثـلـ الـكـلـمـاتـ غـيرـ المـشـكـولـةـ مـشـاهـدـاتـ النـمـوذـجـ المقـرـحـ فيـ عـمـلـ "Gal"ـ بـيـنـماـ تـمـثـلـ الـكـلـمـاتـ المـشـكـولـةـ

الحالات الخفية للنموذج. غير أن عملية التشكيل الآلي للغة العربية في هذا العمل تقتصر على الصوائت القصيرة فقط دون غيرها من علامات التشكيل.

وفي إطار المقاربations الماركوفية للتشكيل الآلي نُسجل أيضاً أعمال كل من (Deltour 2003) و(Elshafei et al. 2006)، حيث استعرضت الباحثة في (Deltour 2003) جملة من الطرق الإحصائية للتشكيل الآلي على مستوى الحروف والكلمات خلصت فيها إلى أن أفضل نتائج التشكيل في بحثها كانت باستخدام نماذج ماركوف الخفية. كما قدم الباحثون في (Elshafei et al. 2006) مقاربة ماركوفية للتشكيل الآلي تتميز عن مقاربة (Gal 2002) في كونها تعامل مع جميع علامات التشكيل، كما تمت عملية التدريب والتقييم فيها على نصوص متعددة لا على النص القرآني فقط.

وقدم (العامدي وأخرون 2006) برنامج التشكيل الآلي المطور في مدينة الملك عبد العزيز للعلوم والتقنية. ويعتمد هذا البرنامج على التسلسلات الرباعية (quad-grams) للحروف حيث يقوم بتشكيل كل حرف من النص استناداً إلى أعلى احتمالية لتشكيله ضمن التسلسلات الرباعية التي يظهر فيها.

2. ج. المقاربations الهجينية (hybrid approaches)

ونعني بها المقاربations التي تزاوج بين القواعد اللغوية والنماذج الإحصائية بغية استغلال نقاط القوة في كلتا الطريقتين. ومن أهم الأعمال التي أنجزت في هذا الإطار نذكر المشكّل الآلي "ArabiDaic" المطور من قبل شركة "RDI" (عطية 2007) حيث يعتمد هذا البرنامج على المحلل الصرفي "ArabMorpho" وعلى واسم تحديد أجزاء الكلام "ArabTagger"، بالإضافة إلى نموذج إحصائي لاختيار التشكيل الأنسب للكلمات باستخدام خوارزمية A*.

كما قدم المؤلفان في (Nelken and Shieber 2005) طريقة تعتمد على تقنية الآلات محدودة الحالات (finite state automata) لتشكيل النصوص العربية. وتزاوج هذه الطريقة بين نماذج إحصائية ثلاثة على مستوى الكلمات

(word tri-gram) ورباعية على مستوى الحروف (quad-gram letter) وبين نموذج صرفي مُبسط يُعرف على سوابق ولوائح الكلمات.

وفي (Zitouni et al, 2006) عَرَض المؤلفون مشكلاً آلياً للنصوص العربية، وذلك باستخدام مصنف إحصائي يعتمد طريقة الأنتروربيا القصوى (Berger et al. 1996) (Maximum Entropy Classifier) حيث تُستخدم الخصائص الصرفية ووسم أجزاء الكلام لاستنتاج التصنيف الأرجح للكلمات.

وفي (Vergyri and Kirchoff 2004) يتم استرجاع علامات التشكيل بمزاوجة التحليل الصرفي والمعلومات السياقية مع نموذج عنونة إحصائية. إذ يذهب المؤلفان إلى اعتبار التشكيل مسألة عنونة غير مسيرة (Unsupervised Tagging) (Buckwalter 2004). وغير بعيد عن هذه الفكرة يقدم المؤلفان (Habash and Rambow, 2007) طريقة للتشكيل الآلي تعتمد من ناحية على محلّ باكوالتر الصرفي وعلى نموذج إحصائي لاختيار أرجح الحلول ضمن مخرجاته.

كما سبق لنا أن قدمنا (بياه وأخرون، 2011) مقاربة صرفية إحصائية للتشكيل الآلي تعتمد على نموذج ماركوفي لتمثيل تسلسل أوزان الكلمات العربية. ولا يفوتنا أن نشير هنا إلى أنها تشتراك مع المقاربة المقدمة في هذا البحث في كونها تعتمدان معاً على برنامج الخليل للتحليل الصرفي، بينما تختلفان في نقاط عدّة، أهمها: مشاهدات النموذج الماركوفي وحالاته الخفية (Observations and Hidden States). وتحتفل أيضاً مدونتا التدريب المستخدمتان فيها، كما أنها أجرينا في عملنا الحالي تعديلات على برنامج الخليل للتحليل الصرفي بإضافة قاعدة معطيات معجمية تتضمن الكلمات الأكثر تكراراً بُغية ضبطها وتسرير مرحلة المعالجة الصرفية.

3. برنامج التشكيل الآلي المُنجز

نعرض في هذه الفقرة لبرنامج التشكيل الآلي الذي قمنا بإنجازه، وقبل أن نتناول الجوانب المتعلقة بالآلية عمل البرنامج ونتائج تقييمه، لا بأس أن نشير بإيجاز إلى دواعي تسميته "برنامج الخليل للتشكيل الآلي".

3.أ. دواعي التسمية

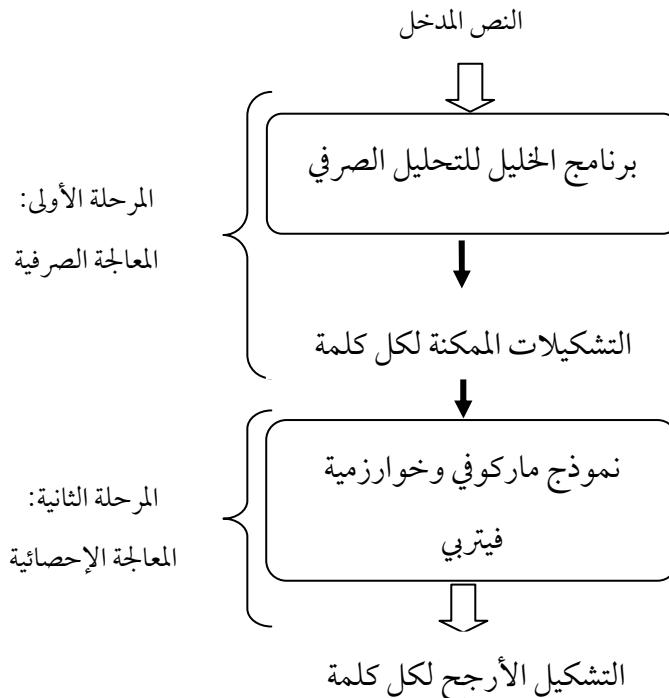
وقع اختيارنا لإطلاق هذه التسمية على برنامج التشكيل الآلي لسبعين: أوهما، للبرنامج علاقة وثيقة ببرنامج الخليل للتحليل الصرفي⁽¹⁾. حيث يعتمد المشكّل الآلي على مُخرجات هذا البرنامج من ناحية وعلى نموذج ماركوفي لاختيار أكثر التشكيلات رجحاناً من ناحية أخرى. أما السبب الثاني فلا يقل وجاهة عن سابقه، إذ إن الفضل في ابتكار علامات التشكيل المعروفة حالياً يرجع إلى الخليل بن أحمد الفراهيدي (ت 170 هـ)، وهي الطريقة التي أثبتت فعاليتها وعمليتها عبر الزمن حتى حصل الاتفاق على كتابة المصاحف بها بدلاً من طريقة النقط المنسوبة إلى أبي الأسود الدؤلي (ت 69 هـ)⁽²⁾.

3.ب. مراحل التشكيل

تمّ عملية التشكيل الآلي في البرنامج المُنجز عبر مرحلتين رئيسيتين كما هو مبين في الشكل 1.

(1) www.sourceforge.net/projects/alkhalil

(2) يذكر أبو عمرو الداني (ت 444 هـ) في كتاب المحكم في نقط المصاحف ما نصه: "وقال أبو الحسن بن كيسان قال محمد بن يزيد الشكل الذي في الكتب من عمل الخليل وهو مأخوذ من صور الحروف فالضمة وأو صغيرة الصورة في أعلى الحرف لثلا تلتبس بالواو المكتوبة والكسرة ياء تحت الحرف والفتحة ألف مبطوحة فوق الحرف.. إلخ".



شكل 1: مراحل المعالجة في برنامج الخليل للتشكيل الآلي

وفي ما يلي نعرض لهذه المراحل بالتفصيل:

1. مرحلة المعالجة الصرافية

وتمّ المعالجة في هذه المرحلة باستخدام برنامج الخليل للتحليل الصرفي، وتهدف إلى استخراج جميع التشكيلات الممكنة لكل كلمة من كلمات النص المدخل. ولمزيد من إيضاح هذه المرحلة نقدم نبذة موجزة عن برنامج الخليل للتحليل الصرفي وأالية عمله قبل أن نتناول التفاصيل المتعلقة بقاعدة المعطيات المعجمية التي أضفناها إليه.

نبذة عن برنامج الخليل للتحليل الصرفي

يعتبر برنامج الخليل للتحليل الصرفي (مزروعي وآخرون، 2011) أحد أهم المحللات الصرفية العربية مفتوحة المصدر. وقد تم إنجازه في إطار مشروع شراكة بين كل من جامعة محمد الأول والمنظمة العربية للتربية والثقافة والعلوم ومدينة الملك عبد العزيز للعلوم والتكنولوجيا. كما روّعيت عند إنجازه التوصيات والمعايير اللغوية والتكنولوجية التي خرج بها المشاركون في اجتماع خبراء المحللات الصرفية العربية (إبريل 2009) بدمشق (حماده والبواض 2009).

ويقوم البرنامج بإعطاء تحليل صرفي للكلمات العربية وذلك بتحديد:

- حالات التشكيل الممكنة للكلمة.
- الزوائد التي تلحق بها (السوابق واللواحق).
- نوع الكلمة (اسم أو فعل أو أداة).
- يشمل نوع الكلمة في حالة الاسم: الأسماء الجامدة، أسماء الأعلام، المصادر (الأصلي، الميمي، المرة والهيئة)، المشتقات (اسم الفاعل، اسم المفعول، اسم الآلة)

◦ وفي حالة الفعل: الماضي مبنياً للمعلوم وللمجهول، المضارع مبنياً للمعلوم وللمجهول (مرفوعاً ومنصوباً ومحظوماً ومؤكداً)، والأمر والأمر المؤكّد.

- الوزن (في حالة الأسماء والأفعال).
- الجذع.
- الجذر (في حالة الأسماء والأفعال).
- الحالة الإعرابية (في حالة الأسماء والأفعال)، وتشمل: الرفع والنصب والجر، الإفراد والتشييد والجمع، التذكير والتأنيث، التعريف والتنكير، التعدي واللزوم.

وتتلخص مراحل المعالجة الصرفية لكل كلمة خارج سياقها من النص في برنامج الخليل في ما يلي:

أولاً : مرحلة التقطيع: وتهدف إلى تحديد السوابق واللواحق المحتملة للكلمة، وذلك من خلال عملية تقطيع (segmentation) يجري خلالها استخدام قواعد معطيات خاصة بالسوابق واللواحق. ومن المعروف أن هذه العملية تؤدي غالباً إلى أكثر من تقطيع نظراً للطبيعة الإلصاقية للغة العربية.

ثانياً : تحليل أسماء الأعلام: ويتم خلالها معالجة كل تقطيع من التقطيعات المتحصل عليها في المرحلة السابقة للتحقق من إمكانية أن تكون الكلمة المدخلة اسمًا علمًا. وفي هذه المرحلة يستعين البرنامج بلائحة من أسماء الأعلام تضم 6000 اسم علم.

ثالثاً : تحليل الأدوات: وفيها يقارن البرنامج الجذوع المحتملة للكلمة المدخلة بلائحة الأدوات. وتضم هذه اللائحة حروف الجر والنصب والجزم والعطف والنداء والاستثناء وأدوات الشرط والاستفهام والتواصخ والضيائير المنفصلة والظروف وأسماء الإشارة وأسماء الموصولة.

رابعاً: معالجة الأسماء: وفيها يتحقق البرنامج من إمكانية أن تكون الكلمة المدخلة اسمًا. وتم هذه المرحلة عبر الخطوات التالية:

- تحديد الأوزان المحتملة للكلمة، وذلك بمقارنتها بالأوزان الموجودة في قاعدة البيانات الخاصة بالأوزان الاسمية.

- استخلاص الجذور المحتملة للكلمة والتحقق من وجودها في لائحة الجذور العربية.

- التتحقق من توافق الجذور والأوزان المحتملة، وذلك بالرجوع إلى قاعدة معطيات تضم الجذور العربية مرفقة بأوزان الأسماء المشتقة منها.

خامساً : معالجة الأفعال: وهي مرحلة مماثلة للمرحلة السابقة حيث تجري عبر الخطوات التالية:

- تحديد الأوزان المحتملة للكلمة، وذلك باستخدام قاعدة البيانات الخاصة بالأوزان الفعلية.

- استخلاص الجذور المحتملة للكلمة والتحقق من صحتها.
- التتحقق من توافق الأوزان والجذور المستخلصة.

هذه بإيجاز أهم مراحل المعالجة في النسخة الرسمية من برنامج الخليل للتحليل الصرفي وهي النسخة المتناثرة مجاناً على موقع "sourceforge"⁽³⁾. غير أن عملية دمج المحلل الصرفي في إطار برنامج التشكيل الآلي تطلبتنا إجراء تعديلات تمثلت أساساً في إضافة قاعدة معطيات على شكل معجم يتضمن الكلمات العربية الأكثر تكراراً في المدونات العربية المتناثرة، كما قمنا بتعديل آلية إخراج نتائج التحليل الصرفي للاحتفاظ بالتشكيلات الممكنة للكلمة دون غيرها من المعلومات الأخرى، كالنوع والوزن والجذر.. إلخ. وفي ما يلي نستعرض تفاصيل هذا المعجم.

معجم الكلمات الأكثر تكراراً:

تم إنشاء هذا المعجم بغية تسريع عملية المعالجة الصرفية من ناحية، ومن أجل ضبط نوعية تشكيل الكلمات التي يكثر استخدامها في مختلف النصوص من ناحية أخرى. ولهذا الغرض فقد اتبعنا الخطوات التالية:

- تجميع مدونة عربية واسعة ومتعددة تضم أزيد من 250 مليون كلمة، وذلك انطلاقاً من المدونات العربية المتناثرة على الشبكة (الإنترنت). وقد بلغ عدد المدونات التي استخرجنا منها معجم الكلمات الأكثر تكراراً ثماني مدونات، هي:

(3) www.sourceforge.net/projects/alkhalil

مدوّنة الباحث أحمد عبد العالى⁽⁴⁾، وهي مدونة غير مشكولة يبلغ حجمها 147 مليون كلمة (Abdelali et al, 2005). وت تكون هذه المدونة من نصوص صحفية متنوعة تم جمعها من 28 موقعًا صحفياً بشكل يغطي أغلب البلدان العربية.

مدوّنة "تشكيلة"⁽⁵⁾، وهي مدونة مشكولة تزيد على 60 مليون كلمة، وت تكون من نصوص تراثية قام بجمعها الباحث طه زروقي انطلاقاً من الكتب الموجودة بالمكتبة الشاملة⁽⁶⁾.

المدونة العربية مفتوحة المصدر⁽⁷⁾ (OSAC: Open Source Arabic Corpora)، وهي مدونة غير مشكولة جمعها الباحثان معتز خالد سعد ووسام عاشور (Saad and Ashour, 2010) بغية استخدامها في تطبيقات التنقيب في النصوص (Text Mining). وتحتوي هذه المدونة على قرابة 20 مليون كلمة جُمعت من موقع إخبارية عربية محلية وعالمية متنوعة.

مدوّنة الباحثة لطيفة السليطي⁽⁸⁾، وهي مدونة غير مشكولة تضمّ نصوصاً معاصرة وتشمل مجالات مختلفة (سياسية، اقتصادية، دينية، رياضية.. إلخ). وتضم هذه المدونة أزيد من نصف مليون كلمة.

مدوّنة "Khaleej-2004"⁽⁹⁾ وهي مدونة غير مشكولة جمعها الباحث مراد عباس انطلاقاً من موقع اليومية البحرينية أخبار الخليج، وذلك لأجل استخدامها في تطبيقات التصنيف الآلي للوثائق. وتضم هذه المدونة قرابة 2,8 مليون كلمة.

(4) <http://aracorpus.e3rab.com/argistestsrv.nmsu.edu/AraCorpus/Data/>

(5) <http://sourceforge.net/projects/tashkeela/>

(6) www.shamela.ws

(7) <https://sites.google.com/site/motazsite/Home/osac>

(8) <http://www.comp.leeds.ac.uk/eric/latifa/research.htm>

(9) <http://sourceforge.net/projects/arabiccorpus/files/>

مدونة قرارات الجمعية العامة للأمم المتحدة⁽¹⁰⁾ (UN Parallel Corpora)، وهي مدونة متوازية تضم 2100 قرار من قرارات الجمعية العامة للأمم المتحدة مكتوبة باللغات الست الرسمية للأمم المتحدة. وقد اقتصرنا بطبيعة الحال على محتواها العربي المكون من 2,5 مليون كلمة، علماً أن هذه المدونة قد جُمعت في الأصل لغرض استخدامها في تطبيقات الترجمة الآلية الإحصائية (Rafalovitch and Dale 2009).

مدونة RDI وهي مدونة مشكولة تتكون أساساً من كتب تراثية، بالإضافة إلى نسبة قليلة من الكتابات المعاصرة. وقد جُمعت هذه المدونة لغرض الاستعمال في مجال التشكيل الآلي. وتحتوي على 20 مليون كلمة.

مدونة نيملار⁽¹²⁾ (NEMLAR Arabic Written Corpus)، وهي مدونة مشكولة ومعنونة صرفيّاً ونحوياً، شارك في إعدادها باحثون من مراكز بحثية مختلفة (Attiya et al, 2005) وذلك في إطار مشروع «NEMLAR». نشير إلى أنها مدونة غير مجانية تقوم بتسويقها الجمعية الأوروبية للموارد اللغوية (ELRA)، وقد اقتصرنا في هذا البحث على محتواها المشكول الذي يضمّ نصوصاً من مجالات مختلفة تحتوي على قرابة نصف مليون كلمة.

- بعد مرحلة التجميع قمنا بتهيئة المدونة، وذلك بتوحيد ترميزها وطريقة تخزينها، حيث قمنا بحفظها في ملفات نصية (.txt). بترميز Cp1256. كما عملنا على تنقيتها من الرموز والكلمات المكتوبة بحروف غير عربية.

- بعد ذلك قمنا بحساب تكرار الكلمات كل مدونة على حدة وترتيبها ترتيباً تناظرلياً. وقد عمدنا في هذه الخطوة إلى عزل كل مدونة على حدة، نظراً لتفاوت أحجام هذه المدونات من جهة، ولكون الكلمات الشائعة تختلف حسب طبيعة

(10) <http://www.uncorpora.org/>

(11) <http://www.rdi-eg.com/RDI/TrainingData/>

(12) http://catalog.elra.info/product_info.php?products_id=873

المدونة من جهة أخرى. ففي حين تكُثر كلمات من قبيل "حدثنا" "قال" و"صلى" في المدونات التراثية، فإننا لا نجد لها أثراً في مدونة قرارات الجمعية العامة للأمم المتحدة على سبيل المثال.

- تمثل الخطوة الموالية في استخلاص اللائحة الأولية لمداخل معجم الكلمات الأكثر تكراراً، وذلك بعد دمج اللوائح المستخلصة من كل مدونة وحذف علامات التشكيل التي تظهر في بعضها، وخصوصاً لوائح المدونات التراثية، ثم حذف المكرر من هذه المداخل. وقد بلغت هذه اللائحة 16200 كلمة.

- بعد ذلك قمنا بتحديد المقابلات المشكولة لمداخل المعجم وقد طلبت هذه الخطوة تقسيم المداخل إلى مجموعتين:

- مجموعة الكلمات التي وجدنا لها مقابلات مشكولة في نصوص المدونات المشكولة (تشكيلة: RDI و NEMLAR). وفي هذه الحالة أرفقناها بالمقابلات المشكولة التي عثنا عليها.

- المجموعة الثانية وكانت في حدود 1200 كلمة، قمنا بتحليلها صرفاًً باستخدام محلل باكتور الصرفي لكونها تحتوي كثيراً من الألفاظ الدخيلة شائعة الاستخدام في النصوص الحديثة وكذلك أسماء الأعلام الأجنبية. ثم أتبعنا ذلك بتدقيق يدوي للكلمات التي حلّلها والكلمات التي لم يتمكن من تحليلها.

بعد استكمال هذه الخطوات تحصل لدينا معجم من 16200 كلمة من الكلمات الأكثر تكراراً في المدونات العربية المتاحة. وقد قمنا بدمج هذا المعجم في إطار عملية التحليل الصرفي لبرنامج الخليل بحيث يتعرف مباشرة على التشكيلات الممكنة للكلمات إن كانت مما يتضمنه معجم الكلمات الأكثر تكراراً مما يُسرّع من عملية المعالجة.

2. مرحلة المعالجة الإحصائية

بعد أن يقوم البرنامج بمعالجة صرفية لكلمات النص المدخل، وبعد أن يحتفظ بالتشكيلات المحتملة لكلّ كلمة، فإنه يمر إلى المرحلة الثانية من مراحل التشكيل. وتمثل هذه المرحلة، كما أسلفنا، في معالجة إحصائية مبنية على نموذج ماركوفي يتم عبره اختيار التشكيل الأرجح لكلّ كلمة باستخدام خوارزمية فيتربي (Neuhoff 1975). وفي ما يلي نعرض لهذا النموذج شيء من التفصيل:

لتكن المجموعة $O = \{o_1, \dots, o_M\}$ مجموعة متهية من "المشاهدات"

(Observations)

ولتكن $S = \{s_1, \dots, s_N\}$ مجموعة متهية من "الحالات الخفية"

(Hidden States)

تعريف:

نُعرّف نموذج ماركوف الخفي من الرتبة الأولى بأنه كلّ زوج من السلاسل $(X_t, Y_t)_{t \geq 1}$ بحيث تكون:

• عبارة عن سلسلة ماركوف متجانسة (homogeneous) تأخذ قيمها في مجموعة الحالات الخفية S بحيث يكون الاحتمال: (Markov chain

$$\Pr(X_{t+1} = s_j / X_t = s_i, \dots, X_1 = s_h) = \Pr(X_{t+1} = s_j / X_t = s_i) = a_{ij}$$

وتعبر a_{ij} عن احتمال المرور من الحالة s_i إلى الحالة s_j .

• بينما تأخذ السلسلة $(Y_t)_{t \geq 1}$ قيمها في مجموعة المشاهدات O بحيث:

$$\begin{aligned} \Pr(Y_t = o_k / X_t = s_i, Y_{t-1} = o_{k_{t-1}}, X_{t-1} = s_{i_{t-1}}, \dots, Y_1 = o_{k_1}, X_1 = s_{i_1}) \\ = \Pr(Y_t = o_k / X_t = s_i) = b_i(k) \end{aligned}$$

وتعبر $(k)_i$ عن احتمال مشاهدة O_k مع العلم بتحقق S_i .

وبخصوص مشاهدات النموذج وحالاته الخفية في عملنا هذا فإننا نعرّفها كالتالي:

أ. مشاهدات النموذج (Observations)

نقترح في هذا العمل تصنيفاً جديداً للكلمات العربية. وتستند عملية التصنيف (classification) هذه على استخلاص ملامح (features) من الكلمة المراد معاجتها تمثّل في:

- طول الكلمة.
- السوابق واللواحق.
- مواضع أحرف الألف والواو والياء إن وجدت في الكلمة.
- حرف أو حرفان من الكلمة يتم اختيارهما تبعاً لطوفها.

ولتوسيع هذا التصنيف نضرب مثالاً بكلمة من قبيل "التشكيل" فلدى تصنيف هذه الكلمة يستنتج البرنامج صنفاً على شكل:

`alS#;0;-1;5;+ش;+ي`

وتؤول هذا الصنف راجع إلى كونه اتحاد الملامح المستخلصة التالية:

- "F7" : وهو الملمح المشترك بين جميع الكلمات المكونة من سبعة أحرف.
- "alS" : ويدلّ على أن الكلمة مبدوءة بأداة التعريف "ال"، وأن أدلة التعريف متبوءة بحرف شمسي (S) وهو التاء في هذه الحالة. إن تمييزنا للملمح السوابق التي تتضمن أدلة التعريف إلى نوعين شمسي (alS) وقمرى (alL) تبعاً للحرف الذي يلي السوابق، راجع بطبيعة الحال إلى تأثير ذلك في تشكيل هذا الحرف وحرف اللام الذي يسبقه.

- "#": وتدل على أن الكلمة ليست لها لاحقة (suffix).
- "0": وتعني أن حرف الألف موجود في كلمة "التشكيل" في الموضع رقم 0 ؛ أي في الحرف الأول من الكلمة. وفي هذا السياق فإنه إذا تعددت مواضع حرف من الحروف (أو ي)، فإن الملمح الخاص بها يكون على شكل متوجهة (vector) تتركب من هذه الموضع. فالملمح الخاص بحرف الواو في كلمة "الموعدون" هو (3,5,7) ويعني أن الواو هو الحرف الرابع والسادس والثامن من هذه الكلمة.
- "-1": ويعني أن حرف الواو لا يوجد في مثالنا أي كلمة "التشكيل".
- "5": مما يعني أن حرف الياء هو الحرف السادس من هذه الكلمة.
- "ش+ي": وهو الحرفان اللذان يقعان في الموضعين الرابع والسادس على التوالي، وهو الموضعان اللذان تختار منها الحروف المميزة للكلمات المكونة من سبعة أحرف.

وتحت الصنف الوارد في المثال السابق تقع الكثير من كلمات اللغة العربية التي تتشابه في كيفية تشكيلها، مثل: "التشبيه" و"التشخص" و"الشرع" و"التشغيل"... إلخ. إذ يشمل، بعبارة أخرى، كل المصادر التي على وزن "التفعيل" إذا كان الحرف الأول من جذرها حرف الشين. وتُعد هذه الخاصية أهم إيجابيات هذا التصنيف؛ إذ تكون "بارامترات" النموذج أقل بكثير من مثيلتها في الأعمال السابقة (Gal, 2002) (Elshafei et al, 2006)، (بawah وأخرون 2011) التي تمثل الكلمات العربية فيها مشاهدات النماذج المارкова.

بـ. الحالات الخفية (Hidden states)

ت تكون الحالات الخفية في النموذج الماركوفي الذي نقترحه من تسلسلاً علامات التشكيل الممكنة في اللغة العربية. وهي عبارة عن متوجهات (vectors)

تختلف أبعادها تبعاً لأطوال المشاهدات. ولتوسيع هذا التعريف أكثر نضرب أمثلة من هذه الحالات الخفية:

٠ (، ، ، ،) : وهو التسلسل المكون من فتحة فسكون فكسرة فضمة، ويعُرَّف حالة خفية تكمن خلف الكثير من المشاهدات عند تصنيف كلمات من أربعة أحرف مثل الكلمات: "يعرض" ، "تعرف" ، "نصير" .. إلخ، بما يسمح بتشكيلها كالتالي: "يَعْرِضُ" ، "تَعْرِفُ" ، "نَصِيرٌ" .. إلخ

٠ (، ، ، ،) : وهو متوجه من خمس مركبات تمثل تسلسل الضمة فالفتحة فالشدة مع الكسر فالفتحة فتنوين الضم. ويقابل هذا التسلسل عدداً من المشاهدات، ما يسمح بتشكيل كلمات مثل "مُعَلَّمَة" ، "مُقدَّمَة" .. إلخ.

٠ (، # ،) : ويدلّ الرمز # على غياب علامة التشكيل على الحرف المقابل. ويسمح هذا التسلسل الثلاثي بتشكيل كلمات مثل: "جاء" ، "نَالَ" ، "فَازَ" .. إلخ.

ولا يخفى أن الطابع التجريدي للحالات الخفية في نموذجنا هذا يستدعي عدداً أقل من "البارامترات" مقارنة بالنموذج الماركوفية السابقة (Gal, 2002)، (Elshafei et al, 2006) التي تتخذ من الكلمات المشكولة حالات خفية.

وفي ما يلي نُبيّن كيف استخدمنا هذا النموذج لغرض التشكيل الآلي للنصوص العربية.

لنفرض على سبيل المثال جملة عربية على الشكل W_1, \dots, W_n بعد القيام بتصنيف هذه الكلمات بالطريقة الواردة أعلاه في النموذج فإننا نحصل على سلسلة من المشاهدات w_1, \dots, w_n ، ولتكن المجموعة $C = \{c_1, \dots, c_N\}$ مجموعة سلاسل التشكيلات الممكنة في اللغة العربية. انطلاقاً من هذه الفرضيات، فإن مسألة تحديد التشكيلات الصحيحة في هذه المرحلة من المعالجة

تتمثل في إيجاد مجموعة سلاسل التشكيلات (c_1^*, \dots, c_n^*) التي تحقق المعادلة التالية:

$$(c_1^*, \dots, c_n^*) = \arg \max_{c_1 \dots c_n \in C} \Pr(c_1 \dots c_n / w_1 \dots w_n)$$

وبما أن:

$$\Pr(c_1 \dots c_n / w_1 \dots w_n) = \frac{\Pr(w_1 \dots w_n / c_1 \dots c_n) \Pr(c_1 \dots c_n)}{\Pr(w_1 \dots w_n)}$$

إن التسلسل (c_1^*, \dots, c_n^*) يتحقق:

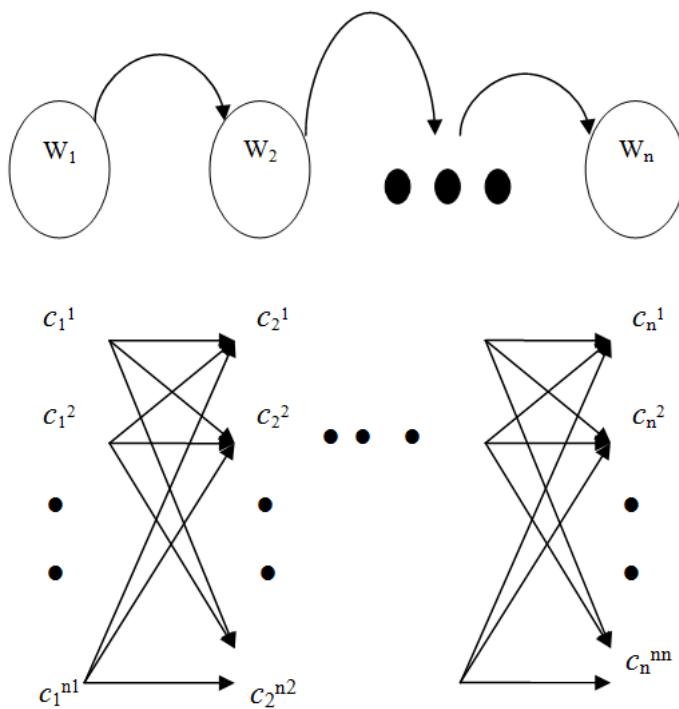
$$(c_1^*, \dots, c_n^*) = \arg \max_{c_1 \dots c_n \in C} \Pr(w_1 \dots w_n / c_1 \dots c_n) \Pr(c_1 \dots c_n)$$

وهي المعادلة التي يمكن أن تكتب على الشكل التالي:

$$(c_1^*, \dots, c_n^*) = \arg \max_{\substack{c_i^{j_i} \in C_i \\ 1 \leq i \leq n}} \Pr(w_1 \dots w_n / c_1^{j_1} \dots c_n^{j_n}) \Pr(c_1^{j_1} \dots c_n^{j_n})$$

حيث ترمز المجموعة $\{c_i^1, \dots, c_i^{n_i}\}$ إلى مجموعة التشكيلات المحتملة الناتجة عن تحليل الكلمة w_i ذات الصنف $.w_i$.

ونقوم بحل هذه المعادلة عن طريق البحث عن المسار الأكثر رجحانًا في شبكة الحلول الناتجة عن التحليل الصري للكلمات خارج سياقها كما هو مبين في الشكل التالي:



الشكل 2. شبكة سلاسل التشكيلات المحتملة الناتجة عن تحليل الجملة

$$W_1, \dots, W_n$$

وتعتبر خوارزمية فيتري وسيلة فعالة للبحث عن المسار الأكثر رجحانًا في هذه الشبكة. وتعتمد الخوارزمية على قيم دالتين ϕ و ψ نعرفهما كما يلي:

$$\phi(t, c_t^k) = \max_{\substack{r_i^{ji} \in \mathfrak{R}_i \\ 1 \leq i \leq t-1}} \left[\Pr(w_1, \dots, w_{t-1}, w_t / c_1^{k_1}, \dots, c_{t-1}^{k_{t-1}}, c_t^k) \times \Pr(c_1^{k_1}, \dots, c_{t-1}^{k_{t-1}}, c_t^k) \right]$$

حيث تُعبر $\phi(t, c_t^k)$ عن قيمة احتمال المسار الجزئي الأرجح الذي يمر من التسلسل c_t^k (c_t^k ينتمي لمجموعة التشكيلات المحتملة للكلمة W_t).

ويمكن أن نكتب المعادلة السابقة على الشكل التالي:

$$\begin{aligned}\phi(t, c_t^k) &= \max_{\substack{r_i^{j_i} \in \mathfrak{R}_i \\ 1 \leq i \leq t-1}} \prod_{i=1}^{t-1} [\Pr(w_i / c_i^{j_i}) \times \Pr(c_i^{j_i} / c_{i-1}^{j_{i-1}}) \times \Pr(w_t / c_t^k) \times \Pr(c_t^k / c_{t-1}^{j_{t-1}})] \\ &= \left(\max_{c_{t-1}^j \in C_{t-1}} \phi(t-1, c_{t-1}^j) \times \Pr(c_t^k / c_{t-1}^{j_{t-1}}) \right) \Pr(w_t / c_t^k) \quad (*)\end{aligned}$$

وتسمح لنا المعادلة الأخيرة بحساب قيم الدالة ϕ بصيغة تراجيعية.

ومن أجل التعرف على المسار الأرجح فإننا نعرف الدالة ψ التي تسمح في كل لحظة t بتخزين التشكيل الذي يتحقق أكبر قيمة للمعادلة السابقة (*).

نُعرف ψ كما يلي:

$$\psi(t, c_t^k) = \arg \max_{c_{t-1}^j \in C_{t-1}} \phi(t-1, c_{t-1}^j) \Pr(c_t^k / c_{t-1}^j) \quad (**)$$

مع ملاحظة أن: $\psi(t, c_t^k) \in C_{t-1}$

هاتان المعادلتان (*) و (**) تسمحان لنا بإيجاد المسار الأرجح عن طريق خوارزمية فيتري التراجيعية التالية:

- المرحلة الأولى الابتداء (Initialization) :

for $1 \leq k \leq n_1$

حساب $\phi(t, c_1^k)$: احتمال أن تبدأ الجملة بالكلمة w_1 ويكون تشكيلها c_1^k .

- المرحلة الثانية الحساب التراجعي (Recursion) :

for $2 \leq t \leq n$ and $1 \leq k \leq n_t$

$$\begin{aligned}\phi(t, c_t^k) &= \left(\max_{c_{t-1}^j \in C_{t-1}} \phi(t-1, c_{t-1}^j) \times \Pr(c_t^k / c_{t-1}^j) \right) \Pr(w_t / c_t^k) \\ \psi(t, c_t^k) &= \arg \max_{c_{t-1}^j \in C_{t-1}} \phi(t-1, c_{t-1}^j) \Pr(c_t^k / c_{t-1}^j)\end{aligned}$$

• المرحلة الثالثة الحالة الأخيرة (Final state):

$$\psi(n+1) = \arg \max_{c_n^j \in C_n} \phi(n, c_n^j)$$

• المرحلة الرابعة استنتاج المسار الأرجح (Deducing the best path):

$$c_n^* = \psi(n+1) \quad \circ$$

$$\text{For } t = n-1 : 1 \quad c_t^* = \psi(t, c_{t+1}^*) \quad \circ$$

4. تدريب البرنامج

تحتاج "بارامترات" النموذج الإحصائي الوارد أعلاه، والمتمثلة في قيم المصفوفتين (a_{ij}) و (b_{it})، إلى عملية تقييم - خلال مرحلة تدريب البرنامج - بالاعتماد على ذخيرة لغوية كافية.

ولو اعتبرنا ذخيرة لغوية $C = \{Ph_1, \dots, Ph_K\}$ مكونة من عدد كبير من الجمل، فإن عملية تدريب النظام تتلخص في تقييم "البارامترات" بطريقة (Manning and Schütze, 1999) (Maximum Likelihood) على الشكل التالي:

$$a_{ij} = \frac{\sum_{n=1}^K (\text{Ph}_n \text{ في الجملة } c_i \text{ التشكيل } c_j \text{ إلى التشكيل } c_i \text{ في الجملة } c_j \text{ عدد مرات الانتقال من التشكيل } c_i \text{ إلى التشكيل } c_j)}{\sum_{n=1}^K (\text{Ph}_n \text{ في الجملة } c_i \text{ التشكيل } c_i \text{ عدد تكرارات التشكيل } c_i)}$$

$$b_{it} = \frac{(\text{عدد مرات أخذ الصنف } w_t \text{ للتشكيل } c_i)}{(\text{عدد تكرارات التشكيل } c_i)}$$

وإنجاز عملية التدريب هذه قمنا بعزل مدوّنة مشكولة بلغ حجمها 1,24 مليون كلمة مشكولة موزعة بين نصوص مدوّنة "NEMLAR"، التي اخترنا 90% منها (460 ألف كلمة) واستبعدينا 10% للاختبار، ومدوّنة "تشكيلة" التي اخترنا منها (780 ألف كلمة). ويرجع السبب في عدم استخدامنا لملفين

النصوص التراثية المشكولة التي لدينا في التدريب إلى محاولة إيجاد توازن في مدوّنة التدريب بين المحتوى المعاصر والمحتوى التراثي. وبالتالي توافي مع النصوص المعاصرة لمدونة "NEMLAR" قمنا بعزل 30 كتاباً تراثياً من مدوّنة "تشكيلة" واختربنا بشكل عشوائي من كل كتاب قرابة 10% من محتواه. كما أجرينا على مدونة التدريب سلسلة من العمليات تمثلت في تقسيم محتوياتها إلى جمل ثم تصنيف جميع كلماتها واستخلاص تسلسلات التشكيل فيها من أجل حساب "البارامترات" باستخدام طريقة الإمكان الأعظم (Maximum likelihood) الواردة أعلاه.

5. اختبار البرنامج

من أجل تقييم كفاءة برنامج التشكيل الآلي فقد اعتبرنا المقاييس المعروفة: نسبة الكلمات الخاطئة (WER : Word Error Rate) ونسبة الخطأ على مستوى الحروف (DER : Diacritic Error Rate). ويتفق هذان المقاييس في الواقع إلى أربعة مقاييس بالنظر إلى اعتبار تشكيل الحرف الأخير في الكلمة من عدمه. فبذلك كانت المقاييس التي استخدمنا في عملية التقييم كالتالي:

• WER1 : ويقيس نسبة الكلمات الخاطئة مع اعتبار تشكيل الحرف الأخير.

• WER2 : ويقيس نسبة الكلمات الخاطئة مع تجاهل تشكيل الحرف الأخير.

• DER1 : ويقيس نسبة الخطأ في التشكيل على مستوى الحروف بها فيها الحرف الأخير.

• DER2 : ويقيس نسبة الخطأ في التشكيل على مستوى الحروف باستثناء الحرف الأخير.

وتمثل تجربة التقييم في اختيار خمس عينات من مدونة الاختبار، تكون كل عينة من 100 جملة، تم سحب كل جملة منها بشكل عشوائي، وذلك بغية

الحصول على القيمة التقريرية لنسب الخطأ الأربع على مستوى كامل مدونة الاختبار.

ويلخص الجدول التالي النتائج التي حصلنا عليها.

العينة	عدد الكلمات	WER1(%)	WER2(%)	DER1(%)	DER2(%)
T1	2601	27,73	13,42	8,09	4,02
T2	2380	27,57	14,79	8,51	4,75
T3	2525	28,04	13,75	8,23	4,19
T4	2863	31,34	16,49	9,68	5,24
T5	2356	27,81	13,97	8,34	4,28
المتوسط	2545	%28,5	%14,5	%8,6	%4,5

جدول 1. نتائج تقييم البرنامج على عينات عشوائية من مدونة الاختبار

تبين هذه النتائج أن نسبة خطأ التشكيل على مستوى الكلمات تقارب 29%， وأن تلك النسبة تهبط إلى النصف في حالة تجاهل الحرف الأخير لتصل 14,5%. ولنا على هذه النتيجة بعض الاستدراكات. فالمقياس "WER1" بالشكل الذي اعتمدناه في هذه التجربة هو مقياس متشدد، يفترض وجود نسخة من مدونة الاختبار مشكولة بشكل صحيح تام. فعند مقارنة الكلمة الأصلية بنتيجة معالجة البرنامج لنسخة غير مشكولة منها لا يعطي المقياس نتيجة إيجابية إلا في حالة التطابق التام بين الكلمتين. والواقع أن المدونات المتاحة لا تسلّم، من

جهة، من الأخطاء الكتابية كالكلمات المتلاصقة والكلمات المفروقة والتساهم في المهزات، وغير ذلك من الأخطاء الإملائية التي تربك المحلل الصري المدمج، كما أنها لا تخلو، من جهة أخرى، من الكلمات المشكولة جزئياً التي قد يقوم البرنامج بتشكيلها تشكيلاً تاماً صحيحاً يبقى رغم ذلك خاطئاً بحسب المقياس . "WER1"

ويدلّ هبوط نسبة الخطأ إلى النصف (14,5%) في المقياس WER2 على أن نسبة مهمة من أخطاء التشكيل هي أخطاء إعرابية كعدم الاتفاق، أحياناً، بين حركة المتابعين (نعتاً و توكيداً و عطفاً و بدلاً) أو عدم جرّ المضاف إليه أو الخطأ في حركة بعض معمولات النواسخ و حروف الجرّ. ذلك بأنّ معيار الترجيح الذي يستند عليه البرنامج هو معيار إحصائيّ يصيب ويخطئ في هذه الحالات الإعرابية تبعاً لقيم الاحتمالات التي استخلصها من مدونة التدريب.

وتصل نسبة الخطأ في تشكيل جميع حروف النص "DER1" إلى 8,6%، وتهبط هذه النسبة أيضاً إلى حدود النصف 4,5% (DER2) لتجاوز بذلك نسبة الدقة على مستوى الحروف 95% في حال تجاهل الحرف الأخير.

6. الخاتمة

قدمنا في هذه البحث برنامجاً جديداً أطلقنا عليه اسم "برنامج الخليل للتشكيل الآلي"، يعتمد هذا البرنامج على مقاربة هجينه لتشكيل اللغة العربية تجمع بين التحليل الصري ونمادج ماركوف الخفية. وتحتفل هذه المقاربة عن غيرها من المقارب المهيمنة على المستويين اللغوي والإحصائيّ.

فعل المستوى اللغوي تمّ توظيف أحد أهم المحللات الصرفية العربية مفتوحة المصدر، حيث يقوم البرنامج باستخدام نتائج تحليل برنامج الخليل للتحليل الصري واستخلاص التشكيلات الممكنة التي يقترحها للكلمات خارج سياقها من النص. وقد تطلب دمج المحلل الصري في إطار التشكيل الآلي إضافة

قاعدة معطيات معجمية تتضمن الكلمات الأكثر تكراراً، وذلك بهدف ضبط تشكيلها وتفادي التعميم (over generation) الذي يظهر في كثرة المخرجات المقترحة من قبل المحلل الصRFي. وقد تم استقاء قاعدة المعطيات المعجمية من مدوّنة تضم أزيد من 250 مليون كلمة جمعت من ثمانٍ مدونات عربية متاحة على الشابكة (الإنترنت).

وعلى المستوى الإحصائي، قمنا باقتراح نموذج ماركوفي ينطلق من تصنيف جديد للكلمات العربية، بحيث تمثل هذه الأصناف مشاهدات النموذج، بينما تمثل تسلسلاً علامات التشكيل الحالات الخفية له. كما رأينا في تدريب البرنامج استخدام مدوّنة متنوعة تجمع بين النصوص الحديثة والتراثية بلغ حجمها أزيد من مليون ومائتي ألف كلمة.

وتعتبر النتائج التي حصلنا عليها جد مشجعة، إذ وصلت نسبة الدقة على مستوى كامل الحروف إلى 91,4%， وإلى أزيد من 95% بتجاهل الحرف الأخير. وهي نتائج نطمح إلى تحسينها بالعمل على مستويين متكمالين: مستوى المدونات، وذلك بتوسيع محتوى المدونة من اللغة المعاصرة وضبط تشكيل أكبر قدر ممكن منها. ومستوى التحليل اللغوي، وذلك بدمج المقاربة في إطار أشمل يتعامل مع باقي المعلومات التي يتيحها محلل الخليل الصRFي سعياً لإيجاد محلل نحوي ذي كفاءة مقبولة، مما سينعكس إيجاباً على نتائج المشكل الآلي. كما نعتزم فتح مصدر هذا البرنامج لتعيم استخدامه وإتاحة المجال أمام اجتهادات المطورين في تحسينه ودمجه في تطبيقات أخرى.

مراجع باللغة العربية

- 1 - سلوى حماده، مروان والباب: "معايير وضوابط تقييم المحللات الصرفية"، تقرير من ضمن توصيات المشاركين في اجتماع خبراء المحللات الحاسوبية الصرفية للغة العربية، 26 - 28 إبريل 2009، دمشق.
- 2 - عثمان بن سعيد الداني أبو عمرو، المحكم في نقط المصاحف، تحقيق: د. عزة حسن، دار الفكر - دمشق، الطبعة الثانية، 1407 هـ.
- 3 - عز الدين مزروعي، عبد الوافي مزيان، عبد الحق لخواجه، عبد الرحيم بودلال، محمد ولد عبد الله ولد بياه: "برنامج الخليل الصرفي"، المؤتمر الدولي لعلوم وهندسة الحاسوب باللغة العربية، ICCA 2011، 31 مايو - 2 يونيو 2011، الرياض.
- 4 - محمد عطية، "التشرع البنائي لمشكّل آلي عربي لتوظيفه في نظام تخليل آلي للصوت المنطوق من النص العربي المكتوب"، ندوة تقنية المعلومات والعلوم الشرعية والعربية، جامعة الإمام محمد بن سعود الإسلامية، الرياض، 2007.
- 5 - محمد ولد عبد الله ولد بياه، عبد الوافي مزيان، عز الدين مزروعي، عبد الحق لخواجة، عبد الرحيم بودلال: "مقاربة صرفية إحصائية للتشكيل الآلي"، اجتماع خبراء التدقيق الإملائي والتشكيل والتحليل النحوي الآلي، 18 - 20 إبريل 2011، دمشق.
- 6 - منصور الغامدي؛ محمد خورشيد؛ مصطفى الشافعي؛ فايز الحرقان؛ أبو أوس الشمسان؛ محمد الكنهل؛ سعد القحطاني؛ سيد زيشان مظفر؛ ياسر التويم؛ عدنان يوسف؛ حسني المحتسب؛ "نظام حاسوبي لتشكيل النص العربي"، التقرير الفني النهائي، مدينة الملك عبد العزيز للعلوم والتقنية، 2006.

مراجع باللغات الأجنبية

- 1 – Abbas M., Smaili K.. Comparison of Topic Identification Methods for Arabic Language, International conference RANLP05 : Recent Advances in Natural Language Processing , 21-23 september 2005, Borovets, Bulgaria.
- 2 – Abdelali A., Cowie J., Soliman H.: “Building A Modern Standard Arabic Corpus”, Workshop on Computational Modeling of Lexical Acquisition, The Split Meeting, Crotia, 25-28 July 2005.
- 3 – Attia. M., Yaseen., M., and Choukri., K. :”Specifications of the Arabic Written Corpus produced within the NEMLAR project”, www.NEMLAR.org. 2005.
- 4 – Berger, A., Della Pietra, S., Della Pietra, V., “A maximum entropy approach to natural language processing”. Computational Linguistics 22 (1), 39–71. 1996.
- 5 – Buckwalter T.. Arabic Morphological Analyzer version 2.0. LDC2004L02. 2004.
- 6 – Debili F., ACHOUR H.: “Voyellation automatique de l’arabe” in Proceedings of the workshop on Computation approaches to Semitic languages, COLING-ACL ‘98., Montréal, 1998.
- 7 – Deltour A., « Méthodes statistiques pour la voyellisation des texts arabes », Projet de fin d’études, ENSIMAG-Universität Karlsruhe, 2003.
- 8 – El-Sadany, T., Hashish, M.,: “Semi-automatic vowelization of Arabic verbs”. In: 10th NC Conference, Jeddah, Saudi Arabia. 1988.
- 9 – Elshafei, M., Almuhtasib H. and Alghamdi M, :”Machine Generation of Arabic Diacritical Marks”, The 2006 World Congress in Computer Science Computer Engineering, and Applied Computing . Las Vegas, USA. 2006.
- 10 – Emam O. and Fischer V.: “Hierarchical Approach for the Statistical Vowelization of Arabic Text. Technical report, IBM Corporation Intellectual Property Law, Austin, TX, US 2005.

- 11 – Gal Y.: “An HMM Approach to Vowel Restoration in Arabic and Hebrew”, In ACL-02 Workshop on Computational Approaches to Semitic Languages 2002.
- 12 – Habash N. and Rambow O.: “Arabic Diacritization through Full Morphological Tagging”. In Proceedings of NAACL/HLT 2007. Companion Volume, Short Papers, Rochester, New York, April. 2007.
- 13 – Maamouri M., Bies A., and Kulick S.: “Diacritization: A Challenge to Arabic Treebank Annotation and Parsing”. In Proceedings of the British Computer Society Arabic NLP/MT Conference, 2006.
- 14 – Manning C. and Schütze H., Foundations of Statistical Natural Language Processing, MIT Press. Cambridge, MA: May 1999.
- 15 – Messaoudi A., Lamel L., and Gauvain J. : “The LIMSI RT04 BN Arabic system”. In Proc. DARPA RT04, Palisades NY, November 2004.
- 16 – Nelken R. and. Shieber S. M.: “Arabic Diacritization Using Weighted Finite-State Transducers. In Proceedings of the ACL 2005 Workshop On Computational Approaches To Semitic Languages”, Ann Arbor, Michigan, USA. 2005.
- 17 – Neuhoff D.L., The Viterbi Algorithm as an Aid in Text Recognition, IEEE Transaction on Information Theory, pp 222-226, March 1975.
- 18 – Rafalovitch A, Dale R.: “United Nations General Assembly Resolutions: A Six-Language Parallel Corpus”. In *Proceedings of the MT Summit XII*, pages 292-299, Ottawa, Canada, August.
- 19 – Saad M. K., Ashour W.: “OSAC: Open Source Arabic Corpora”, 6th ArchEng International Symposiums, EEECS’10 the 6th International Symposium on Electrical and Electronics Engineering and Computer Science, pp. 118-123, European University of Lefke, Cyprus, 2010.
- 20 – Schlippe T., Nguyen T. and Vogel S.: “Diacritization as a Machine Translation Problem and as a Sequence Labeling Problem,” 8th AMTA conference, Hawaii, 21-25 October 2008.

- 21 – Vergyri D. and Kirchhoff K., ‘Automatic diacritization of Arabic for Acoustic Modeling in Speech Recognition’. In Proceedings of the Workshop on Computational Approaches to Arabic Script-based Languages. COLING 2004, Geneva: 66-73.
- 22 - Zitouni I., Sorensen J. S. and Sarikaya R : ‘Maximum Entropy Based Restoration of Arabic Diacritics’. In Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics. Workshop on Computational Approaches to Semitic Languages, Sydney, Australia. July 2006.

