

An Application Oriented Arabic Morphological Analyzer

م. رياض سنبل د. ندى غنيم د. محمد سعيد دسوقي

يقوم نظام التحليل الصرفي بتحليل الكلمة العربية والمعالجة الصرفية لها. وهو يقبل في مدخله كلمةً مشكولةً كلياً أو جزئياً أو غير مشكولة، ويعطي في مخرجه جميع الوجوه الممكنة للسوابق واللواحق والجذور والجدوع والأوزان، ويرتّب النظام المنجز هذه الحلول حسب "جودتها" بعد أن يقدرّ درجة وثوقيّة معينة لكل حل.

يعتمد النظام على خوارزمية مبتكرة جعلته يصل إلى وثوقية تتجاوز 97.7% دون الاعتماد على مواد مخزنة، وبسرعة معالجة يمكن أن تصل إلى أكثر من خمسة أضعاف الأنظمة المشابهة (تصل إلى أكثر من ٥٠ ألف كلمة بالثانية في بعض الحالات). ويحقق النظام المنفّذ مرونة كبيرة تجعله قابلاً للاستخدام في مختلف المجالات ومختلف أنواع التطبيقات، مثل تطبيقات المعالجة اللغوية كالمعاجم وأنظمة تعليم اللغة العربية، وتطبيقات البحث عن المعلومات، ومحركات البحث، وتطبيقات التنقيب في النصوص والمعطيات، وتطبيقات الترجمة الآلية، وغيرها. وقد تم تحقيق ذلك من خلال التحكم بنقطة توازن النظام بين سرعة المعالجة ووثوقية النتائج وشموليّة الحلول (مدى الحصول على كل التحليلات الممكنة للكلمة)، وتفعيل أو عدم تفعيل معالجة الحالات الخاصة كالإعلال (بالقلب والحذف) والإدغام (الشدة).

The screenshot shows a window titled 'hiast0.txt.txt' with a menu bar containing 'ملف', 'تحرير', 'عرض', 'تنسيق', 'مساعدة'. Below the menu bar is a toolbar with various icons. The main content area displays the following text:

يهدف المعهد إلى إعداد فرق مؤهلة للبحث العلمي والتقني في جميع مجالات العلوم التطبيقية والتكنولوجيا.

نتيجة التحليل الصرفي

الكلمة: يهدف

الوزن	الجذر	السوابق	الجذع	اللواحق
يفعل	هـدـف	—	يهدف	—

الكلمة: المعهد

الوزن	الجذر	السوابق	الجذع	اللواحق
المفعّل	عـهـد	ال	معهد	—

يلخّص الجدول التالي وسائط التحكم المستخدمة في المحلل الصرفي وعددها ٩، والتي ينتج عنها ١٨ نسخة من نظام التحليل الصرفي. تضيف كل نسخة إمكانية معالجة جديدة على النسخة السابقة لها في الترقيم. فالنسخة Rn في الجدول تعني أن الوسائط من ١ .. n مفعّلة (تأخذ قيم true). أما الإشارة + فهي تعني أنّ النسخة تأخذ الحل الأفضل (ذو الاحتمال الأعلى) من الحلول الممكنة.

اسم النسخة	+	8	7	6	5	4	3	2	1	وصف النسخة
	Bst	Elm	Shd	Ela	Ebd	Rot	Pat	Frg	Stp	
R0										النسخة الأولى دون معالجة الحالات الخاصة
R0+										اختيار الحل الأفضل من بين الحلول الممكنة في النسخة R0
R1										إضافة اختيار الكلمات الخاصة stop words
R1+										اختيار الحل الأفضل من بين الحلول الممكنة في النسخة R1
R2										إضافة اختيار الكلمات الأجنبية والمعرّبة
R2+										اختيار الحل الأفضل من بين الحلول الممكنة في النسخة R2
R3										إضافة اختيار وجود الجذر في قائمة الجذور
R3+										اختيار الحل الأفضل من بين الحلول الممكنة في النسخة R3
R4										إضافة اختيار الوزن
R4+										اختيار الحل الأفضل من بين الحلول الممكنة في النسخة R4
R5										إضافة معالجة الإبدال
R5+										اختيار الحل الأفضل من بين الحلول الممكنة في النسخة R5
R6										إضافة معالجة الإعلال بالقلب
R6+										اختيار الحل الأفضل من بين الحلول الممكنة في النسخة R6
R7										إضافة معالجة الإدغام (الشدّة)
R7+										اختيار الحل الأفضل من بين الحلول الممكنة في النسخة R7
R8										إضافة معالجة الإعلال بال حذف
R8+										اختيار الحل الأفضل من بين الحلول الممكنة في النسخة R8

يجب أن لا ننسى أن هناك معايير أخرى مثل سرعة المعالجة، وعدم الاعتماد على الموارد اللغوية المخزنة، إضافةً إلى جودة الحلول الناتجة ووثوقية النتائج. وهذا ما دفعنا إلى اعتماد مبدأ وسائط التحكم التي نستطيع بها الموازنة بين هذه المعايير واختيار الأفضل منها حسب التطبيق.

يوفر النظام عدداً متنوعاً من الحالات، فمثلاً الحالات (R3 ... R6) تتوافق مع متطلبات الأنظمة العالية الوثوقية (كالمراحل الأخرى من المعالجة الحاسوبية، الشكل الآلي للنصوص،...). وقد بينت الاختبارات أن هذه الحالات تحقق وثوقيةً تجاريةً أو تفوق وثوقية الأنظمة الموجودة في هذا المجال، إضافةً لكونها تتمتع بسرعات معالجة تصل لما بين ضعف إلى خمسة أضعاف الأنظمة الأخرى (التي تمت دراستها). مع وجود احتمال كبير لوجود الحل الموافق للسياق كأول حل من الحلول الناتجة.

أما بالنسبة لمجال تطبيقات التتقيب في النصوص، فإن الحالات $Rn+$ تعتبر الحل الأكثر مثالية بين الحلول المطروحة، لأنها سرعات معالجة عالية تصل إلى أكثر من ٥٠ ألف كلمة بالثانية في بعض الحالات. ومن المميزات الهامة لهذه الحالات أن الحل الأول يتوافق مع السياق في أكثر من ٩٥% من الكلمات المختبرة.

أما بالنسبة للتطبيقات اللغوية (المعاجم الحاسوبية، أنظمة تعليم اللغة العربية،...) فهي تطبيقات تتطلب عادةً وثوقيةً "شبه كاملة"، وسرعة المعالجة ليست ذات أولوية فيها، لذلك يمكن أن نلجأ لمشاركة النظام مع قاعدة معطيات تحوي جذور الجذوع العربية، وبالتالي فإن ما يفترض أن يقدمه نظام التحليل الصرفي هو وجود كل الحلول الممكنة للكلمة، وإن لم يكن ذلك بوثوقية عالية لأن وجود قاعدة المعطيات سيرفع الوثوقية للمستوى المطلوب. وبناء على هذا التوصيف فالنظام المنجز يُخدم هذا النوع من التطبيقات أيضاً (مثل الحالة R7, R8).

النتائج:

نورد في الجدول التالي النتائج التي تعطيها النسخة R8 على الكلمات الواردة في وثيقة " مواصفات نظام التحليل الصرفي في اللغة العربية" للأستاذ مروان البواب. وهي - كما أسلفنا- النسخة الأفضل للنظام من حيث شمولية الحلول (توليد جميع الحلول الممكنة) وهي الحالة الأهم في تطبيقات المعالجة اللغوية كالمعاجم وأنظمة تعليم اللغة. علماً أن الحل الأفضل (إحصائياً) هو الحل الأول، ويمكن الحصول عليه مباشرةً باستخدام النسخة R8+.

يعطي المحلل الصرفي هذه النتائج دون الاعتماد على قاعدة معطيات المعجم أي باستخدام أقل تخزين ممكن للمعطيات. مع إمكانية الربط مع قاعدة معطيات المعجم (عند توفره).

الحل الثامن	الحل السابع	الحل السادس	الحل الخامس	الحل الرابع	الحل الثالث	الحل الثاني	الحل الأول	الكلمة
							سمع	سمع
							علم	تعلم
							سحب	سحب
					هيم	همم	وهم	وهم
				رمي	ورم	رمم	فرم	فرمت
						كيب	كتب	وليكتب
					سعي	سوع	لسع	لسعت
							ضرب	ضربت
							ضرب	ضربها
							ضرب	ضربتها
						زجج	زوج	زوجناكها
				دلا	دول	دلل	دلك	دلك
							أخذ	أخذنا
							افتقر	افتقر
							كمش	ينكمش
							ابتعد	ابتعد
					رمي	رمم	كرم	أكرم
			سرا	سدر	سور	سرر	يسر	يسر
						صيد	صيد	اصطادوا
				صير	صور	صرر	نصر	نصر
							علم	علم
							صفح	تصفح
							قبل	استقبل
							عرقل	عرقل
							زلزل	تزلزل
								اطمأن
							علم	علم
					قمم	قوم	سقم	استقام
					قلل	قيل	قول	يقال
			رثا	ريث	روث	رثث	ورث	يرث

					رث	ريث	روث	تراث
							أول	آلاء
							أوب	مآب
					تب	توب	كتب	كتاب
							جهل	مجهول
							غرب	مغرب
						فوح	فتح	مفتاح
							جلس	جلسة
							ذهب	مذهب
							جبر	جبار
							شعب	شعب
							علم	عالم
							عظم	عظيم
						كب	كتب	مكتب
					رحي	روح	فرح	فرح
		نسا	نوس	سنن	سنا		أنس	إنسانية
							نظر	نظرة
							حمد	محمد
			بور	بر	برا		كبر	كبرى
							عسل	عسل
							فوح	تفاح
							درس	مدارس
							عطش	عطشان
							كلمة_جامدة	إبراهيم
							كلمة_	
							كلمة_	عثمان
						سد	سود	سوداء
							قمر	قمر
		رضو	رضا	روض	رضض		أرض	أرض
							طرق	طريق
							شمس	شمس
		يوم	منن	يم	مون		يمن	يمين

							عنق	عنق
							خضراوان	خضراوان
		علو	عيل	عول	علل	علن	علو	أعلون
							صحراوات	صحراوات
						عوص	عصا	عصوان
			رضا	روض	رضض	رضو	أرض	أرضون
				هوا	موه	أمم	أمه	أمهات
							عطا	معطيات
					عوم	عمم	عما	أعماوي
					غيا	لغا	لغا	لغية
						عوص	عصا	عصية
						ظبي	ظبي	ظبي
				ديم	دوم	دمم	دما	دمي
				خيل	خول	خلل	خلا	خلوي
							بيض	بيضاوي
					فوت	فتت	فتي	فتوي
						أوب	أبا	أبا
	حنا	حوا	حنن	لحح	حين	لوح	لحن	فلاحون
							صبح	بمصاييح
		فوت	لفا	ليف	فتت	لفف	فتي	للفتي
				منا	مون	منن	كلمة_خاصة	من
							حذر	حذار
					كفف	كيف	كلمة_خاصة	كيف
							درس	المدرسة
						حث	كلمة_خاصة	حيث
					سفف	سوف	كلمة_خاصة	سوف
				بلا	بول	بلل	كلمة_خاصة	بل
				عدا	عود	عدد	وعد	وعدت
							فصل	ومفصلك
	كمي	كوم	منا	مون	يمن	كمم	منن	كلمة_خاصة