

مشكلات استخدام اللغة العربية فى نظم استرجاع المعلومات الببليوجرافية

أسامة لطفى محمد

مدرس مساعد: قسم المكتبات

كلية الآداب جامعة المنوفية - مصر

email: OS - LOTFY @ FRCU. EUN. EG

أولاً: المعايير

أ: شفرات تمثيل الحروف العربية.

منذ اختراع الحاسب فى بداية الأربعينيات، وحتى منتصف السبعينات اقتصرت معالجة الحاسبات على البيانات المعتمدة على الحروف الرومانية، وبدأ فى منتصف السبعينات يظهر اهتماماً كبيراً فى استخدام الحاسب مع النصوص غير الرومانية، وتنقسم الحروف غير الرومانية إلى قطاعين:

1 - القطاع الأول: لا تزيد حروفه عن 256 حرف، ومنها الحروف العربية والحروف الروسية

2 - القطاع الثانى: الذى تزيد حروفه عن 256 حرف، مثل: الحروف الصينية والكورية.

ومما يسر عملية التعريب هو عدم حاجة الشفرة العربية إلى أكثر من 256 حرف⁽¹⁾.

ولكن ظهرت مشكلة حيث اختلفت الأكواد Code Pages باختلاف شركات الإنتاج لبرامج التعريب، ولكن منذ عام 1976 بدأت مجموعة من الجهود العربية فى محاولة تطوير كود عربى موحد، نتج عنها ما يعرف بالأسمو 449 (ASMO 449)، وهى الشفرة العربية للحروف ذات السبع محارف (7 Bits)، والمستخدم فى الحاسبات الصغيرة (Mini - Computers)، و (ASMO - 708) ككود عربى للثمانية محارف⁽²⁾.

ويلاحظ أن استخدام ASMO - 708 فى بناء قواعد البيانات الببليوجرافية كاد أن يقتصر على قواعد البيانات المبنية باستخدام الإصدار الثالثة من برنامج CDS / ISIS وذلك لعدم التزام مصممو النظم فى العالم العربى باتباع هذه الشفرة، بل تم استخدام الشفرات الخاصة بنظم التعريب المستخدمة مثل (نافذة ومساعل العربى وفجر... إلخ).

هذا إلى أن ظهرت بيئة Windows العربية ونتيجة لطغيان Windows كبيئة للتشغيل أصبحت الشفرة الحروف العربية الخاصة بها هى معيار الأمر الواقع للثمان محارف.

ومن أهم التطورات الأخيرة هو ظهور الشفرة الموحدة للحروف UNICODE والتي تغطى جميع حروف

اللغات الحية ومنها العربية - حيث تسمح بتمثيل 65000 حرف وقد ظهر منها الطبعة الثانية حتى الآن إلا أنها لم تنتشر الانتشار المتوقع حتى الآن سواء على المستوى العالمى أو على المستوى العربى ويخشى أن تلقى نفس مصير 708 - ASMO .

ب - بنية التسجيلة البيولوجرافية :

أما بالنسبة لاستخدام الحاسب الآلى، نجد أن أهم الجوانب التى تعتمد على المعايير: هو جانب إنشاء بنية التسجيلات البيولوجرافية، وذلك لأهمية هذا النوع من المعايير عند التعاون بين المكتبات، أو عند شراء التسجيلات البيولوجرافية سواء على الخط المباشر، أو على أقراص مدموجة، أو على أقراص ممغنطة .

وتنقسم هذه المعايير إلى ثلاث مستويات، هى :

أولاً: مستوى الهيكل العام للتسجيلة البيولوجرافية:

ويتناوله معيار واحد حالياً وهو ISO 2709 .

ثانياً: المستوى التفصيلى لتيجان وأسماء الحقول والحقول الفرعية

وتغطيه مواصفة MARC و UNIMARC و CCF .

ثالثاً: مستوى محتويات حقول التسجيلة البيولوجرافية

والذى تغطيه قواعد الفهرسة الأنجلو أمريكية وخطط التصنيف المختلفة وقوائم رؤوس الموضوعات أو الكاتز .

وسوف نتناول معايير المستوى الأول والثانى فيما يلى:

1 - معيار ISO 2709 :

والخاص بتبادل التسجيلات البيولوجرافية على الأشرطة الممغنطة، حيث يعطى صيغة عامة لتبادل التسجيلات البيولوجرافية وتسجيلات ذات العلاقة بها، مثل تسجيلات الاستناد، وهو لا يحدد طول أو محتويات التسجيلات بينما هو يختص بتقديم معيار دولى يصف هيكل عام، أو إطار مصمم خصيصاً ليستخدم فى الاتصال بين نظم تجهيز البيانات، وليس للاستخدام فى التجهيز داخل تلك النظم، وعلى الرغم من أنه صمم مبدئياً ليعمل على الشرائط الممغنطة، إلا أنه يمكن استخدامه مع أى من وسائط البيانات الأخرى Data Carriers، وهو بصفة عامة يقسم التسجيلة البيولوجرافية إلى:

معرف التسجيلة (الفتاح) Record Labels

دليل المحتويات Directories

حقول البيانات Data Fields

فواصل التسجيلات Record Separators

ولا توجد مشكلة للتعامل مع البيانات العربية عند هذا المستوى حيث أن المعيار لا يقتصر على شفرة محارف معينة .

2- مواصفة UNIMARC و USMARC :

ثاني هذه المعايير: هي مواصفة Universal MARC UNIMARC، وهذه الصيغة تحدد المعارف Tags ومؤشرات Indicators، ورموز الحقول الفرعية Subfield Codes المستخدمة في التسجيلات البيولوجرافية، وهي تعمل تحت مظلة المعيار ISO 2709، ولكنها تقوم بوظيفة أكثر تفصيلاً وهي تحديد حقول ثابتة التيجان والمواصفات للتسجيلات البيولوجرافية، وقصد بها أن تستخدم كصيغة عامة لتبادل التسجيلات البيولوجرافية على المستوى الدولي أو الوطني.

3- تركيبة التراسل المشتركة CCF :

ثالث هذه المعايير وأحدثها، هو مواصفة CCF، والتي بنيت اعتماداً على المعايير السابقة لها، ومنها UNIMARC، وتتضمن أيضاً تعريف المعارف، والمؤشرات، ورموز الحقول الفرعية، وهي تتميز عنه بتغطيتها الأشمل لكل من المنفردات والأوعية غير المستقلة التي لا تغطيها MARC، وقد صدر منها ثلاث طبعات، آخرها تمتد تغطيتها إلى قواعد البيانات غير البيولوجرافية للأشخاص والمؤسسات ومشروعات البحوث، وقد صدرت الطبعة الثالثة عام 1992 في جزئين الأول التسجيلات البيولوجرافية، والثاني البيانات الحقائقية Factual Data.

وعند هذا المستوى تختفى المعايير العربية، وما يستغرب له أن أكثر المعايير استخداماً في العالم العربي هو معيار USMARC والذي فرضته النظم التجارية العربية.

ثانياً: خصائص اللغة العربية

ينبغي أن نؤكد أن عملية تعريب برنامج ما لا تعنى فقط مجرد ترجمة للرسائل والقوائم التي يعمل من خلالها النظام، حيث أن للغة العربية خصائص تجعلها متميزة ومختلفة عن اللغات الأخرى الأجنبية، وتؤثر هذه الاختلافات على قدرة برامج استرجاع النصوص بصفة خاصة، أي يجب عند التعريب مراعاة هذه الخصائص للغة العربية حتى نقلل من تأثيرها في جودة الاسترجاع⁽³⁾.

على سبيل المثال لا الحصر :

1- ترتبط أداة التعريف بالكلمة على عكس اللغة الإنجليزية، فبالتالي يسهل في الإنجليزية استبعاد أداة التعريف The عن طريق ملف الكلمات الموقوفة، ونتيجة أيضاً لوجود مسافة بينها وبين الكلمة سوف نجد أن جميع تكرارات الكلمة سواء المعرف منها أو غير المعرف تحت مدخل واحد في الملف المقلوب أو في الكشافات المطبوعة، أما بالنسبة للغة العربية، فلو اكتفى بمجرد ترجمة البرنامج، سوف يوجد لدينا كم كبير من الكلمات المكررة في ترتيبها الهجائي تحت حرف الألف واللام مثلاً وكذلك تحت الحرف التالي لها، مثال: العلم، علم، حيث نجد أنها تتكرر مرة تحت حرف الألف، ومرة أخرى تحت حرف العين، وهكذا تشتت تسجيلات المصطلح الواحد تحت مدخلين.

2- يوجد عدد من البادئات ترتبط بالكلمة، هي: همزة الاستفهام، الواو، الكاف، الفاء، اللام، الباء، السين، الياء، بالتالي يمكن أن تتكرر الكلمة نفسها في الملف المقلوب تحت حروف تلك البادئات بالإضافة إلى حروفها الأصلية، مثل: بمصر، كمصر، فمصر، . . . إلخ، وذلك بالإضافة إلى مصر، فيؤدي ذلك إلى تشتت المصطلح الواحد إلى عدة تكرارات في الملف المقلوب مما يؤدي إلى زيادة في الحجم الخاص به، والتأثير على جودة الاسترجاع نتيجة لتشتت المصطلح.

3- لا تختلف اللغة العربية كثيراً عن اللغة الإنجليزية في وجود اللواحق، والتي يتم معالجتها عن طريق استخدام أسلوب البتر، ولكنها بالتأكيد تختلف في تغير حروف الكلمة تبعاً لموقعها من الإعراب، وتختلف أيضاً في وجود الكثير من صيغ الثنية والجمع للمصطلح على عكس اللغة الإنجليزية التي في أكثر من 90% من الحالات يكون الجمع بإضافة حرف S أو es لنهاية الكلمة.

4- تختلف معاني بعض الكلمات في اللغة العربية باختلاف علامات التشكيل المستعملة معها عِلْم، عِلْم، عِلْم أو كَتَبَ، كُتِبَ... وفي حالة استخدام علامات التشكيل فتواجهنا مشكلة عند الفرز - حيث يتم التعامل معها كحروف مضافة للكلمة ويجب دائماً البحث بالتشكيل المطلوب.

نتعرف على أكثر هذه المشكلات تأثيراً على نظم الاسترجاع للتسجيلات البيولوجرافية العربية، والتعامل معها، - ونركز هنا على التسجيلات البيولوجرافية، وذلك لاختلاف اللغة المستخدمة في وصف أوعية المعلومات عن تلك التي تستخدم في نصوص أوعية المعلومات ذاتها -، وقد تم عمل مجموعة من التجارب للمقارنة بين طبيعة البحث في كلمات حقل العنوان لتحديد أهم هذه المشكلات:

حيث تم إنشاء ملف مقلوب لجميع الكلمات الواردة في حقل العنوان لعدد 2500 تسجيلة عربية من قاعدة البيانات لنظام الفهرسة الخاص بمكتبة كلية الآداب جامعة المنوفية، وتم أيضاً عمل ملف مقلوب لعدد 2500 تسجيلة باللغة الإنجليزية من قاعدة بيانات الفهرسة لمكتبة مركز بحوث التنمية والتخطيط التكنولوجي ووجد الآتي:

بالنسبة للغة العربية:

- 1- الكلمة بدون «ال» في 37% من الحالات.
- 2- تظهر الكلمة معرفة «ال» في 62% من الحالات.
- 3- تظهر الكلمة مع بادئة من البادئات في أقل من 1% من الحالات.
- 4- أكثر البادئات تواتراً (باستثناء الواو) في الظهور في عناوين الكتب هي «لل» و«الباء»، فعلى سبيل المثال تظهر كلمة المكتبات في الملف المقلوب بعدة أشكال هي:

المكتبات

مكتبات

بالمكتبات

للمكتبات

المكتبة

مكتبة

بينما يظهر نفس المصطلح في الملف المقلوب لقاعدة البيانات الإنجليزية بالأشكال التالية:

Library

Libraries

Librarian

Librarianship

ويلاحظ أنه في حالة استخدام اللغة الإنجليزية فإن استخدام أسلوب البتر في البحث سوف يسترجع جميع المصطلحات، حيث أن جميع التغييرات تظهر في اللواحق.
أما في حالة استخدام اللغة العربية لن يصلح أسلوب البتر في البحث لاسترجاع جميع الأشكال التي ظهر بها المصطلح لأن هناك تغييرات على مستوى البادئات، وفي المثال المذكور لن يتم استرجاع كل من:

المكتبة

المكتبات

بالمكتبة

فيما يلي الحلول التي لجأ إليها الباحث لحل هذه المشكلات من خلال استخدام نظام CDS / ISIS .
أولاً: أداة التعريف «ال»: -

وهي أهم هذه المشكلات نظراً للنسبة العالية لتكرار الكلمات المعروفة بال ويمكن التعامل مع هذه المشكلة بعدة من الطرق:

1 - عن طريق وضع أداة التعريف بين علامتي « > ... < » عند إدخال البيانات - حيث أن برنامج- CDS ISIS يتعامل مع الحروف الموضوعة بين هاتين العلامتين « > ... < » بأن يستخدمهما في العرض والطباعة فقط، ولا يستخدمهما سواء في بناء الملف المقلوب أو في عملية الفرز - ولكن هذا يمثل جهداً كبيراً على مدخلى البيانات ويستهلك الكثير من الوقت.

2 - عمل برنامج يقوم بهذه العملية آلياً، وفي هذه الحالة تظهر مشكلة الكلمات التي يكون فيها حرفا الألف واللام جزءاً أصلياً من الكلمة: مثل: ألمانيا، ألبانيا، ألم... إلخ.

3 - يتم استخدام برنامج لإضافة علامتي « > < » آلياً مع الاعتماد على ملف يتضمن الكلمات المستثناة أو الموقوفة عن هذه العملية هو يختلف بالطبع عن ملف الكلمات الموقوفة لقاعدة البيانات، حيث يمنع الثاني الكلمات من الدخول في الملف المقلوب كلية بينما يمنع الأول حذف الألف واللام من الكلمات الموجودة به.

وقد استخدم هذا الحل في إنشاء فهرس المجموعة العربية لمكتبة جامعة الملك فهد للبترول والمعادن، وذلك بوضع ملف إيقاف للكلمات التي تتضمن «ال» أصلية يتضمن حوالي 200 كلمة⁽⁴⁾، ولكن من الأفضل أن يتم بناء هذا الملف تدريجياً تبعاً لاحتياجات قاعدة البيانات دون الاعتماد على ملف سابق التجهيز، مما يؤكد ذلك هو أن عدد الكلمات التي تضمنت «ال» أصلية في عناوين جميع الكتب العربية في مكتبة كلية الآداب جامعة المنوفية هو ثلاث كلمات هي: «الف / الفية / الله».

4 - تعديل برنامج البحث بحيث يقوم المستفيد بإدخال المصطلح ثم يقوم البرنامج بتوليد جميع الأشكال الممكنة للمصطلح وإدخالها كبداية للبحث باستخدام الرابطة المنطقية «أو».

مثال:

يقوم المستفيد بإدخال مصطلح حاسب فيقوم البرنامج بتوليد الأشكال التالية:

«الحاسب أو يحاسب أو فحاسب أو لحاسب أو للحاسب أو كحاسب أو بالحاسب... إلخ».

ولكن هذا الأسلوب قد ينتج عنه كلمات بلا معنى أو أخطاء صرفية مثل: المصّر، هذا بالإضافة إلى أنه لا يؤثر على عملية الفرز، حيث سوف تخرج الكلمات بنفس الترتيب للبادئات.

5 - استخدام البتر في اتجاه اليمين ولكن هذا الحل يعترضه ندرة برامج استرجاع النصوص التي تدعم البتر في الاتجاهين والبتر الداخلي، ولكن هذا البديل أيضاً لن يؤثر في عمليات الفرز.

ثانياً: مشكلة البادئات:

لا يمكن استخدام الطريقة الثانية أو الثالثة في علاج هذه المشكلة، حيث لا يمكن جمع جميع الكلمات التي تحتوي على هذه البادئات ووضعها في ملف إيقاف، أى معنى ذلك أن نضع جميع الكلمات العربية التي تبدأ بتلك الحروف التسعة، وأيضاً يمثل الحل الأول، والذي يتمثل في وضع علامتى «>» حول البادئة «مثل: <ب> مكتبة <ب> المكتبات» عبئاً كبيراً على مدخلى البيانات، وبهذا نجد أنه:

- 1 - يفضل استخدام الأسلوب الثالث بالنسبة للألف واللام لأهميته في الفرز، واستخدام نفس الأسلوب بالنسبة ل «لل».
- 2 - بالنسبة للواو، فيتم كتابتها مفصولة عن الكلمة بمسافة ويتم وضعها في ملف الإيقاف، كى لا تؤثر على حجم الملف المقلوب المستخدم فى البحث.
- 3 - ويفضل استخدام الأسلوب الرابع بالنسبة لباقي البادئات حيث لن تؤثر بشكل كبير على جودة الفرز، حيث لم يوجد عنوان واحد من ٢٥٠٠ عنوان يبدأ ببادئة.

أما بالنسبة إلى البديل الخامس وهو استخدام محلل صرفى ليؤدى نفس الوظيفة التي يقوم بها البتر في اللغة الإنجليزية، وقد تم إنشاء ما يعرف بالمعالج الصرفى متعدد الأطوار يستخدم فى الاسترجاع لكلمات القرآن الكريم، وبعض التطبيقات الأخرى مثل ضغط (compression) الحروف العربية والتدقيق الإملائى العربى⁽⁵⁾.

ولكن هنا يجب أن نشير إلى استخدام جميع الجذور للكلمة العربية بدلاً من البتر أو معالجة البادئات سوف يؤدى نتائج غير حميدة فى مجال استرجاع التسجيلات البيولوجرافية، ولتخيل كم الكلمات التي سوف تسترجع بالجذر: علم أو كتب بدلاً من معلومات أو مكتبات، فسوف يسترجع البرنامج عدد كبير من الكلمات التي لا تمت بصلة موضوعية للبحث أكبر من عدد التسجيلات ذات الصلة. حيث أن المحللات الصرفية للغة العربية لم تخرج بعد من إطار توليد الكلمات من الجذر لاستخراج الصيغ الصرفية والتي تستخدم فى البحث.

ثالثاً: استخدام اللغة العربية فى البيئة متعددة اللغات

والمقصود بالبيئة متعددة اللغات فى هذا السياق هو تواجد تسجيلات بيولوجرافية بأكثر من لغة طبيعية مثل العربية والإنجليزية، وقد اعتمدت مراكز المعلومات فى أغلب الأحيان على الفصل بين تسجيلات اللغتين، ولكن هذا الحل غير عملى بالنسبة للمستفيد الذى يريد إجراء بحث موحد فى جميع التسجيلات

الموجودة في المكتبة في موضوع معين بغض النظر عن اللغة، وتمثل هذه الطريقة في المعالجة بأن يتعامل النظام المتكامل مع قاعدة بيانات بيلوجرافية موحدة تتضمن التسجيلات البيلوجرافية باللغتين. ولكن استخدام قاعدة بيانات موحدة لتسجيلات اللغتين يؤدي إلى عدد من المشكلات عند التطبيق تتمثل في:

أولاً: وجوب توحيد محددات الحقول الفرعية لشكل الاتصال المستخدم:

فعلى سبيل المثال تستخدم الحروف الإنجليزية كمحددات للحقول الفرعية للنسخة الأصلية من CCF، وتستخدم الحروف العربية في النسخة العربية من CCF، وهنا تظهر مشكلة عند تصميم جدول اختيار الحقول وصيغ العرض، حيث إذا اتبعنا معيارين في نفس قاعدة البيانات، سوف نحتاج إلى عمل نسختين من كل من صيغ العرض، ومضاعفة عدد المداخل في الملف المقلوب.

وبهذا يقترح تعديل مؤشرات الحقول الفرعية في كل من النسخة العربية والإنجليزية من CCF إلى استخدام الأرقام بدلاً من الحروف الهجائية، سواء الإنجليزية أو العربية، وذلك لأن الأرقام لها شفرات موحدة في اللغتين، مما يؤدي إلى إمكانية توحيد صيغ العرض لتسجيلات اللغتين وعدم تكرار مداخل الملف المقلوب وسهولة البرمجة، بالإضافة إلى توحيد محددات الحقول الفرعية عند إدخال البيانات لتسجيلات اللغتين.

ثانياً: مشكلة توحيد أشكال العرض:

حيث يؤدي استخدام صيغة للعرض تتضمن نصاً مضافاً لبيانات التسجيلة - عنوان الحقل مثلاً - مع تسجيلة من لغة مخالفة إلى اضطراب وتشويش في شكل المخرجات على الشاشة أو على الطابعة. وحلاً لهذه المشكلة قام الباحث بالاستعانة بلغة الصياغة لبرنامج CDS / ISIS في إعداد مواصفات لأشكال العرض ذات حساسية للغة، حيث تعرض النصوص المضافة العربية مع التسجيلة العربية وتتغير أياً إلى الإنجليزية في حالة التسجيلة الإنجليزية.

ثالثاً: سياسة التحليل الموضوعي المستخدمة:

ونظراً لطبيعة قواعد البيانات المتعددة اللغات، يجب أن تتبنى المكتبة سياسة محددة تجاه رؤوس الموضوعات حيث توجد خيارات عدة، هي:

1 - استخدام رؤوس موضوعات أو واصفات باللغتين، بمعنى وجود رؤوس موضوعات عربية للتسجيلات العربية، ورؤوس موضوعات إنجليزية للتسجيلات الإنجليزية، مما قد يسبب مشكلة في اختيار قائمتين مختلفتين لرؤوس الموضوعات تعتمد كل منهما على فلسفة خاصة بها في تصنيف المعرفة، مما يؤثر على التوافق Consistency في سياسة التحليل الموضوعي، ويفصل بين التسجيلات ذات الموضوع الواحد المختلفة اللغة.

2- استخدام رؤوس موضوعات بلغة واحدة العربية أو الإنجليزية لتجميع التسجيلات باللغتين، وسوف يحرم ذلك المستفيدين من البحث باللغة الأخرى والتي قد تكون اللغة الأم لهم.

3- استخدام مكتز متعدد اللغات: وهنا يجب أن يتقبل البرنامج الأساليب المستخدمة فى إدارة المكتاز المتعددة اللغات، والتي تبنى أساساً على قدرة النظام على استدعاء التسجيلات لرأس موضوع بلغة معينة، بغض النظر عن لغة الرأس الذى استخدم فى الوثيقة، أى أن للمكتبة الحرية فى إدخال الرأس بأى لغة من اللغتين، ويقوم البرنامج بالاسترجاع باللغتين.

وقد اختير الحل الثالث، حيث تم تطوير البرنامج الخاص بالمكتاز ليسمح بإدارة المكتاز متعددة اللغات، حيث يتم إدخال الواصفات باللغتين العربية والإنجليزية إلى المكتز، ويقوم البرنامج باستدعاء المكتز عند الإدخال، حيث يسمح بإدخال المصطلح بأى من اللغتين، وفى حالة البحث يقوم البرنامج بالبحث بالمصطلح وبديله باللغة الأخرى.

وهناك أكثر من تجربة لإنشاء مكتاز ثنائية اللغة فى الوطن العربى، منها:

«المكتز العربى للنشاطات الاجتماعية والاقتصادية والسياسية»

والذى صمم بواسطة مجلس الوزراء الكويتى فى بداية الثمانينات.

وهناك أيضاً:

المكتز العربى للبترول

والذى انتهى العمل فيه عام 1987 فى المعهد العربى للبترولى «للتدريب» التابع لمنظمة الأوبك⁽⁶⁾.

وهذا بالإضافة إلى: مكتز جامعة الدول العربية ثلاثى اللغة: عربى، إنجليزية، فرنسى.

رابعاً: مشكلات إتاحة قواعد البيانات العربية من خلال شبكة إنترنت:

ظهر عدد من المشكلات عند إتاحة قواعد البيانات العربية من خلال شبكة إنترنت من أهمها ما يلى:

1- الاختيار بين الإتاحة من خلال برنامج متصفح يعمل من خلال mswindows العربية مما يحجب البيانات العربية عن العالم الخارجى لعدم توافر النسخة العربية خارج حدود الوطن العربى أو الاعتماد على برنامج مساعد Plugin يسمح بعرض الحروف العربية على جميع إصدارات Windows الأجنبية ولكن هذا يزيد العبء على المستخدم من ناحية ولا توجد معايير لهذه البرامج من ناحية أخرى.

2- زيادة عبء عملية التعريب على الأجهزة الخادمة Servers والتي سوف تستخدم لعرض البيانات العربية بشفرة لاتوجد على نظام التشغيل الأسمى لها لكى تعرض على متصفحات موجودة على أجهزة عملاء clients تعمل على أكثر من نظام تشغيل.

3- معيار Z 39.50.

يضيف معيار Z 39.50 بعداً جديداً لتعقيدات تعريب قواعد البيانات الجيوجرافية على شبكة إنترنت يتمثل فى تعريب كل البرامج الوسيطة بين الخادم والعميل والتي ترأسل بشكل معيارى من خلال Z 39.50 وذلك لأن أغلب هذه البرامج الوسيطة تتعامل مع النصوص بشفرات ذات 7 محارف مما يشوهه البيانات العربية ذات الثمانية محارف.

المراجع

(1) - Richard Lee, (1978). "MINISIS: a multilingual information management System". Automated systems for access to multilingual and multiscrypt library materials problems and solutions. - Munchen: IFLA, 1987. p. 215

(2) - Aman, Mohammed M. (1987). Use of Arabic Script in computerized information systems: Automated systems for access to multilingual and multiscrypt library materials problems and solutions. - Munchen: IFLA, 1987, p. 129.

(3) تجدر الإشارة في هذا الصدد إلى مجموعة من الجهود التي تمت لدراسة تأثير اللغة العربية على كفاءة الاسترجاع في نظم المعلومات البليوجرافية.

ونستطيع معالجة هذه الدراسات من حيث كيفية تعرفها للمشكلة - فهناك الدراسات التي قامت بالرصد النظرى لبعض خصائص اللغة العربية ومدى تأثيرها المتوقع على مدى كفاءة الاسترجاع.

فقد تعرض البعض لاختلاف بنية الكلمة العربية عن الإنجليزية من نواح عدة، ومنها:

بناء الكلمات، الكلمات المتصلة، السوابق، اللواحق... إلخ

ومن هذه الدراسات:

(أ) - ييار فيرميل. معالجة المعلومات في اللغة العربية. «الإعلامية والتعريب». بيروت - لبنان: الوكالة الإعلامية الفرنسية، جمعية معالجة اللغة العربية في الإعلامية، 1984. ص ص 4 و 8.

وهناك بعض الدراسات التي تناولت تأثير بعض صفات اللغة العربية في الاسترجاع الفعلى نتيجة لتجارب عملية، مثل:

(ب) - ناصر محمد السويدان. «الاسترجاع الموضوعى بواسطة كلمات العنوان». السجل العلمى لندوة استخدام اللغة العربية في تقنية المعلومات. الرياض: مكتبة الملك عبد العزيز العامة، 1993. ص ص 533 و 568.

والذى تعرض في دراسته لبعض مشكلات بناء كشافات الكلمات المفتاحية الناتجة عن «ال»، وحروف الجر والتصاقها بالكلمات في اللغة العربية.

كذلك الدراسة المقدمة من:

(ج) - سعد عبد العزيز. «نظام الوثائق: نحو نظام بليوجرافى عربى للوثائق الحكومية فى مكتبات معهد الإدارة العامة». السجل العلمى لندوة استخدام اللغة العربية فى تقنية المعلومات. الرياض: مكتبة الملك عبد العزيز العامة، 1993. ص ص 337-344.

حيث تعرض لبعض المشكلات الناتجة عن الخلط بين أشكال كل من: ء، أ، ال، ت، ه، ي، الالف المقصورة.

والدراسة المقدمة من :

(ء) بحيث سليمان البخيت. «البحث فى العنوان فى قواعد البيانات العربية: دراسة تطبيقية على خدمة برمجيات CDS / ISIS». السجل العلمى لندوة استخدام اللغة العربية فى تقنية المعلومات. الرياض: مكتبة الملك عبد العزيز العامة، 1993. ص ص 569 - 580 .

والذى قام بحصر تكرارات (الكاف، الباء، اللام فى عناوين 5000 (خمسة آلاف) تسجيلة عربية مخزنة على برنامج CDS / ISIS، وقد اقترح إضافة البتر من اليمين لتحسين كفاءة الاسترجاع فى برنامج CDS / ISIS.

وأخيراً توجد الدراسات التى تعرض لبعض الحلول العلمية التى تم استخدامها فعلياً عند استخدام اللغة العربية ومن هذه الدراسات:

(هـ) - Zahiruddin Khurshid. (1992). "Arabic Online Catalog" Information Technology and Libraries. September 1992, pp. 244 - 251.

حيث قام بعرض بعض الأساليب المستخدمة فى نظام DOBES / LIBIS والتى اعتمدت على وضع علامات خاصة قبل وبعد البادئات واستخدام ملف إيقاف لمعالجة (ال). والدراسة التى قدمها:

(و) - سريع محمد السريع. «نظام ابن النديم فى مكتبات معهد الإدارة العامة». السجل العلمى لندوة استخدام اللغة العربية فى تقنية المعلومات. الرياض: مكتبة الملك عبد العزيز العامة، 1993. ص ص 315 - 336.

والذى عرض لبعض مشكلات بناء كشافات الاسترجاع للكلمات العربية واتبع أسلوباً يعتمد على وضع رمز من أربعة رموز قبل بعض الكلمات عند الإدخال لأخذ القرار فى شكل دخول هذه الكلمات إلى الكشف.

(4) - Zahiruddin Khurshid. (1992). "Arabic Online Catalog". Information Technology and Libraries. September, 1992. p. 249.

(5) نبيل على. (1988). اللغة العربية والحاسوب . - القاهرة: دار تعريب للنشر، 1988. ص ص 330 - 332..

(6) - Shawky Salem. (1991). "Computerized Bilingual Thesauri". Microcomputer for information management. Mar 1991. p. 29.

(*) تجدر الإشارة إلى وجود البرنامج لدعم إنشاء المكانز العربية تم تطويره فى المملكة العربية السعودية ((راجع «عبد الجبار عبد الرحمن العبد الجبار». «استخدام نظام المستشار فى بناء المكانز العربية». السجل العلمى لندوة استخدام اللغة العربية فى تقنية المعلومات. الرياض: مكتبة الملك عبد العزيز العامة، 1993. ص ص 631 - 647)).