

## منظومة قراءة آلية لمخطوطات التراث العربي

أ.د. محمد يونس الحملاوي  
قسم هندسة النظم والحاسبات  
كلية الهندسة، جامعة الأزهر

د.محمد يسرى النحاس  
قسم هندسة النظم والحاسبات  
كلية الهندسة، جامعة الأزهر

### نبذة

واسعا فى الآونة الأخيرة [٢]، [٣]، [٤]، [٥]، [٦]. ويرجع ذلك إلى الأهمية التى تمثلها هذه المخطوطات والحاجة إليها فى بناء قواعد المعرفة المتكاملة، ومن ناحية أخرى فإن التطور السريع الذى حدث فى تقنيات المعلومات أتاح للباحثين فى مجال التعرف على الأنماط من تنفيذ الكثير من التطبيقات بنجاح مثل التعرف على الحروف المطبوعة والتعرف على الأصوات، الخ...

إن معالجة الوثائق والتعرف عليها يحتوى على الكثير من المشاكل التى مازالت قيد البحث والتجريب مثل ترميم صور الوثائق، والقصل بين الأنواع المختلفة من المعلومات المطبوعة سواء أكانت حروفا أو رسومات أو جداول، الخ...، وهذا التعدد فى الأنماط والانتشار الكبير للكائنات الممثلة لها يمثل تحديا حقيقيا من الناحية النظرية مما يتطلب إيجاد نظرية متكاملة للتحليل والتعرف على هذه النوعية من الصور بمعدلات عالية وثابتة. كما أن المنظومات التقنية المطلوب استخدامها لابه من أن توفر درجة عالية من المعالجات المتوازية حتى يمكن التعرف على الوثائق فى الوقت الفعلى بالنسبة للمستخدم.

وتتميز مخطوطات التراث المكتوبة بالعربية عن الكتابات المطبعية بتسوع أسلوب الخط

يقدم هذا البحث منظومة آلية للتعرف على مخطوطات التراث العربى. تعتمد هذه المنظومة منهج التعلم بدون إشراف للتعرف على نمط الحروف المكتوبة. وقد صممت إجرائية تعرف وتعلم تعمل "باستراتيجية تعلم كلما مضيت"، حيث تتم عمليات التعرف والتعلم بصورة مستمرة وديناميكية، أى أنه بدءاً من أساسيات مقترحة لأشكال الحروف يقوم الحاسوب بتحليل الصفحة ثم البحث فى نافذة معزولة من الصفحة عن الحروف وفى هذه الأثناء يقوم بتطوير معرفته لهذه الأساسيات بما يستنبطه من الأساسيات الموجودة فى صورة الحروف. وباستمرار عمليات التحليل والتعرف يتم تطوير منظومة الحروف ويقبل معدل أخطاء التعرف، وعند الوصول إلى نسبة عالية من الإستقرار فى التعلم تعيد المنظومة تحقيق نتائجها فى بدايات التعرف. طبقت هذه المنظومة على مجموعة من الحروف العربية فى مخطوطات التراث، وقد أظهرت نتائج التجارب الأولية بالإضافة إلى إرتفاع معدل التعرف، ثبوت هذا المعدل بتغير أسلوب الكتابة داخل المخطوطة ذاتها.

### ١ مقدمة

يعتبر موضوع القراءة الآلية لمخطوطات التراث من الموضوعات التى تشهد نشاطا بحثيا

على منهج التعلم بدون إشراف الذي يتيح للمنظومة تعلم نمط الكتابة في المخطوطة وأساسيات الحروف بالمواعمة.

## ٢ تعريف المشكلة

إن منظومة التعلم والتعرف الآلى على الوثائق يمكن تمثيلها بالشكل ١. تعمل هذه المنظومة عادة على مرحلتين: مرحلة التعلم، ثم مرحلة التعرف. وتتكون مرحلة التعرف من العناصر والعمليات الآتية:

١- المصورة ، ٢- معالجة الصورة ، ٣- تحليل الصورة ، ٤- تصنيف كائنات الصورة.

وتتكون مرحلة التعلم من العناصر والعمليات الآتية:

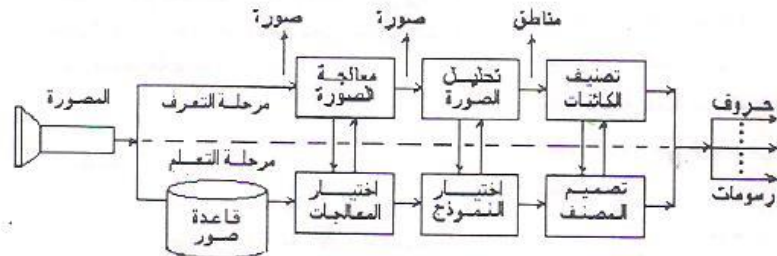
١- قاعدة صور ، ٢- اختيار معالجات الصورة ، ٣- اختيار نموذج الصورة ، ٤- تصميم المصنف.

إن مشكلة تصميم منظومة القراءة الآلية هي وضع نماذج العناصر، وتوصيف العمليات، وترتيب عملها في هذه المنظومة. وحل مشكلة التصميم هذه يتحدد بالمدخل المستخدم في التعرف على الأنماط. ونحن في هذا البحث نعتد المدخل

الذي خط يد في المقام الأول ، مع افتقاد النمطية في أسلوب عرض النصوص في الصفحات ، هذا بالإضافة إلى وجود بعض التزييلات والحواسي في من الصفحة. كل هذه الخصائص تمثل صعوبات نظرية وتقنية في عمليتي التعلم والتعرف على هذا النوع من الوثائق بالنسبة لمثيلاتها من الوثائق النمطية.

يهدف هذا البحث إلى توصيف المشاكل المرتبطة بالقراءة الآلية لمخطوطات التراث العربية وتصميم منظومة تحقق القراءة الآلية في الوقت القليل للمستخدم. والبدائل المطروحة في منظومات التعرف على الأنماط كثيرة مثل مصنفات 'بايز' [١] أو شبكات العصبية [١] أو سلاسل ماركوف المخيأة [٥] والاختيار بين أي من هذه المنظومات يفرضه نوع للتطبيق والمشاكل النظرية والعملية الخاصة به. كما أن أسلوب عمل المنظومة قد يختلف من النصف آلياً إلى الآلية الكاملة أيضاً حسب القيود التي يحددها التطبيق.

وأمام كثرة الأنماط الموجودة في صفحات مخطوطات التراث والانتشار الواسع لكل نمط منها وتغير خواصه الإحصائية من صفحة إلى أخرى فقد وجدنا أنه من الأنسب وضع منظومة تستخدم آلية عضوية محدودة في التعرف على الأنماط ، وتعتمد



شكل ١ منظومة القراءة الآلية للمخطوطات

المؤتمر السادس من الحاسب الآلى بين النظرية والتطبيق، الاسكندرية، ٣-٥ سبتمبر ١٩٩٦م

منظومة قراءة آلية لمخطوطات التراث العربي، محمد يمري الفعاصي، محمد يونس الحماوي

أرجاء هذه المرحلة لما بعد التعرف على الحروف الواضحة في المخطوطة.

أما المشكلتان الأساسيتان في معالجة هذا النوع من الصور فهما كيفية فصل مناطق الكتابة في الصفحة، ثم تليها مشكلة التعرف على حروف الكتابة، ونلاحظ في هذه الصورة وجود نوعين من المعلومات المرئية تتمثل في الكتابة والجدول. كما أنه توجد ثلاثة أنماط من الكتابة تتمثل في: الكتابة الأفقية، والكتابة العمودية، والكتابة المائلة. أما الخط المستخدم فهو خط يد حروفه متصلة ومتداخلة، وإن كان يتميز بالتمطية.

وهذا البحث يقدم بالتفصيل كيفية فصل مناطق الكتابة من الصفحة ثم تجزئتها إلى أسطر ثم إلى كلمات وحروف. سواء كانت الكتابة تمثل سطرا أفقيا أو سطرا مائلا، وسواء كانت الكتابة شرحا أو تنبيلا للصفحة. أما الجداول، أو الرسوم بصفة عامة، والكتابات الموضحة عليها فتمثل بصورتها كما هي. ومشكلة التعرف على الحروف هي تجزئتها كل سطر وتصنيف هذه الأجزاء إلى إحدى طبقات الحروف التسع والعشرين الممثلة لحروف الهجاء العربية. ومن المعروف أن عملية التصنيف هذه

الإحصائي-النحوي في عمليتي التعلم والتعرف كما سنوضح في تصميم عناصر هذه المنظومة.

هناك ثلاث مشكلات جزئية على التحديد يجب حلها عند التصميم: ١- معالجة عيوب الصورة، ٢- فصل مناطق الكتابة، ٣- التعرف على الحروف. ولتحديد نوعية هذه المشاكل الجزئية المطروحة سنستعين بصورة الصفحة الممثلة في الشكل ٢ وهو يمثل إحدى صفحات كتاب الكاشي [١١].

إن النظرة الفاحصة لهذه الصورة تبين أنها تحتاج إلى نوعين أساسيين من المعالجات:

١- ترشيح الشوشرة، ٢- وترميم الأجزاء الناقصة في الحروف والكلمات. ومشكلة ترشيح الشوشرة في هذه الحالة يمكن حلها بأحد المرشحات التقليدية المعروفة [٩]. أما مشكلة ترميم الأجزاء الناقصة فهي أكثر صعوبة ولن نتعرض لها في هذا البحث لأنها تتطلب المزج بين التعرف على أشكال الحروف والتعرف على الكلمات المكتوبة، الذي يتطلب استخدام معارف الصرف والنحو والمضمون [١]. والحل الأمثل في هذه الحالة هو



شكل ٢. صورة رقمية لصفحة من مخطوطة الكاشي [١١]

المؤتمر السادس من الحاميم الألي بين النظرية والتطبيق، الامنحدرية، ٢-٥ سبتمبر ١٩٩٦م

كل منها. وهذا يتطلب وضع نموذجاً لوصف مكونات الصورة، والنموذج الذي نستخدمه هنا لوصف مكونات الصورة هو النموذج الهرمي. أي أننا نفترض أن الصفحة مكونة من مجموعات تنقسم بصورة مكررة إلى مجموعات أصغر فأصغر من مستوى لآخر كما هو ممثل في الشكل ٣.

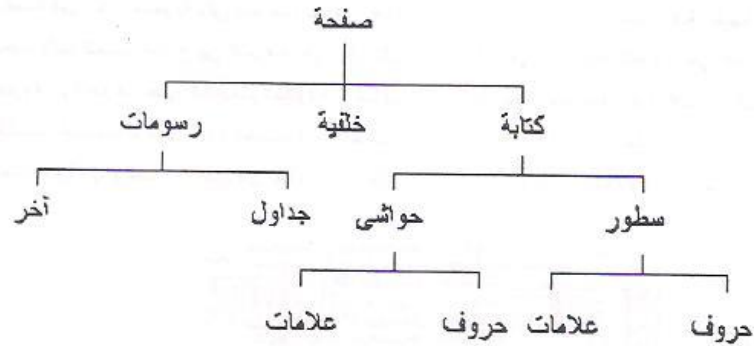
وبتطبيق عملية التحليل والتقسيم حتى المستوى الثاني للمخطوطة المصورة في الشكل ٢ ينشأ عنها التقسيم الموضح في الشكل ٤. في هذا الشكل أحيطت كائنات كل منطقة بناطقة أفقية مستطيلة لسهولة العرض وإن كان شكل هذه النافذة وميلها يعتمد على الشكل الهندسي للمنطقة وميل محاوره الأساسية بالنسبة لمحاور الصفحة.

تعتمد على دقة عملية التحليل السابقة وعلى مدى تكامل المعارف الخاصة باللغة في منظومة التعرف سواء بالنسبة لشكل الحرف، أو بالنسبة لقواعد الصرف والنحو إلى حد كبير.

ولحل هاتين المشكلتين فقد رأينا من الأسب وضع آلية تعمل بأسلوب المواعمة [٨] أثناء عمل المنظومة مما يسمح بثبات أداء المنظومة مع تغير نوع المخطوطات. كما أن هذه الآلية يجب أن تستخدم بصورة متكررة في حل المشكلتين.

### ٢ نموذج الصفحة

يهدف تحليل الصورة إلى تقسيمها إلى طبقات متباينة من الكائنات وكشف وقياس سمات



شكل ٣. النموذج الهرمي لصفحة من المخطوطة.



شكل ٤. تقسيم صورة الصفحة إلى مناطق من الكتابة والجداول

الإشراف. عند غياب مثل هذه المعلومات لكل منطقة على حدة يصبح منهج التعلم بدون إشراف هو البديل الأفضل في هذه الحالة. ومفهوم التعلم هنا يقصد به إيجاد التقسيم الأمثل لمجموعة من الكائنات ممثل كل منها بمتجه السمات بحيث تكون درجة التشابه بين كائنات المجموعة الواحدة أكبر ما يمكن ودرجة التباين بين المجموعات أكبر ما يمكن [١]. ويمكننا تعريف مشكلة تقسيم صورة الصفحة وضعبها على النحو التالي:

عند كل نقطة من نقاط الصورة يعرف متجه السمات،  $s$ ،

$$s = (s_1, s_2, \dots, s_n) \quad (1)$$

وتمثل نقاط الصورة فئة كائنات التعلم المطلوب توزيعها على الطبقات المختلفة، ويرمز لفئة كائنات التعلم بـ  $S$ ،

$$S = \{s_1, s_2, \dots, s_n\} \quad (2)$$

وتعرف مشكلة التقسيم على أنها كيفية تقسيم فراغ السمات،  $S$ ، إلى عدد  $m$  من المناطق،  $S_l$ ، تحقق الشرطين الآتيين:

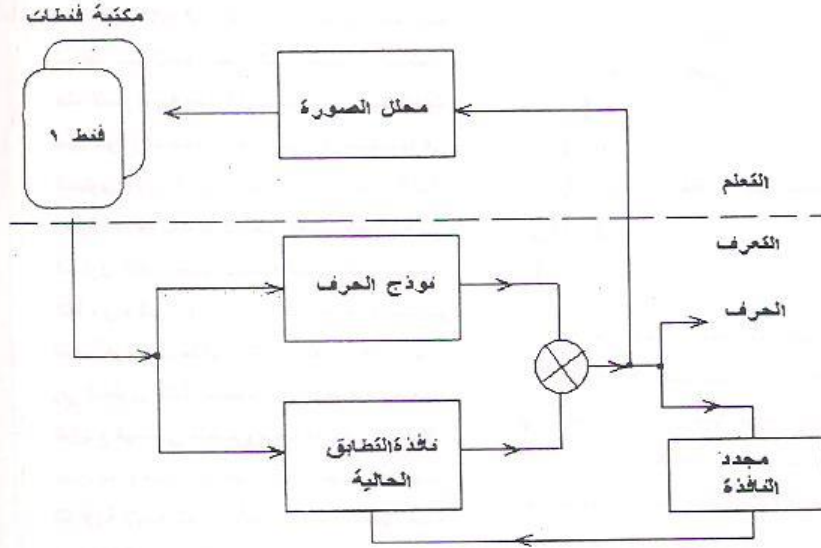
$$S_l \cap S_j = \emptyset \quad 0 \leq l < j < m \quad (3)$$

$$S = \bigcup_{l=1}^m S_l$$

#### ٤ تقسيم صورة الصفحة

يمكن تعريف عملية تقسيم الصورة بأنها عملية تصنيف نقاط الصورة، حيث تقسم الصورة إلى مناطق بتوزيع النقاط على الطبقات التي تنتمي إليها. ولتحقيق عملية التقسيم نفترض أن الصورة مكونة من مناطق تختلف في خواصها مثل المكانية، أو الهندسية، أو القياسية، أو شدة إضاءتها أو نسيجها أو في كل هذه الخواص معا [٩]. وعند توافر المعلومات الإحصائية عن هذه المناطق مسبقا نشرع في قياس كل من هذه الخواص عند كل نقطة من نقاط الصورة ونضعها في هيئة متجه يعرف بمتجه السمات، والسمات التي يمكن استخدامها في تقسيم الصفحة هي: مستوى الرمادية للنقطة، القيمة المتوسطة لمستوى الرمادية في نافذة محلية، إلخ... ثم يصنف متجه السمات باستخدام منظومة تصنيف مصممة على أسس قريبة من مصنف بايز المثالي. وتوفير المعلومات الإحصائية عن مناطق الصورة يتم من خلال المسح الشامل لصورة أو أكثر تحتوي على المناطق المفترضة فيما يعرف بالتعلم تحت





شكل ٥. منظومة تعلم أساسيات الحروف.

٢- يصنف الحرف إلى الطبقة الأقرب في النسق الأول، وتتأ طبقة جديدة في النسق الثاني نواتها الحرف المكتشف، ويلغى استخدام هذا الحرف في النسق الأول.

٣- لا يمكن تصنيف محتويات النافذة فيوجد اتخاذ القرار حتى يتوافر عدد مناسب من الكائنات الممثلة لكل طبقة في النسق الثاني. أو ينشأ نسق جديد يمثل فنطاً جديداً.

واختيار مسافة ماهاالاتوبيس هنا للتصنيف يضمن ثبات معدل التصنيف مع تغير موضع أو حجم أو ميل الحرف في الصفحة، وإن كان يتطلب عدداً أكبر من الكائنات الممثلة لشكل الحرف قبل استقرار نتائج التصنيف.

التعلم تمتلئ فئة النسق الثاني وينعدم استخدام فئة النسق الأول.

عند كل نقطة من نقاط مناطق الكتابة تقارن محتويات نافذة التطابق الحالية مع أشكال الحروف في الفنط القياسي ويصغر أو يكبر إلى حجم نافذة التطابق الحالية. ثم تصنف محتويات نافذة التطابق الحالية بقياس مسافة ماهاالاتوبيس بين نموذج الحرف والشكل الموجود بداخل النافذة، ونتيجة عملية التصنيف هي إحدى الخطوات الآتية:

١- يصنف الحرف إلى الطبقة الأقرب في النسق الثاني، وتعديل السمات الممثلة لشكل الحرف المكتشفة لاستخدامها في عملية التعرف التالية، وكذلك مكان وحجم نافذة التطابق الحالية.

## ٦ إجراءات التعلم والتعرف

إن الإجرائية المقترحة للتعلم والتعرف تستخدم استراتيجية "تعلم كلما مضيت" وتستخدم هذه الاستراتيجية عند كل مستوى من مستويات تمثيل نموذج الصفحة، وهذا يعني أننا نستخدمها في المستوى الأول لتقسيم الصفحة إلى مناطق كتابة، وخطية، ورسومات ثم تستخدم نفس الإجرائية في المستوى الثاني لتقسيم منطقة الكتابة إلى سطور كتلة، وحواشي، ولكننا في هذه المرحلة نستخدم سمات أخرى مثل مكانية الكائنات للفصل فيما بينها. وفي المستوى الثالث نستخدم نفس الإجرائية وسمات توزيع الهندسى للتعلم والتعرف على أشكال لحروف، وهذه الإجرائية هي تنفيذ للطريقة تكرارية لإيجاد الحل الأمثل [١٠]. وسنمثل هذه الإجرائية بالبرمجية الآتية:

إجرائية للبحث عن الطبقات (س، ق)

منخلات: فئة الكائنات المجهولة، س،

مخرجات: نسق التصنيف، ق،

طالما (حصل\_كائن(س، س) == لاثنين)

صنف - خطأ؛

طالما ((صنف == حقيقي) و او (ل > عدد

لطبقات))

ف - مسافة\_ماهاالاتوبيس(س، س، ل)؛

إذا (ف > ف٠) { حدث\_الطبقة (س،

س، ل)؛ صنف - حقيقي؛

وإلا إذا (ف > ف١)

}

إذا( الرصة == مملوءة)

صنف\_المخزون (رصة، ق)؛

دفع\_الكائن(س، رصة)؛

صنف - حقيقي؛

{

؛

إذا (صنف - خطأ) انشأ\_طبقة(س،

س، ل+١، ق)؛

{

}

هذه الإجرائية تعمل كما يلي: بداية تعتبر فئة التعلم حاوية وكذلك يعتبر نسق الطبقات حاوية، ق، . عند دخول كائن جديد إلى فئة التعلم تقوم الدالة "حصل\_كائن" باكتسابه ويعتبر مجهول الطبقة أي غير مصنف ثم تقاس مسافة ماهاالاتوبيس بين متجه السمات، س، ومتجه سمات الوسط الممثل لمنطقة الطبقة، م، فإذا كانت هذه المسافة أقل من الحد الأصغر، ف،، صنف الكائن مع كائنات هذه الطبقة، ثم تجدد معلومات هذه الطبقة ممثلة في متجه الوسط، ومصنوفة الانتشار، ش، . ماعدا ذلك إذا كانت هذه المسافة أقل من الحد الأكبر، ف١، يعتبر الكائن الجديد غير محدد الطبقة مؤقتا ويدفع به إلى رصة؛ بالدالة "دفع\_الكائن"؛ عندما تمتلئ هذه الرصة تقوم الدالة "صنف\_المخزون" بتصنيف جميع كائناتها لأقرب جار من الطبقات الموجودة حين ذلك، وتفرغ الرصة في نفس الوقت؛ أما إذا لم يمكن تصنيف هذا الكائن أي أنه يوجد خارج منطقة الحد الأكبر لجميع الطبقات تنشأ طبقة جديدة نواتها هذا الكائن.

ونلاحظ هنا أن إجرائيتي التعرف والتعلم

اندمجتا في إجرائية واحدة تحقيقا للاستراتيجية

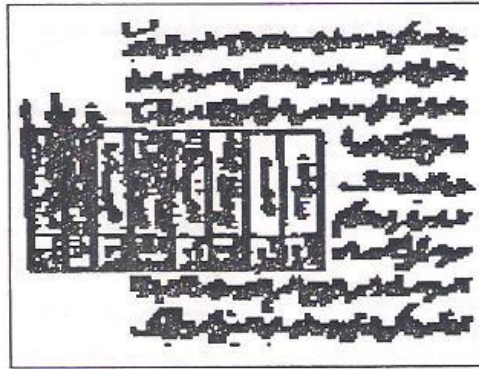
وقد استخدمت المنظومة المقترحة في التعرف على صفحات من مخطوطة كتاب الكاشي [١١]. وفي هذه المرحلة من التصميم والاختبار تعمل المنظومة بطريقة نصف آلية حيث يتم تشغيل مراحل تقسيم الصفحة إلى مناطق وسطور والتعرف على الحروف بعد تدخل المستخدم. وينحصر تدخل المستخدم هنا في اختيار واحدا من عدة اختيارات معروضة عليه تمثل حجم نافذة التطابق والحدود الصغرى والكبرى لتصنيف كائنات كل طبقة.

يمثل الشكل ٦ مراحل معالجة صورة إحدى عينات الصفحات المستخدمة في اختبار المنظومة. عند تشغيل المرحلة الأولى للتقسيم تخرج صورة الصفحة وبها مناطق الكتابة، والخلفية، والرسومات، كما هو موضح في الشكل ٦-أ. وفي المرحلة الثانية تقسم مجموعة الكتابة إلى سطور بناء على خاصية المكانية لينشأ عنها الصورة الممثلة في الشكل ٦-ب.

المقترحة في البداية وهي تعلم كلما مضيت". وهذه الاجرائية هي نواة البحث عن النسق الأمثل للأنماط الموجودة في الصفحة؛ أي أنها تعمل داخل حلقة مستمرة للبحث تبدأ بنسق افتراضى وتنتهى عند الوصول إلى النسق الأمثل والذي يحقق أصغر معدل للخطأ في التعرف على الكائنات الجديدة.

#### ٨ النتائج التجريبية

نفذت منظومة التعلم والتعرف المقترحة على حاسوب شخصى. استخدمت لغة الـ "سى++" في تنفيذ عناصر هذه المنظومة، كما اعتمدت في التصميم منهجية التصميم بالكائنات مما يسمح ببناء النواة الأولى بسرعة واعادة استخدام عناصرها في مراحل التطوير اللاحقة.



شكل ٦-أ صورة الصفحة بعد مرحلة التقسيم الأولى



Images", graphical models and image processing, vol 55, no.3, may, pp 203-217, 1993,

[5] Gary E Kopec, Philip A Chou, "Document Image Decoding Using Markov Source Models", IEEE transactions on Pattern Analysis and Machine Intelligence, vol 16, no 6, june 1994,

[6] Su Liang and M.Ahmadi, "A morphological Approach to Text String Extraction from Regular Periodic Overlapping Text/Background Images", graphical models and image processing, vol 56; no.5, september, pp 402-418, 1994,

[7] M. Y. Mahmoud El Nahas, Mohamed A. El Hamalaway & Ali Ali Fahmy, "A Statistical Approach for Arabic Character Recognition", 12th International Congress for Statistics, Computer Science, Social and Demographic Research, Cairo; pp. 243-250; March 28th-April 2nd, 1987,

[8] M.Y.Mahmoud (El Nahas) : " Object Decomposition by a Nondeterministic finite automaton ", Fifth ORMA Conference, MTC Cairo, November 23-25, 1993,

[9] Azriel Rosenfeild , Avinach C. Kak "Digital Picture Processin", Academic Press, London, 1982,

[10] Stephen P. Banks, "Signal Processing, Image Processing and Patter Recognition", Prentice Hall, London, 1990,

[11] جمشيد غياث الدين الكاشي ، مفتاح الحساب

والرسالة المحيطة ٣ ، دار الطبع والنشر للكتب

الفني والعلمي للدولة ، موسكو ١٩٥٦م.

ولقد حققت المنظومة معدلا مقبولا من

التعرف على الحروف يصل إلى ٧٧٪ . ويجرى

تطوير المنظومة حاليا بتوصيف نماذج أساسيات

الحروف بصورة أدق تسمح برفع معدل التعرف

على الحروف . ويهدف البحث مستقبلا لمعالجة

أنماط أخرى أكثر تعقيدا في المخطوطات كالجداول

مثلا:

#### إمتان

يتقدم الباحثان بجزيل الشكر للأستاذ الدكتور/ محمود

خشان لتوفيره صورة المخطوطة محل التطبيق.

#### المراجع

[1] R.O.Duda , P.E. Hart, " Pattern Classification and Scene Analysis", Wiley, New York, 1973,

[2] Fujihara, H.; Babiker, E.; Simmons, D.B., "Fuzzy Approach to Document Recognition" Second IEEE International Conference on Fuzzy Systems , p.980-5 vol.2,

[3] Jianying Hu; Pavlidis, T., "Hierarchical curve retrieval with application to aerial image alignment and cursive handwriting recognition", Proceedings. Second Annual Symposium on Document Analysis and Information Retrieval, p. 411-27,

[4] Mohamed Kamel and Aiguo Zhao, "Extravtion of Binary/Graphics Images from Grayscale Document



المجلس الأعلى للغة العربية



اللغة العربية  
في

تكنولوجيا المعلومات

تطور واعد... و تطوير متواصل



منشورات المجلس  
2005

## دراسة مقارنة بين أشكال الحروف العربية والحروف الإنجليزية

أ. د. محمد يونس الحملوى

أستاذ هندسة الحاسبات، كلية الهندسة، جامعة الأزهر، القاهرة  
5 شارع سليمان محمد أباطة، الحى الثانى، مصر الجديدة، القاهرة

هاتف 6321465 ناسوخ (فاكس) 6377446

بريد إلكترونى mhamalwy@hotmail.com

الكلمات المفتاحية: الحروف العربية، الحروف اللاتينية، الحروف الإنجليزية، التعرف الآلى على الحروف المطبوعة OCR، المعالجة الآلية للغة العربية.

لبحث يدخل ضمن المحور التالى من محاور المؤتمر كما فى نشرة باللغة الإنجليزية:

Natural Processing of the Arabic Language: e-tools, OCR,

C:\Documents and Settings\محمد الحملوى\My Documents\استندات الجزائر\المجلس الأعلى للغة العربية\الحروف العربية والإنجليزية.doc

### نبذة:

يدرس هذا البحث أساساً أشكال الحروف المكتوبة للغة العربية وللغة الإنجليزية. وقد تم استنباط سمات أساسية للأشكال المختلفة للحرف العربى الواحد بحيث يمكن الاستفادة منها فى عملية التعرف على الحروف. كما تمت دراسة الشكلين الأساسيين للحرف الإنجليزي الواحد، واتضح أنه لا يوجد للشكلين الأساسيين للحرف الإنجليزي الواحد فى المجمع سمة أساسية تربط الشكلين، بل أن بعض هذه الحروف لا توجد علاقة البتة بين شكلى الحرف المختلفين. ومن المأمول أن تساعد هذه الدراسة فى وضع برمجيات للتعرف الآلى على الحروف العربية المكتوبة بطريقة سهلة تؤدى إلى الحصول على نتائج جيدة. كما أظهر البحث منطقية استنباط الأشكال المختلفة للحرف العربى تبعاً لموقعه من الكلمة على عكس الحروف الإنجليزية التى لا يوجد فى الغالب رابط منطقى بين الشكل الكبير (المستعمل فى بداية أسماء الأعلام) والشكل الصغير (المستعمل فى وسط الكلمة) لها.

## 1. مقدمة:

يقوم الكثيرون بعمل محاولات للتعرف على الحروف المكتوبة في اللغات المختلفة وقد تم وضع برمجيات كثيرة في لغات مختلفة تقوم بهذه الوظيفة بدرجات نجاح متفاوتة، بل ولقد تم وضع برمجيات للتعرف على الخطوط المكتوبة باليد أيضاً لهذه اللغات؛ وان لم يزل الكثير منها في طور التطوير<sup>1</sup>. على أنه بالنسبة للغة العربية فمازالت نتائج تلك البرمجيات غير مرضية وتحتاج لعمل دؤوب للوصول إلى نتائج أفضل<sup>2-3-4-5</sup>.

وفي هذا الإطار فإن أول خطوة للتعرف على الحروف هي دراسة أشكالها وأنماط كتابتها المختلفة. بالنسبة للغة العربية توجد خطوط عديدة للكتابة وللطباعة منها النسخ، الرقعة، الكوفي، الثلث... وهي تشترك جميعها في سمات يسهل معها تمييز الحرف الواحد مهما كان نوع الخط المكتوب به. والدراسة الحالية لأشكال الحروف تيسر عملية التعرف الآلي عليها من خلال وضع قواعد منطقية تسير عليها عملية التعرف.

ويركز البحث على حروف اللغة العربية بالإضافة إلى حروف اللغة الإنجليزية والتي هي مجموعة جزئية من حروف اللغة اللاتينية<sup>6-7-8</sup>. وبالرغم من أن الحروف العربية والتي تأخذ أشكالاً متعددة باختلاف موقعها من الكلمة؛ مما كان يشكل صعوبة في عملية التعرف؛ إلا أنه قد أمكننا التوصل إلى أنه ذلك توجد سمات أساسية لكل حرف تجري عليه إضافات بشكل ما للحصول على الأشكال الأخرى لنفس الحرف على حسب موقعه من الكلمة. كما تمت دراسة علاقة شكلي الحرف الإنجليزي ببعضهما البعض بهدف الوصول إلى النسق الذي يربط بين هذين الشكلين.

<sup>1</sup> Ching Y. Suen; Advances in Optical Character Recognition; Canadian Computer Conference; Edmonton, Canada; May 23-25<sup>th</sup>, 1978.

<sup>2</sup> H. Y. Abdelazim and M. Y. Hashish; Arabic Reading Machine; 10<sup>th</sup> National Computer Conference; Riyadh, Saudi Arabia; 1988; pp. 733-744.

<sup>3</sup> M. A. Sharkawy et al; Fourier Descriptors for Printed Arabic Character Recognition; 13<sup>th</sup> International Conference on Statistics and Computer Science; Cairo, Egypt; March 26-31<sup>st</sup>, 1988.

<sup>4</sup> M. F. Tolba et al; A Recognition Algorithm for Arabic Printed Characters; 11<sup>th</sup> International Congress for Statistics and Computer Science; Cairo, Egypt; 1986

<sup>5</sup> M. Y. Mahmoud et al; A Statistical Approach for Arabic Character Recognition; Twelfth International Congress for Statistics and Computer Science; Cairo, Egypt; March 28<sup>th</sup> - April 2<sup>nd</sup>, 1987.

<sup>6</sup> ISO/IEC 10646 International Standard; ISO/IEC; Geneva, Switzerland; 1993.

<sup>7</sup> Mohamed A. El Hamalawy; AFGUST: A Standard for Coding Arabic Character Sets; 8<sup>th</sup> International Congress for Statistics and Computer Science; Cairo, Egypt; March 26-31<sup>st</sup>, 1983; also in Egyptian Computer Science Journal; Vol. 6, No. 1; Jan. 1983.

<sup>8</sup> ISO/IEC 8859 International Standard; ISO/IEC; Geneva, Switzerland; 1999.

## 2. السمات العامة لحروف اللغة العربية:

من المعروف أن حروف اللغة العربية تتخذ أشكالاً تختلف باختلاف مواقعها من الكلمة أي بحسب كونها في حالة إفراد أو في أول الكلمة أو وسطها أو نهايتها<sup>9</sup>. ويوضح جدول (1) هذه الأشكال لجميع الحروف العربية. في عملية مثل عملية التعرف الآلي على الحروف لا بد من الأخذ في الاعتبار هذا التغيير في الشكل<sup>10</sup>. على أنه بفحص هذه الأشكال يمكن التوصل إلى أن لكل حرف من الحروف العربية سمة أساسية يجرى التغيير فيها بالإضافة أو الحذف في أولها أو آخرها لإيجاد شكل مختلف فإذا أخذنا الحرف "ا" مثلاً نجد أن السمة الأساسية له جدول رقم 1.

الإضافة	السمة الأساسية للحرف	الحرف في نهاية الكلمة	الحرف في وسط الكلمة	الحرف في بداية الكلمة	الحرف مفرداً
—	ا	ا	ا	ا	ا
—	ب	ب	ب	ب	ب
—	ت	ت	ت	ت	ت
—	ث	ث	ث	ث	ث
/C	ج	ج	ج	ج	ج
/C	ح	ح	ح	ح	ح
/C	خ	خ	خ	خ	خ
	د	د	د	د	د
	ذ	ذ	ذ	ذ	ذ
	ر	ر	ر	ر	ر
ز	ز	ز	ز	ز	ز

<sup>9</sup> M. C. Vanwormhoudt and Mohamed A. El Hamalaway; Remarks about printing and Displaying some Non- Latin Characters; International Federation of Automatic Control (IFAC) Conference; Cairo, Egypt; Nov. 26-29<sup>th</sup>, 1977.

<sup>10</sup> Mohamed A. El Hamalaway and Salwa H. El Ramly; A Language Dependent Arabic Character Recognition Approach; 14<sup>th</sup> International Congress for Statistics and Computer Science; Cairo, Egypt; March 25-30<sup>th</sup>, 1989.

س	س	س	س	س	س
ش	ش	ش	ش	ش	ش
ص	ص	ص	ص	ص	ص
ض	ض	ض	ض	ض	ض
ط	ط	ط	ط	ط	ط
ظ	ظ	ظ	ظ	ظ	ظ
ع	ع	ع	ع	ع	ع
غ	غ	غ	غ	غ	غ
ف	ف	ف	ف	ف	ف
ق	ق	ق	ق	ق	ق
ك	ك	ك	ك	ك	ك
ل	ل	ل	ل	ل	ل
م	م	م	م	م	م
ن	ن	ن	ن	ن	ن
هـ	هـ	هـ	هـ	هـ	هـ
و	و	و	و	و	و
ي	ي	ي	ي	ي	ي

هي "ا" ويتم إضافة جزء أفقي "ـ" إذا ما وقع الحرف في وسط الكلمة. وبالنسبة للحرف "ب" تجد أن له سمة أساسية هي "بـ" دـ يضاف إليها جزء على شكل "ـا" في آخر الحرف وذلك في حالة إفراد الحرف أو وجوده في آخر الكلمة. وتطبيق هذه الملاحظة على الحروف تـ، ثـ، فـ. أما الحرف "جـ" فإن له سمة أساسية هي "جـ" يضاف إليها الجزء "ـا" في حالة إفراد الحرف أو وقوعه في نهاية الكلمة. وتطبق نفس الملاحظة على الحروف "حـ"، "خـ"، "عـ"، "غـ". ونلاحظ أنها نفس القاعدة المطبقة في المجموعة "بـ"، "تـ"، "ثـ"، "فـ" مع تغيير شكل الإضافة.

بالإضافة إلى ذلك فإنه يمكن الملاحظة أن الحروف "سـ"، "شـ"، "صـ"،

"ض"، "ق"، "ل"، "ن" لها سمات هي "س"، "ش"، "ص"، "ض"، "ة"، "ل" ويجري إضافة إما جزء على شكل "ـ" إذا كتب الحرف بصورة مفردة أو إذا كان في نهاية الكلمة، أو إضافة جزء أفقي "ـ" في حالة كتابة الحرف في أول الكلمة أو في آخرها.

كما يمكن ملاحظة أن الحرف "د"، "ذ"، "ر"، "ز"، "ط"، "ظ"، "م"، "و" تحفظ بشكل واحد أيا كان موقعها في الكلمة وهو نفس شكل السمة الأساسية لها.

يبقى ثلاثة حروف عربية ينفرد كل منها بخصائص ولا يشترك مع أي مجموعة أخرى وهي الحرف "ك" حيث له سمة أساسية "ـ" يتم إضافة جزء علوي "/" في حالة وجود الحرف في أول أو وسط الكلمة، أما إذا وقع الحرف في نهاية الكلمة أو جاء منفردا فتتم عليه إضافتين هما "ء"، "ـ". أما الحرف "هـ" وهي نفس شكل الحرف في حالة الإفراد، أما إذا كان الحرف في نهاية الكلمة فيتم إضافة "ـ" قبله، أما إذا جاء الحرف في أول أو في وسط الكلمة فيتم شكل "هـ" بالإضافة إلى المد "ـ" فيصبح الشكل "هـ". يبقى من الحروف العربية الثمانية والعشرين حرف واحد هو الحرف "ى" الذي لم يكن استنباط سمة أساسية له لكن يبقى الحرف على شكله "ى" في حالة وجوده منفردا أو وقوعه في آخر الكلمة، بينما يصير شكله "يـ" في حالة وروده في أول الكلمة أو في وسطها.

#### - السمات العامة لحروف اللغة الإنجليزية:

بتفحص الحروف الإنجليزية المعروضة في الجدول "2" نجد شكلين لكل حرف شكل صغير وشكل كبير. ونجد أن مجموعة من الحروف تحتفظ بنفس شكلها في المجموعتين مع تغيير جدول رقم "2".

## أشكال حروف اللغة الإنجليزية

الإضافة	السمة الأساسية	الحرف الصغير Small	الحرف الكبير Capital
		a	A
	b	b	B
	c	c	C
		d	D
	E	e	E
-	F	f	F
		g	G
	h	h	H
	I	i	I
	J	j	J
	K	k	K
-	l	l	L
		m	M
		n	N
	O	o	O
	P	p	P
		q	Q
		r	R
	S	s	S
	T	t	T
	U	u	U
	V	v	V
	W	w	W
	X	x	X
	Y	y	Y
	Z	z	Z

حجمها من صغير إلى كبير وهي المجموعة "O", "S", "U", "V", "W", "C", "Z", "X" ويمكن اتخاذ هذا الشكل للحرف على أنه السمة الأساسية له.

وعلى خلاف ذلك فتوجد مجموعة من الحروف لا يمكن أستنباط سمة أساسية لها حيث لا توجد علاقة بين الشكلين مثال ذلك مجموعة الحروف "Q", "R", "N", "M", "G", "D", "A" كما يمكن ملاحظة أن مجموعة من الحروف هي "K", "F"، جرى تحريف في شكلها الأصلي لتصبح "k"، "f"، وأخرى مثل الحرفين "T"، "E"، جرى لها تحريف في الشكل الأصلي مع إضافة. بالإضافة إلى ذلك يمكن ملاحظة أن الحرفين التاليين "Y"، "P" قد جرى تحريف في شكلهما الأصلي وإنزالهما للمستوى الأسفل. أما الحرف "J" فقد جرى له تحريف في الشكل الأصلي مع تحريك للمستوى لأسفل مع إضافة نقطة ليصبح "j"، وكذلك الحرف "T" فتم تصغيره مع إضافة نقطة ليصبح "t". كما توجد ثلاثة حروف يمكن استنباط سمة أساسية لها وهي الحرف "B" (سمة أساسية b مع إضافة جزء علوي لتصبح "B") والحرف H (سمة أساسية h مع إضافة جزء علوي لتصبح "H") والحرف "L" (سمة أساسية l مع إضافة "l" لتصبح "L"). وبهذا فإنه توجد سمات أساسية لأحد عشر حرفاً فقط من ستة وعشرين حرفاً هي مجموعة حروف اللغة الإنجليزية على عكس ما تبين بالنسبة لحروف اللغة العربية التي كان لجميع حروفها ما عدا حرفاً واحداً سمات أساسية.

ويمكن من الجدول "2" ملاحظة أن الحروف الإنجليزية لا يوجد في الغالب رابط منطقي بين شكلها الكبير (Capital) وشكلها الصغير (Small). كما يمكننا ملاحظة وجود لبس منطقي بين الشكل الكبير لحرف (أى "I" والشكل الصغير لحرف آخر هو حرف (إل) "l".

### 3- ملاحظات على صوتيات الحروف فى اللغة العربية مقارنة باللغة الإنجليزية

باعتبار الحروف المرسومة تعبر عن أصوات بعينها لذلك من المفيد إثبات بعض الملاحظات على أصوات حروف اللغة العربية وحروف اللغة الإنجليزية كى تتضح صورة ما نحن بصدد رسمه. ونلاحظ بالتالى تميز الحروف العربية بأن لكل منها صوتاً واحداً لا يتغير بموقع الحرف فى الكلمة وبتتابع الحروف قبله وبعده. كما توجد حركات محددة يمكن عن طريقها تنعيم الحرف وهى الفتحة والكسرة والضمة والشدة والسكون، بالإضافة إلى أن التنوين يلعب دوراً فى تنعيم الكلمة. وعلى خلاف ذلك فإننا نجد أن اللغة الإنجليزية لا يوجد لكل حرف صوت واحد، بل أكثر من صوت مثل الحرف "S" والذي ينطق بطرق متعددة على حسب ما قبله

وما بعده من حروف. وينطبق نفس الشيء على حرف "C" الذي ينطق س، ش، ك على حسب الكلمة بدون قواعد محددة لهذا.

أضف إلى ذلك أنه توجد حروف لا تنطق ولا تتبع قاعدة محددة لذلك كما في حرف "G". كما يلاحظ أن حرفاً مثل حرف "U" له ثلاثة عشر صوتاً ليس لها قاعدة منطقية في توليدها.

ولم يتعرض البحث لرسم الأصوات المختلفة في كلتي اللغتين انطلاقاً من تلك الملاحظات اللغوية الصوتية رغم أنها تستحق البحث إلا أننا اكتفينا برسم الحروف المكتوبة الموثقة في كلتي اللغتين مع تسجيل أن هذه النقطة تستحق بحثاً مستقبلياً.

#### 4- الخلاصة:

يقدم هذا البحث دراسة على أشكال الحروف المكتوبة للغة العربية ولغة الإنجليزية. وقد تم استنباط سمات أساسية للأشكال المختلفة للحرف العربي الواحد بحيث يمكن الاستفادة منها في عملية التعرف على الحروف. ومن المأمول أن تساعد هذه الدراسة في وضع برمجيات للتعرف الآلي على الحروف العربية المكتوبة بطريقة سهلة تؤدي إلى الحصول على نتائج جيدة.

تشترك مختلف خطوط اللغة العربية للكتابة في سمات مشتركة يسهل معها تمييز الحرف الواحد مهما كان نوع الخط المكتوب به. والدراسة الحالية لأشكال الحروف تضع قواعد منطقية تسيّر عليها عملية التعرف. ولقد تم توليد السمة الأساسية لكل حرف عربي التي يجرى عليها إضافات بشكل ما للحصول على الأشكال الأخرى لنفس الحرف على حسب موقعه من الكلمة. وبالتالي فإن عملية التعرف تفضى إلى التعرف على تلك السمات الأساسية للحرف وليس على شكله، أخذين في الاعتبار وجود قواعد لتكوين الكلمة العربية لا يمكن الحيود عنها.

كما أظهر البحث منطقية استنباط الأشكال المختلفة للحرف العربي تبعاً لموقعه من الكلمة على عكس الحروف الإنجليزية التي لا يوجد في الغالب رابط منطقي بين الشكل الكبير (المستعمل في بداية أسماء الأعلام) والشكل الصغير (المستعمل في وسط الكلمة) لها.

كما اتضح من دراسة الشكلين الأساسيين للحرف الإنجليزي الواحد، أنه لا يوجد للشكلين الأساسيين للحرف الإنجليزي الواحد في المجمل سمة أساسية تربط الشكلين، بل أن بعض هذه الحروف لا توجد علاقة البتة بين شكلي الحرف المختلفين. كما يتضح وجود لبس منطقي بين الشكل الكبير لحرف (آي) 'i' والشكل الصغير لحرف آخر هو حرف (إل) 'l'.

5- عرفان:

يشكر الباحث الأستاذة الدكتورة سلوى الرملى على مساعدتها القيمة فى إنجاز هذا البحث.



الجمعية المصرية لهندسة اللغة

المؤتمر الخامس لهندسة اللغة

مجلد الأوراق البحثية

١٤-١٥ سبتمبر ٢٠٠٥م  
القاهرة - جمهورية مصر العربية

## التعرف على حروف اللغة العربية باستخدام خوارزمات التحريف والتقطيع

أ.د/محمد يونس الحملأوى أ.د/محمد يسرى النحاس أ.د/إسماعيل عبد الغفار إسماعيل م/محمد لبيب يوسف  
كلية الهندسة-جامعة الأزهر كلية الهندسة-جامعة الأزهر الكلية الفنية العسكرية إدارة نظم المعلومات

### ملخص الورقة :

تقدم في هذه الورقة منظومة تجريبية للتعرف على حروف اللغة العربية. الإطار المستخدم في تصميم هذه المنظومة هو المنهج التركيبي للتعرف على الحروف. هذا المنهج يتعامل مع الحروف على أساس أنها أشكال هندسية مكونة من مقاطع أساسية وإضافية لها سمات قياسية وعلاقات مكانية مميزة لكل مقطع وبنية تركيبية لكل حرف. البرنامج المقترح للتعرف على الحروف العربية يقوم بمعالجة أولية للنص المأخوذ من الماسح الضوئي ثم يقوم بعملية تحريف لهذه الحروف و بالتالي تحويلها إلى سلسلة متصلة من النقاط ثم تصنيفها إلى مجموعات من المقاطع، و بتحليل العلاقات والمسافات بين نقاط المقاطع يتم عملية تمييزها ثم تركيبها للتعرف على الحروف. في هذا الإطار تم تطوير خوارزم التحريف "المحور الأوسط" وخوارزم تقطيع الكتانية باستخدام الوصف التركيبي لحروف اللغة العربية. تشير النتائج التجريبية الأولى إلى تحقيق دقة عالية في تحريف خطوط الكتابة وتقطيعها و معدل جيد من التعرف على حروف اللغة العربية.

### الكلمات الدليلية :

التعرف على الحروف -المنهج التركيبي-خوارزم التحريف - المحور الأوسط-تقطيع الكتابة- القارئ الآلى- البرمجة الشبئية.

### ١-مقدمة:

يعتبر تاريخ البحث في مجال التعرف على الحروف العربية قصير نسبياً مقارنة بنظيره في التعرف على الحروف اللاتينية. إذ سجلت المحاولات الأولى للتعرف على الحروف العربية في بداية الثمانيات [١،٢]. قدم الباحثون (أمين وآخرون) طرق تفاعلية للتعرف على الحروف العربية [٢]. واستخدم الباحثون (المعلم وآخرون) المدخل التركيبي لتقطيع الحروف إلى مقاطع والتعرف عليها باستخدام السمات الهندسية والمكانية للمقاطع [٢]. وقدم الباحثون (التكريتى وآخرون) منطق الغموض في التعرف بعض الحروف العربية المكتوبة باليد [٢]. واستخدم (الشيخ) تحويلات فورييه في توصيف حدود الحروف العربية واستخدم (ماجنت) عزوم زرنيك في توصيف الحروف العربية [٣].

وبالرغم من هذه الجهود فما زال التعرف على حروف اللغة العربية يمثل صعوبة نسبية والمنظومات التي نجحت ما زالت محدودة التطبيق لنوع معين من الكتابات العربية. وأهم الصعوبات التي تواجه البحث في التعرف على الحروف العربية هي:

- أن حروف الكتابة العربية متصلة سواء المطبوع منها أو المكتوب باليد
- أن كثيراً من هذه الحروف يحتوى على مكونات مركبة من النقاط والهمزات وعلامات التشكيل.
- أن أشكال الحروف العربية تختلف باختلاف موقعها من السطر.
- كثرة ونوع فظطات وأساليب الكتابة العربية. فقد قدرت عدد الفظطات المستخدمة في الكتب بـ ٥٠ ؛ ، وعدد الفظطات المستخدمة في الصحف بـ ١٥٠ [٢].
- ويمكن تصنيف المنهجيات المتبعة في التعرف على حروف اللغة العربية الطبوعة أو المكتوبة باليد إلى:
- المنهج التركيبي: وهي تتميز بالدقة والمرونة في التعرف على الفظطات المختلفة للحروف ولكنها تعتمد إلى حد كبير على التحليل الدقيق للغة والتوظيف الجيد لطرق التطبيق
- المنهج الإحصائي: وهي تتميز بالدقة والمرونة في التعرف على الحروف في المضامين المختلفة وأشهرها الطرق المبنية على سلاسل ماركوف الخفية التي أثبتت نجاحات أخرى في مجال التعرف على الكلام.
- الشبكات العصبية: و هي تتميز بسهولة بالتنفيذ وخاصة بالمعالجات المتوازية وإن كانت تحتاج إلى إمكانيات حسابية فائقة ووقتاً طويلاً للتعلم
- مناهج توفيقية بين أكثر من منهج مما سبق ذكرها

ونحن في هذه الورقة نقدم منظومة تجريبية للتعرف على الحروف المطبوعة للغة العربية يمكن تطبيقها في مجالات متعددة من تقنيات المعرفة. صممت هذه المنظومة في إطار المنهج التركيبي للتعرف على الحروف بحيث لا يعتمد

معدل التعرف على الحروف على نوع أو حجم الفنت المستخدم. والورقة البحثية منسقة كالاتي. في المقطع الثاني نعرض الإطار المنهجي للتعرف على حروف اللغة العربية. في المقطع الثالث نقدم الإجراءات التفصيلية لبرنامج التعرف على الحروف العربية. في المقطع الرابع نعرض النتائج التجريبية الأولية وتقييم معدلات التعرف في الحالات المختلفة. والمقطع الخامس يلخص نتائج البحث واتجاهات البحث المستقبلية.

## ٢- الإطار المنهجي للتعرف على الحروف:

الإطار المستخدم في تصميم منظومة التعرف على الحروف العربية هو المنهج التركيبي للتعرف على الحروف وهو ينقسم إلى ثلاث مراحل هي:

المرحلة الأولى: مرحلة ما قبل المعالجة:

في هذه المرحلة تتم العمليات الآتية :

١. قراءة ملف النص المأخوذ من الماسح الضوئي
٢. تنقية بيانات النص بمعنى إزالة النقاط السوداء الزائدة التي تنتج أثناء عملية المسح الضوئي وذلك باستخدام مرشح "ملح وفلفل" على سبيل المثال.
٣. تحويل الصورة إلى التمثيل الثنائي: الغرض من هذه الخطوة تخصيص أحد اللونين فقط لكل نقطة من النقط المضنية التي يتكون منها الملف ، ويتم ذلك بتحديد القيمة المتوسطة لشدة إضاءة جميع النقط ويخصص اللون الأسود لجميع النقط التي يزيد وزنها عن هذه القيمة المتوسطة ، ويخصص اللون الأبيض لباقي النقط.
٤. إزالة الحزرات التي تظهر أحيانا على أطراف الحروف على هيئة دوائر صغيرة جدا يمكن التعرف عليها بأنها نقطة واحدة بيضاء محاطة من جميع الجهات بنقط سوداء وسيستخدم لذلك مرشح "الحزرات".

المرحلة الثانية: مرحلة المعالجة والتحليل:

في هذه المرحلة تتم العمليات الآتية :

١. تحجيف النص: تطبيق خوارزم التحجيف على تجمعات النقط السوداء التي يتكون منها النص فينتج عن ذلك تقليص سمك هذه التجمعات ليصبح كل منها عبارة عن سلسلة متلاصقة من النقط السوداء لايزيد سمكها عن نقطة واحدة
٢. تحديد مواصفات السطر-المحور الرئيسي: تحديد مواصفات السطر وأهمها السمك ويقصد به الفرق بين أكبر وأصغر إحداثي أفقي لنقطتين لهم نفس الإحداثي الرأسي ، وبفرض أن جميع سطور النص مكتوبة بحجم ثابت فإنه يمكن إدراك حدود السطور التالية حتى وإن كانت لاتتضم الحد الأدنى من سلاسل النقط المتلاصقة
٣. تقطيع النص: تمييز كل سلسلة من النقط المتلاصقة واعتبارها شيء منفصل بذاته ، وكل واحدة من هذه السلاسل قد تعبر عن حرف منفصل أو مقطع من كلمة يضم حرفين أو أكثر أو نقطة منعزلة (وهي التي تظهر أعلى أو أسفل الحروف ب ، ج ، خ ، ن ، ض ، ط ، غ ، ف ، د ، ز) أو نقطتين متلاصقتين (وهي التي تظهر أعلى أو أسفل الحروف ت ، ي) أو ثلاثة نقط متلاصقة (وهي التي تظهر على الحروف ث ، ش ،).

المرحلة الثالثة: مرحلة التعرف:

في هذه المرحلة تتم العمليات الآتية :

١. تقدير سمات المقاطع: بالنسبة لكل مقطع يتم تحديد المحور الرئيسي، والمستطيل المحدد لكل تقاطع سلسلة المقطع، و مركز المقطع أي النقطة المتوسطة في هذه السلسلة و الجزء الواقع منه على المحور الرئيسي المذكور أعلاه ، ويجب الأيقل عرض هذا الجزء عن حد معين ، وبمعلومية هذه الأجزاء لجميع مقاطع السطر الواحد يتم تحديد المحور الرئيسي له ، وفي حالة تعذر ذلك يستفاد من معرفة سمك السطور السابقة.
٢. التعرف على الحروف: والفكرة الأساسية لأسلوب تمييز الحروف هي أن الأشكال الهندسية للحروف المتصلة في اللغة العربية عبارة عن تراكيب يجمعها محور أفقي ، يلتصق بعضها بهذا المحور من أعلاه وبعضها من أسفله ، ويتميز كل حرف بالطول النسبي لهذه التراكيب واتجاه ميلها على هذا المحور ، وكذا موقع الحرف في المقطع سواء في بدايته أومنتصفه أو آخره ، أو أن يأتي الحرف منعزلا بمفرده يتم التعرف على الحروف المحتوية لكل مقطع من خلال مواصفاته الهندسية والمكانية التي يمكن إيجازها في المحددات الآتية:

١. علاقة النقطة الطرفية للحرف بنقطة التلاقي مع المحور الرئيسي (مقارنة الإحداثيات)
٢. تحديد النبرات (حروف ب ، ت ، س ، ن ، ..... ) وعصا حرف الطاء والظاء بمقارنة طولها بحرف الألف.
٣. تحديد الدوائر المغلفة ( حروف ص ، ض ، ط ، ظ ، ف ، ق ، م ، هـ ، و ) .
٤. بناء على هذه المحددات يمكن تحديد المجموعة التي ينتمي إليها هذا الحرف ( ب ، ت ، ث ) أو ( ج ، ح ، خ ) وهكذا ، وتميز هذه المجموعات بأكواد خاصة.
٥. يستكمل التصنيف طبقا للقواعد الآتية:

- ربط النقط المنعزلة بالحروف التي تقع في الجوار فتتحول النبرة إلى ب أو ن ، والدائرة المغلقة إلى ف أو ض ، وحرف ح إلى ج أو خ.
  - ترقية بعض الحروف مثل ن لتصبح ت أو ب لتصبح يد أو س لتصبح ش.
  - ربط الهمزات بالحروف التي تقع في الجوار لتتحول ا إلى أ والنبرة إلى نـ والواو إلى و وهكذا.
- والقواعد التي نطبقها في توصيف حروف اللغة العربية محددة في الجدول (١) والجدول (٢).
- جدول (١) توصيف حروف اللغة العربية: الحروف البسيطة

الحرف	موقعه في المقطع	مكوناته	علاقة قير بالبحور	علاقة قير بالبحور الإحداثي السيني
ا	منفصل		أعلى	-
	أخره			
ب	منفصل	نبره	أعلى	-
	أوله			
	منتصفه			
	أخره			
ج	أوله		أعلى	>
	منتصفه			
د	منفصل	نبره	أعلى	-
	أخره			
ر	منفصل		أسفل	>
	أخره			
س	أوله	ثلاث نبرات	أعلى	-
	منتصفه			
ص	أوله	دائرة مغلقة	ليس لها نقطة نهاية	لها تقطع
	منتصفه			
ط	منفصل	دائرة مغلقة	ليس لها نقطة نهاية	لها تقطع
	أوله			
	منتصفه			
	أخره			
ع	أوله	نبره	أعلى	<
	منتصفه			
ف	منفصل	دائرة مغلقة	ليس لها نقطة نهاية	لها تقطع
	أخره			
ق	أوله	دائرة مغلقة	ليس لها نقطة نهاية	لها تقطع
	منتصفه			
ك	منفصل	ا	أعلى	-
	أخره			
	أوله			
	منتصفه			
ل	منفصل	ا	أعلى	-
	أخره			
	أوله			
	منتصفه			
م	أوله	دائرة مغلقة	ليس لها نقطة نهاية	لها تقطع
	منتصفه			
ن	منفصل	نبره	أعلى	-
	أوله			
	منتصفه			
هـ	منفصل	دائرة مغلقة	ليس لها نقطة نهاية	لها تقطع
	أخره			
ي	أوله	نبره	أعلى	-
	منتصفه			

جدول (٢) توصيف حروف اللغة العربية: الحروف المركبة

الحرف	موقعه في المقطع		الجزء العلوي		الجزء السفلي	
	ع	ح	علاقة ق، بالمحور	علاقة ق، الإحداثي السيني	علاقة ق، بالمحور	علاقة ق، الإحداثي السيني
ج ح خ	ح	منفصل أخره	أعلى	>	أسفل	<
س ش	س ش	منفصل أخره	أعلى	=	أسفل	>
ص ض	ص ض	منفصل أخره	ليس لها نقطة نهاية	لها نقطتي تقاطع	أسفل	>
ع غ	ع غ	منفصل أخره	أعلى أعلى	< =	أسفل	<
ق	ق ق	منفصل أخره	ليس لها نقطة نهاية	لها نقطتي تقاطع	أسفل	>
ل	ل	منفصل	أعلى	أ		
	ل	أخره	أعلى	نيره		
	ل	أوله	أعلى	أ		
	ل	منتصفه	أعلى			
م	م م	منفصل أخره	دائرة مغلقة		أسفل	=
ن	ن	أخره	أعلى	=	أسفل	>
هـ	هـ	أوله	دائرتان مغلقتان			
	هـ	منتصفه	دائرة مغلقة			دائرة مغلقة
و	و	منفصل	دائرة مغلقة		أسفل	>
	و	أوله				
	و	أخره				
ي	ي	منفصل	نيره		أسفل	>
	ي	أخره			أسفل	>

وسنقدم فيما يلي من البحث شرحاً للطرق المستخدمة في المرحلتين الثانية والثالثة

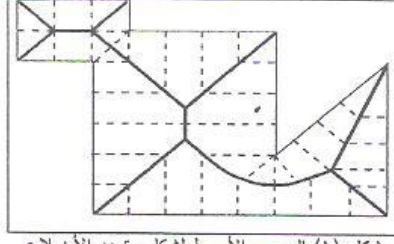
## ٢-١-٢ تحييف النص:

عملية تحييف النص تعني تحويل خطوط الكتابة السميكة إلى خطوط ومنحنيات سمك كل منها نقطة واحدة. وتعتبر عملية التحييف من العمليات الأساسية في طرق التعرف على الحروف باستخدام سماتها الشكلية الثابتة مع تغير فظف الكتابة أو درجة ميل سطور الكتابة. والمهم بمعرفة هذه الطرق سيجد في البحث المنشور [٥] مسح مرجعي وافى لطرق تحييف الحروف. وتختصر طرق التحييف في: طرق تتابعية [٥،٦،٧،٨]، و طرق متوازية [٩]. وتعتمد الطرق التتابعية على استخدام القوالب والنوافذ في تحييف وتقليم الأشكال. وتعمل الطرق المتوازية باستخدام نفس أساسيات التحييف السابقة ولكنها تطبق على مراحل متوالية وبالتالي فيمكنها استخدام المعالجات المتوازية مما يزيد من سرعة المعالجة. وكذلك توجد طرق تعتمد على الشبكات العصبية [١٠]. وما زالت كفاءة أداء هذه الطرق ودقتها وسرعتها محل بحث في مجالات تحليل صور الوثائق. والكثير منها يتأثر بدوران الأشكال بالنسبة لمحاور الصفحة. واقترحت مؤخراً طرق لا تعتمد على دوران أشكال الحروف [١١].

في هذا البحث نقترح استخدام خوارزم "المحور الأوسط" لتحيف أشكال حروف النصوص العربية. فقد وجد في دراسة سابقة عن طرق توليد المسارات لحركة الروبوت أن طريقة المحور الأوسط لها سمة أساسية وهي البحث عن المسار الوسيط وتحقق دقة عالية في إيجاد هذا المسار أياً كانت درجة دوران الأشكال في المستوى أو مدى ضيق المسافة بين العوائق مما يجعلها مناسبة لتحيف الحروف العربية التي تتميز بتعقيد أشكالها.

## خوارزم المحور الأوسط:

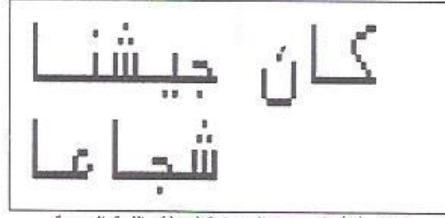
يعرف المحور الأوسط لشكل متعدد الأضلاع بأنه المحل الهندسي لجميع النقط التي يتساوى بعدها عن جميع أضلاع هذا الشكل [١٢]، ويوضح الشكل (١) مثال لشكل متعدد الأضلاع و المحور الأوسط له، ويلاحظ أن هذا المحور الأوسط يمثل "تحيف" لهذا الشكل وإذا نستطيع أن نطلق عليه خوارزم التحييف.



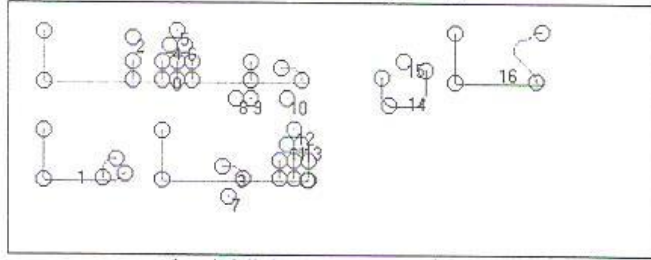
شكل (١) المحور الأوسط لشكل متعدد الأضلاع

يمكن استخدام "خوارزم المسافات الإقليدية" أو "خوارزم الهيكلية المورفولوجية" في إيجاد المحور الأوسط للأشكال الهندسية. وهي عمليات متكافئة في النتيجة. وفي هذا البحث تم تطوير "خوارزم تحريف مورفولوجي" يعمل تكرارياً على حذف نقاط حدود الشكل مستخدماً نافذة مورفولوجية من  $3 \times 3$  عناصر حتى يتوقف التغيير في هيكل الشكل.

ويوضح الشكل (٢) أن حروف اللغة العربية بعد أن نحصل على الصورة الضوئية لها تكون عبارة عن تجمعات غير منتظمة من الأشكال ، ولذا فإنه يمكننا أن نتعامل معها بخوارزم التحريف بحيث تتحول إلى سلسلة متصلة من النقاط كما يبين الشكل (٣)



شكل (٢) الصورة الضوئية لجملة باللغة العربية



شكل (٣) ناتج التحريف لجملة باللغة العربية

## ٢-٢ تقطيع النص:

يعتبر تقطيع الحروف من العمليات الهامة والأكثر صعوبة في التعرف على حروف اللغة العربية. والغرض من عملية التقطيع تقسيم كلمات الكتابة إلى أشكال مميزة يمكن استخدامها في التعرف على الحروف والكلمات. وعملية التقطيع تتلخص في الخطوات العامة الآتية:

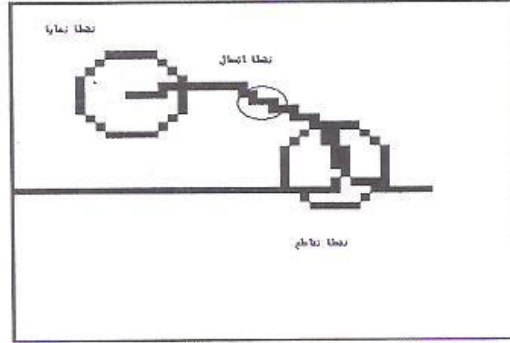
١. بدءاً من نقطة في صورة الصفحة
٢. أوجد مكان الشكل التالي
٣. استخراج السمات المميزة لصورة الشكل الحالي
٤. أوجد الرمز الأكثر تماثلاً من مجموعة الرموز مع سمات الشكل الحالي
٥. استمر حتى نهاية صور الأشكال في الصفحة

و يوجد ثلاث استراتيجيات متبعة في تقطيع النصوص [١٣] وهي:

١. الطريقة التقليدية والتي تقطع الصورة إلى مكونات مميزة بناءً على خواص الحروف
٢. التقطيع بالتعرف على أشكال الحروف
٣. الطريقة الكلية والتي تبحث عن الكلمات ككل

ونحن في هذا البحث نهج الطريقة التقليدية مع استخدام الخواص المميزة لمقاطع الحروف العربية في عمليات تقطيع صورة النص العربي. وتتخلص هذه الطريقة في الخطوات التالية:

١. تصنف النقط المكونة لكل مقطع إلى أربعة أصناف كما هو موضح بالشكل (٤)
  - نقطة منعزلة : وهي تلك التي ليس لها أى اتصال بنقط أخرى
  - نقطة نهاية : وهي تلك التي لها اتصال بنقط أخرى في اتجاه واحد فقط قـ هـ
  - نقطة اتصال : وهي تلك التي لها اتصال بنقط أخرى في اتجاهين.
  - نقطة تقاطع : وهي تلك التي لها اتصال بنقط أخرى في أكثر من اتجاهين
٢. بمقارنة الإحداثيات الرأسية لنقط الإتصال في كل مجموعات النقط ، يتم تحديد المحور الرئيسي ثم فصله
٣. نتيجة عملية الفصل تبقى أجزاء من مجموعات النقط المتلاصقة السابق تحديدها كأشياء منفصلة كما يوضح الشكل (٥).
٤. بدءاً من كل نقطة من نقط النهاية يتم تتبع النقط الملاصقة لها إلى أن نصل إلى آخر نقطة كانت تربط هذه السلسلة بالمحور الرئيسي قـ رـ .
٥. ترقيم كل سلسلة من النقط المتلاصقة لترتيبها وتصنيفها فيما بعد.



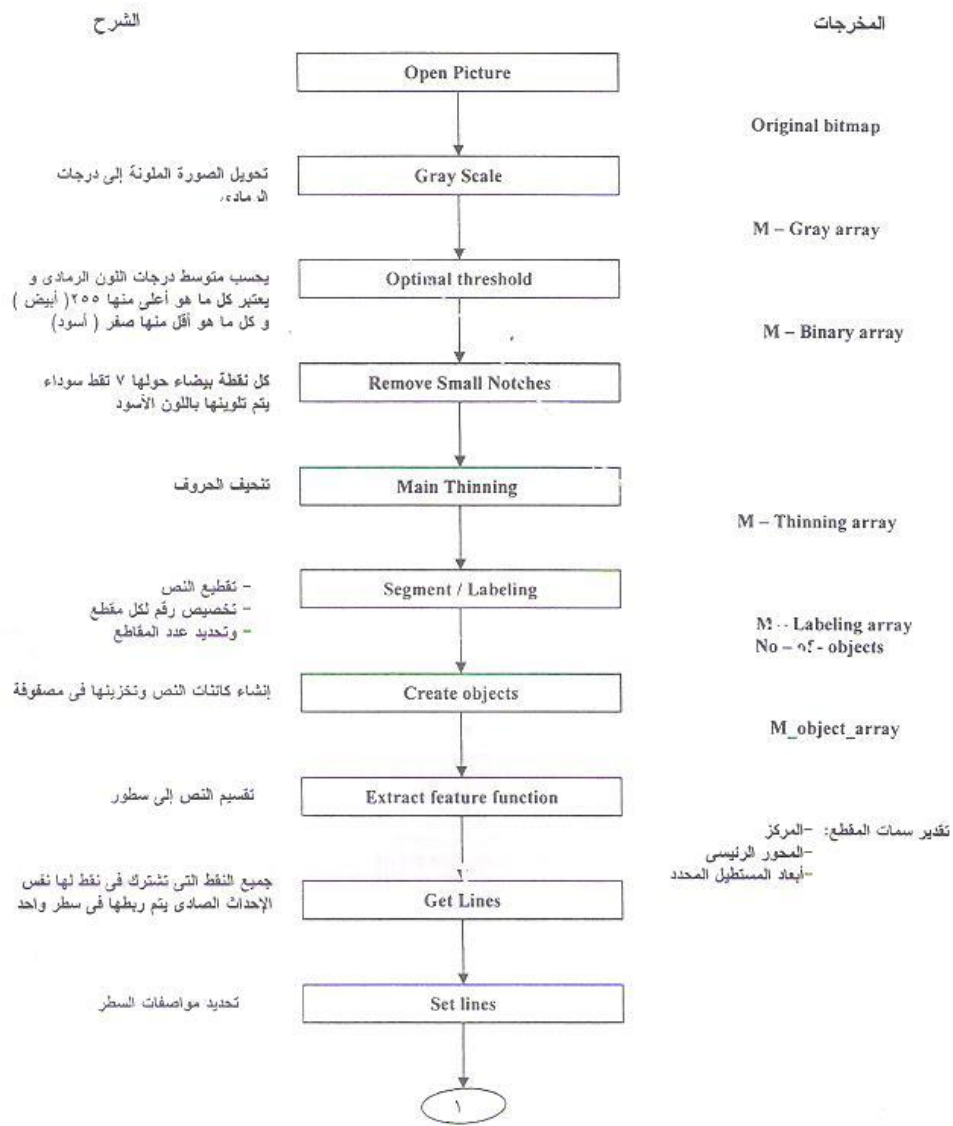
شكل (٤) أنواع النقط التي يتكون منها الحرف بعد التحريف



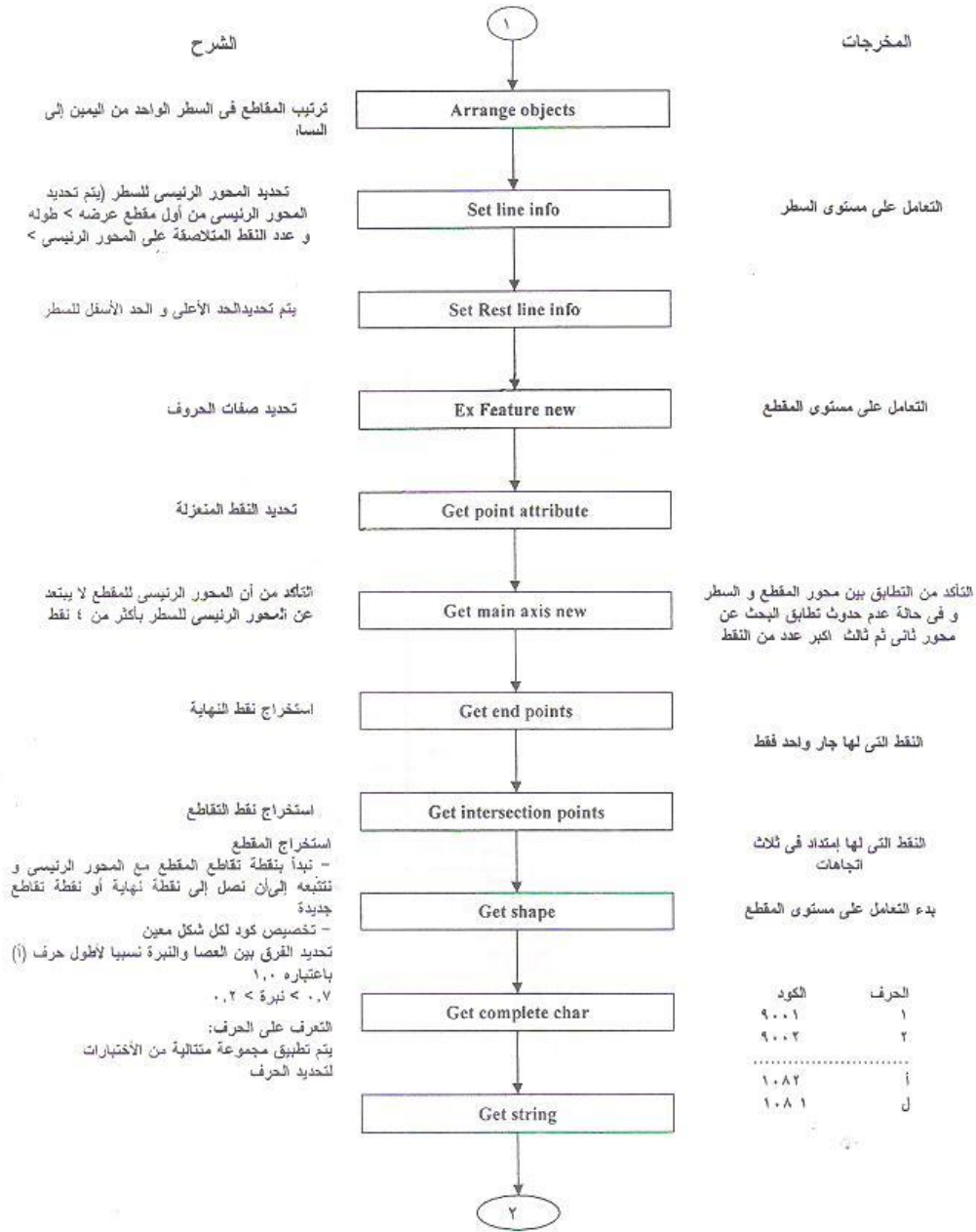
شكل (٥) شكل الجملة بعد التحريف وفصل المحور الرئيسي

### ٣-برنامج التعرف على حروف اللغة العربية :

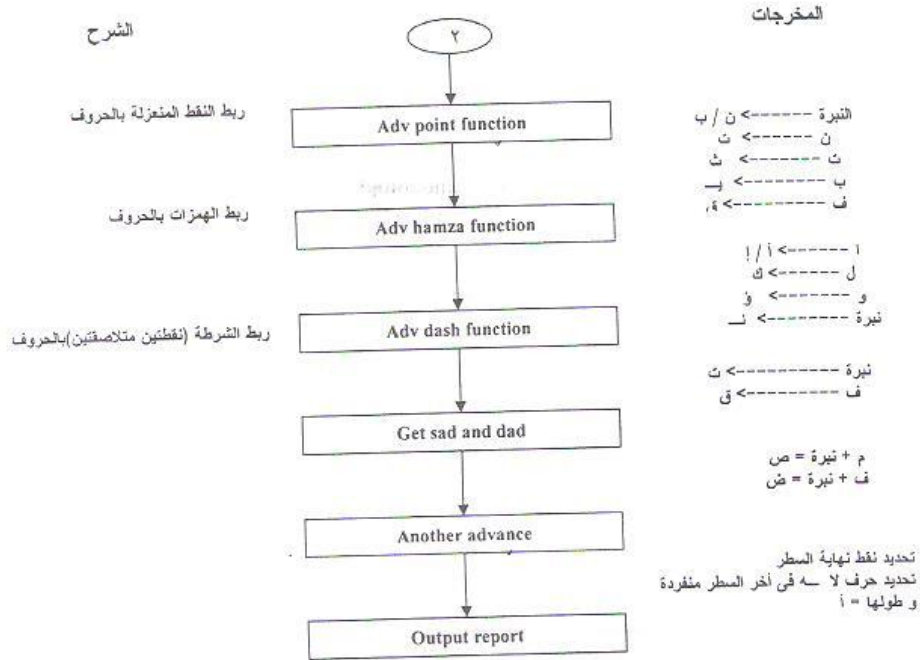
- بعد قراءة صورة الصفحة ومعالجتها يعمل برنامج التعرف طبقاً للتسلسل التالي :
١. تحريف النص باستخدام خوارزم التحريف ليتحول النص إلى سلاسل من النقط المتلاصقة
  ٢. تمييز كل سلسلة واعتبارها "كائن" منفصل قائم بذاته .
  ٣. تقطيع كائنات السلاسل الناتجة من عملية التحريف إلى مقاطع فرعية باستخدام خوارزم التقطيع
  ٤. استخراج سمات كل مقطع وهي:
    - المحور الرئيسي
    - المستطيل المحدد لكل نقاط سلسلة المقطع
    - مركز المقطع أي النقطة المتوسطة في هذه السلسلة قـم
    - ترتيب المقاطع في السطر الواحد من اليمين إلى اليسار
    - ٦. يصنف كل مقطع إلى حرف من الحروف طبقاً للقواعد الآتية:
      - المقطع ينتمي إلى المجموعة الأولى من الحروف البسيطة و يوضحها جدول (١) إذا كان يقع بأكملها أعلى أو أسفل المحور الرئيسي
      - المقطع ينتمي إلى المجموعة الثانية من الحروف المركبة و يوضحها جدول (٢) إذا كان يقع جزء منه أعلى المحور الرئيسي و جزء منه أسفل المحور الرئيسي
      - بناء على اتجاه النقطة قـم بالنسبة للمحور الرئيسي و بمقارنة الإحداثي السيني للنقطتين قـم ، قـم و كذا علاقة كل مقطع بما قبله و ما بعده يتم التعرف على الحرف طبقاً للتوصيف الموضح بالجدول (١) و الجدول (٢).
- وشرح الإجراءات التفصيلية لهذا البرنامج ومخرجات كل إجراءاته موضح في مخطط التدفق بالشكل (٦)



شكل (6) مخطط تدفق البيانات لبرنامج التعرف على حروف اللغة العربية



تابع شكل (١) مخطط تدفق البيانات لبرنامج التعرف على حروف اللغة العربية



تابع شكل (٦) مخطط تدفق البيانات لبرنامج التعرف على حروف اللغة العربية

#### ٤- النتائج العملية:

للتحقق من صحة المنهج المقترح في التعرف على الحروف العربية واختبار أداء البرنامج تم اختيار عينات مختلفة مختلفة للنصوص العربية المطبوعة وتحتوي جميع التراكيب الممكنة للحروف العربية داخل الكلمات وكتابتها بثلاث فنطحات شائعة الاستخدام (Simplified Arabic-Arabic Transparent-Times New Roman) ثم إدخالها إلى الحاسب باستخدام الماسح الضوئي. وقد استخدمت هذه العينات في مرحلة التعلم وكانت نتائجها كما في الجدول (٣).

الفقرة	عدد الحروف	الحروف المقروءة	نسبة التمييز
١	٣٢١	١٨٦	%٥٨
٢	٢٧٧	٢١٣	%٧٦
٣	٢٩٧	١٠٧	%٦٩
٤	٢٢٩	١٢٥	%٥٤
٥	١٦٥	١٢٠	%٧٢
٦	١٣٠	١٠٤	%٨٠
٧	١٦١	١٣٠	%٨٠
الإجمالي	١٥٨٠	١٠٨٥	%٦٨

وفي مرحلة تقييم أداء البرنامج اختيرت عينات للنصوص المستقطعة من الصحف اليومية والكتب المدرسية. والشكل (٧) به صورة نص مستقطع من إحدى الصحف اليومية المصرية. والشكل (٨) به نتيجة لتحريف وتقطيع هذا النص. ونتيجة التعرف موضحة في الشكل (٩).

سجلت أسعار الجملة للسلع الزراعية انخفاضا ملحوظا خلال شهر مايو الماضي لهذا العام مقارنة بالأسعار في نفس الفترة من العام الماضي بلغت نسبته ١٢,٢% أشار الي ذلك تقرير الأرقام القياسية لأسعار الجملة الذي يصدره الجهاز المركزي للتعينة العامة والاحصاء بصفة شهرية.

وأشار التقرير الي تراجع أسعار مجموعة من المنتجات من أهمها الخضراوات والفواكه والبيض وغيرها حيث انخفضت اسعار الخضراوات بنسبة ٨١,١% والفواكه بنسبة ٦,٨% والبيض بنسبة ١١,٣% .

وصرح اللواء أبو بكر الجندي رئيس الجهاز المركزي للتعينة العامة والاحصاء بأن أسعار بعض السلع الاخرى في اسواق الجملة سجلت ارتفاعا قدره ٤,٢% خلال شهر مايو الماضي مقارنة بالاسعار في نفس الفترة من العام الماضي

مثل المواد الغذائية والوقود ومنتجاته فقد ارتفعت اسعار الوقود بنسبة ١٣,٦% والمواد الخام بنسبة ١,٨% والسلع الوسيطة بنسبة ٦,٨% والسلع الاستهلاكية غير المعمرة بنسبة ٣,١% والسلع الاستهلاكية المعمرة بنسبة ٣% والسلع الاستثمارية

شكل (٧) صورة لنص المستقطع من إحدى الصحف اليومية

سجلت أسعار الجملة للسلع الزراعية انخفاضا ملحوظا خلال شهر مايو الماضي لهذا العام مقارنة بالأسعار في نفس الفترة من العام الماضي بلغت نسبته ١٢,٢% أشار الي ذلك تقرير الأرقام القياسية لأسعار الجملة الذي يصدره الجهاز المركزي للتعينة العامة والاحصاء بصفة شهرية.

وأشار التقرير الي تراجع أسعار مجموعة من المنتجات من أهمها الخضراوات والفواكه والبيض وغيرها حيث انخفضت اسعار الخضراوات بنسبة ٨١,١% والفواكه بنسبة ٦,٨% والبيض بنسبة ١١,٣% .

وصرح اللواء أبو بكر الجندي رئيس الجهاز المركزي للتعينة العامة والاحصاء بأن أسعار بعض السلع الاخرى في اسواق الجملة سجلت ارتفاعا قدره ٤,٢% خلال شهر مايو الماضي مقارنة بالاسعار في نفس الفترة من العام الماضي

مثل المواد الغذائية والوقود ومنتجاته فقد ارتفعت اسعار الوقود بنسبة ١٣,٦% والمواد الخام بنسبة ١,٨% والسلع الوسيطة بنسبة ٦,٨% والسلع الاستهلاكية غير المعمرة بنسبة ٣,١% والسلع الاستهلاكية المعمرة بنسبة ٣% والسلع الاستثمارية

شكل (٨) نتيجة تحييف وتقطيع صورة النص المستقطع من إحدى الصحف اليومية

سحلت اسعار لحملة السلع لزر عية تحقا ضا ملحوظا .حن.ل شهر مايو لما في لهذا العان مفار نه نالاسعار فنفنن ا لفترة مد العان المامن نلعن نسفته ١٢ ا ثرى ذلك تقرير فهرق لقي نية فهنعر لعملة لاذع يمدره لعهرز لمركز ء للتعينة العامة والاحصاعن بصفة شهرية . و ثر لتقريرى تر جمع نعر معموعة ق لمننتق طممه لعفروت والقواكه والبيفن و غير طما حيث ا تحفضت اسعارا لحضراوات بنسبة ٥٨١١ و القواكه بنسبة ٥٦١٥ والبيفن بنسبة ٥١١٣ . و صرح اللواعن ابو بكر الحنذع ر سنن الحهاز المركز ع للتعينة العامة والاحصاعن با ذ اسعار بعفن السلع الا .جر ع ق اسوا ا لحملة سحلت ارتقا قدر ١٢٥ ه .حن.ل شهر مايو لماض مقارنة بالاسعار ق نفنن لفترة مد لعان لماض مثل المواد لغذا نية والوقود و مننتحاته فقد ار تقعت اسعار الوقود بنسبة ١٣٦ ه و المواد الحان بنسبة ١٨ ه و السلع الو سيطرة بنسبة ٦٨ ه السلع الاستهل؛ كية غير المعمرة بنسبة ٣١ ه و السلع الا ستهل. كية المعمرة بنسبة ٥٣ ه و السلع الا ستهلاره . .

شكل (٩) نتيجة التعرف على الحروف في صورة النص المستقطع من إحدى الصحف اليومية وقد جاءت نتائج تقييم معدلات التعرف على الحروف العربية بناء على هذه العينات كما هو مبين في الجدول (٤).  
جدول (٤) نسبة تمييز الحروف على عينات الاختبار

عدد الحروف	الحروف المقروءة	نسبة التمييز
٧٢٣	٥٧٢	٧٨,٧٨٧٨٨

هذه النتيجة ٧٨,٧٨٧٨٨% أفضل من الأداء المتوسط المحسوب ٦٨% لأنها تمثل عينة واحدة من عينات الاختبار مقارنة بمتوسط المحسوب بالنسبة لـ ٢٧ عينة من النصوص. ومعظم الأخطاء نجمت من الخلط بين الحروف المنقوطة ومن الخلط بين حروف الـ س والـ ش مع الثبرات المتكررة لحروف بـ ، نـ ، تـ والخلط بين الهمزات وحرف ع . وكذلك تعريف بعض الرموز الرياضية مثل % هذه الأخطاء تتطلب الأخذ في الاعتبار السمات المورفولوجية للحروف بجانب السمات التوبولوجية . وهو اتجاه التطوير المقترح.

## ٥-الخلاصة :

هذه الورقة تلخص نتائج تصميم منظومة تجريبية للتعرف على حروف اللغة العربية تعتمد على المنهج التركيبي للتعرف على الحروف. تم تنفيذ البرنامج بلغة (C#). للتعرف على الحروف العربية يقوم بمعالجة أولية للنص المأخوذ من الماسح الضوئي ثم يقوم بعملية تحريف لهذه الحروف وبالتالي تحويلها إلى سلسلة متصلة من النقاط ثم تصنيفها إلى مجموعات من المقاطع، وتحليل العلاقات والمسافات بين نقاط المقاطع تتم عملية تمييزها ثم تركيبها للتعرف على الحروف. في هذا الإطار تم تطوير خوارزم التحريف "المحور الأوسط" وخوارزم تقطيع الكتابة باستخدام الوصف التركيبي لحروف اللغة العربية. تشير النتائج التجريبية الأولى إلى تحقيق دقة عالية في تحريف خطوط الكتابة وتقطيعها و معدل جيد من التعرف على حروف اللغة العربية يصل إلى ٧٨% عمليا. هذه النتائج أقل من النتائج المنشورة [١٤] عن أداء برنامج صخر ٩٠.٣٣% وبرنامج "أومنيبيج" ٨٦.٨٩%. هذا البرنامج متاح على شبكة الإنترنت وهو تحت التطوير ومحاو التطوير المقترحة هي تدقيق قواعد تقطيع النص وتطوير قواعد توصيف الحروف العربية لتحتوى السمات المورفولوجية بالإضافة إلى السمات التوبولوجية.

## المراجع :

- [1] Gheith A. Abandah , Mohammed Zeki Khedher, "Printed and Handwritten Arabic Optical Character Recognition – Initial Study", University of Jordan, Faculty of Engineering and Technology, Amman, August 2004,
- [2]H. Al-Yousefi and S. S. Udpa, "Recognition of Arabic Characters", IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 14,no. 8, Aug.. 1992,
- [3] Maget Mahomed Mahmoud Fahmy\*, Hytham El-Messiry\*\*, "Automatic Recognition of Typewritten Arabic Characters Using Zernike Moments as a Feature Extractor", \*Computer Science Department, College of Science, University of Bahrain, Isa Town,

- \*\*Informatics Research Institute, Mubarak City for Scientific Research and Technological Applications, Mansheyat Al-Olama, Alexandria, Egypt,  
web:[http://www.ici.ro/ici/revista/sic2001\\_3/art4.html](http://www.ici.ro/ici/revista/sic2001_3/art4.html),
- [4]Issam Bazzi, Richard Schwartz, and John Makhoul , "An Omnifont Open-Vocabulary OCR System for English and Arabic", IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 21, no. 6, Jun. 1999,
- [5] L. Lam, S. Lee, and C. Suen, "Thinning Methodologies—A Comprehensive Survey," IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 14, no. 9, pp. 869-885, Sept. 1992,
- [6] Paul C.K.Kwok, "A Thinning Algorithm by Contour Generation", Department of Computer Science, University of Calgary, 2500 University Drive NW, Calgary, Canada T2N 1N4,
- [7] Tal Steinerherz, Nathan Intrator and Ehud Rivlin, "A Special Skeltonization Algorithm for Cursive Words,
- [8] Jung-Me Park, Hui-Chuan Chen, Shu T. Huang., "A New Gray Level Edge Thinning Method", University of Alabama, Tuscaloosa, AL 35487,
- [9] L. Lam and C. Suen, "An Evaluation of Parallel Thinning Algorithms for Character Recognition," IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 17, no. 9, pp. 914-919, Sept. 1995,
- [10] M. Altuwajri and M. Bayoumi, "A New Thinning Algorithm for Arabic Characters Using a Self-Organizing Neural Network," Proc. IEEE Int'l Symp. Circuits and Systems, vol. 3, pp. 1824-1827, 1995,
- [11]Maher Ahmed and Rabab Ward, "A Rotation Invariant Rule-Based Thinning Algorithm for Character Recognition", IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 24, no. 12, Dec. 2002,
- [12]Steven S. Skiena, "The algorithm design manual", Springer Telos, NY. 1998.
- [13]Richard G. Casey and Eric Lecolinet, "Survey of Methods and Strategies in Character Segmentation", IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 18, no. 7, Jul. 1996,
- [14] Tapas Kanungo, Gregory A. Marton, Osama Bulbul, " Performance Evaluation of Two Arabic OCR Products", Center for Automation Research University of Maryland, College Park, MD 20742, web:<http://www.cfar.umd.edu/~kanungo>



# المؤتمر الثاني لهندسة اللغة

تنظمه

الجمعية المصرية لهندسة اللغة

تحت رعاية

الأستاذ الدكتور / حسن أحمد غراب

رئيس جامعة عين شمس

١٨ أبريل ١٩٩٩

القاهرة - جمهورية مصر العربية

الأبحاث المقبولة



*The Second Conference  
on Language Engineering  
18 April 1999*

عنوان البحث

تجانس شكلي منظومة الأرقام العربية مع أشكال الحروف العربية  
وأشكال حروف لغات أخرى

أ.د / محمد يونس الصملاوي

د / محمد يسرى النحاس

## تجانس شكلي منظومة الأرقام العربية مع أشكال الحروف العربية وأشكال حروف لغات أخرى

د. محمد يسرى النحاس  
قسم هندسة النظم والحاسبات  
كلية الهندسة جامعة الأزهر  
القاهرة

أ. د. محمد يونس الحملاوي  
قسم هندسة النظم والحاسبات  
كلية الهندسة جامعة الأزهر  
القاهرة

### نبذة

يهدف هذا البحث إلى دراسة مدى التجانس بين فئتي منظومة الأرقام العربية (العربية المشرقية والغبارية المغربية) وأشكال الحروف في اللغة العربية وبعض اللغات الأخرى. ولقد أجرينا دراسة على كل من أشكال حروف اللغات الآتية: اللغة العربية، واللغة اللاتينية، واللغة اليونانية، واللغة السنسكريتية (الهندية) ومقارنة ذلك مع أشكال منظومة الأرقام العربية بفئتيها وصنفت الأشكال من خلال شجرة وراثية. وهذا الاختيار يتيح لنا تمحيص وجه الحقيقة في الجدول الدائر حول الارتباط التاريخي بين شكل الأرقام العربية والحضارة الهندية القديمة والحضارة الإسلامية العربية.

بالنسبة للعينة الأولى والمكونة من أشكال الحروف العربية، أظهرت النتائج تجانساً أكبر لفئة الأرقام العربية المشرقية مع أشكال الحروف، بينما كان التجانس بين فئة الأرقام الغبارية المغربية وأشكال الحروف العربية أقل. ولقد تم عمل شجرة وراثية تضم أشكال فئتي الأرقام وحروف اللغة العربية مع حروف اللغة اللاتينية وصنفت الأشكال من خلال تلك الشجرة الوراثية وأدى ذلك إلى ظهور تجانس للأرقام العربية المشرقية مع الحروف العربية عند المستوى الأدنى من الشجرة الوراثية مما يبين أن أصولها التاريخية أقرب إلى العربية. وعلى النقيض من هذا ظهر تجانساً بين الأرقام الغبارية المغربية وأشكال الحروف العربية عند المستوى الأعلى من الشجرة مما يبين أن أصولها التاريخية أبعد عن العربية.

أما بالنسبة للعينة الثانية والمكونة من أشكال الحروف السنسكريتية (الهندية)، فقد كانت النتائج على العكس من التجربة السابقة حيث أظهرت النتائج تجانساً أكبر لفئة الأرقام الغبارية المغربية مع أشكال الحروف السنسكريتية الهندية، بينما كان التجانس بين فئة الأرقام العربية المشرقية وأشكال الحروف السنسكريتية الهندية أقل.

وبالنسبة للعينة الثالثة والمكونة من أشكال الحروف اللاتينية، فقد أظهرت النتائج تجانساً أكبر

المؤتمر الثاني لهندسة اللغة، القاهرة، 18 إبريل، 1999ء

لفئة الأرقام الغبارية المغربية مع أشكال الحروف اللاتينية، بينما كان التجانس مع فئة الأرقام العربية المشرقية أقل.

وبالنسبة للعينة الرابعة والمكونة من أشكال الحروف اليونانية، فقد أظهرت النتائج تجانساً أكبر لفئة الأرقام الغبارية المغربية مع أشكال الحروف اليونانية، بينما كان التجانس مع فئة الأرقام العربية المشرقية أقل.

## ١ مقدمة:

يهدف هذا البحث إلى عمل دراسة مقارنة بين درجة التوافق الشكلى لكل من فئتي منظومة الأرقام العربية المستخدمتين في منظومة الأرقام العشرية مع منظومة الكتابة العربية. هاتان الفئتان تسميان خطأً فئة الأرقام الهندية وفئة الأرقام العربية<sup>١</sup>. بينما الأصح تسميتهما بالأرقام العربية المشرقية والأرقام المغربية الغبارية. والبحث يستهدف وضع الأسس العلمية لتوصيف منظومة الكتابة العربية آخذين في الاعتبار كل الخواص والمعايير.

والخاصية التي نبحثها هنا هي درجة توافق بنية أشكال الأرقام في كل فئة مع منظومة الكتابة العربية وكتابات عض اللغات الأخرى. إن علم التعرف على الأنماط يتيح لنا استخدام القياس الكمي في دراسة الخواص البنيوية الشكلية للأنماط من خلال عدة معايير موضوعية. والحكم الموضوعي على درجة التوافق يخضع لعدة معايير منها:

١- درجة تنافر نمط أشكال الأرقام في كل فئة مع نمط أشكال الحروف العربية.

٢- درجة توافق نمط أشكال الأرقام في كل فئة مع نمط أشكال الحروف العربية.

....

في هذا البحث نستخدم المدخل الإحصائي لعلم التعرف على الأنماط كوسيلة للمقارنة بين منظومتى كتابة الأرقام: منظومة الكتابة العربية المشرقية، ومنظومة الكتابة الغبارية المغربية. هذا المدخل الإحصائي يعتمد على الاختيار الأفضل لثلاث أركان للمنهج الموضوعي في المقارنة وهي: التمثيل، المعيار، والإحصائيات. وإيضاح هذا المنهج نبدأ بتعريف نمط الأشكال والكائنات المنتمية له وكيفية تمثيله من خلال مدبج السمات، ونحدد المعيار المستخدم في عملية المقارنة بحيث تكون نتائج عملية المقارنة موضوعية غير متحيزة ومعتمدة عملياً. ثم نعرض الإحصائيات المختلفة الممكن استخدامها في المقارنة المطروحة. وبعد أن نحدد أركان منهج المقارنة نعرض كيفية تطبيق هذا المنهج لإجراء المقارنة الكمية المطلوبة

<sup>١</sup>محمود فهمي حجازي ومحمد يونس الحملوى ومحمد يسرى النحاس؛ أرقامنا العربية: الأرقام المشرقية والأرقام المغربية؛ المؤتمر السنوى الثامن لتعريب العلوم؛ القاهرة؛ ٢٠-٢١ مارس ١٩٩٦م

المؤتمر الثامن لمنظمة اللغة، القاهرة، ١٨ أبريل ١٩٩٩م

ونعرض نتائج التطبيق ممثلة في الإحصائيات التلخيصية ونوضح كيفية تفسير هذه النتائج. ولتدعيم نتائج البحث فلقد تم تصنيف أشكال الأرقام والحروف في شكل شجرة وراثية معيار التفرع فيها هو مدى التجانس بين أشكال الحروف والأرقام. للبحث عن هذه الشجرة رتبنا بطريقة عشوائية خليطاً من جميع أشكال فنتى منظومة الأرقام العربية (المشرقية والمغربية) مع أشكال الحروف العربية (خط النسخ) والحروف اللاتينية ثم اتبعنا منهج التعرف على الأنماط لإيجاد التجمعات المتجانسة من الأشكال في هذا الخليط العشوائي، وترتيب هذه التجمعات في مستويات هرمية عند كل درجة تجانس.

ولقد أجرينا هذه الدراسة على أربع عينات تمثل أشكال حروف الخط العربي (خط النسخ) و الخط الهندي (السنسكريتي) والخط اللاتيني والخط اليوناني.

## ٢ منهاج الدراسة:

### ٢-١ تعريف المشكلة:

يمكن صياغة مشكلة إيجاد التجمعات المتجانسة على النحو التالي: بدءاً من فئة كائنات مجهولة الهوية؛ ش؛

$$ش = \{ش_١، ش_٢، \dots، ش_ن\}$$

مطلوب البحث عن بنية التجمعات المتجانسة من هذه الكائنات؛ ج؛

$$ج = \{ج_١، ج_٢، \dots، ج_ص\}$$

لحل هذه المشكلة يجب تعريف معياراً لقياس درجة التجانس داخل كل تجمع ثم البحث عن فئة التجمعات التي تحقق القيمة المثلى لهذا المعيار. بالرغم من أن هذه المسألة قد تجد حلاً مثالياً ومباشراً في بعض المشاكل العملية، إلا أنه من الصعب إيجاد حلاً لها في الكثير من المشاكل العملية الأخرى وذلك للمشاكل الناجمة عن كبر حجم فئة الكائنات المطلوب تجميعها وكبر حجم متجه السمات الممثل لكل كائن منها. ويكفي لمعرفة مدى صعوبة البحث عن هذه الحلول حساب عدد التقاسيم الممكنة؛ ت؛ من العلاقة التوفيقية التالية<sup>٢</sup>:

$$ت = \frac{ص}{ن} \frac{ص-١}{ن-١} \dots \frac{ص-٢}{ن-٢} \dots \frac{ص-١}{٢} \frac{ص}{١} = \frac{ص!}{(ن-ص)!}$$

هذه العوامل يجب مراعاتها بدقة عند البحث عن أفضل تقسيم للتجمعات المتجانسة والذي

<sup>٢</sup> R.O. Duda & P. Hart, "Pattern Classification and Scene Analysis", A. Wiley Interscience Publication, New York, 1974

المؤتمر الثاني لمنحة اللغة، القاهرة: ١٨ إبريل، ١٩٩٩ء

يحقق درجة عالية من المصادقية في الحكم على الكائنات الجديدة المطلوب دمجها في هذه التجمعات. ويقدم هذا البحث طريقة غير متحيزة للبحث عن العلاقات بين أشكال الحروف والأرقام كما يلي شرحه.

#### ٢-٢ تمثيل الأشكال:

حيث أن أشكال الأرقام والحروف هي التي تهتمنا في هذا البحث فإننا سنمثل منظومة الأرقام والحروف بفئة الأشكال النمطية المستخدمة؛ ز؛ والمكونة من تسع وثلاثون شكلاً كالاتي :

$$ز = \{ ش١، ش٢، ش٣، ش٤، ش٥، ش٦، ش٧، ش٨، ش٩ \}$$

وحيث إن الاختلافات المطروحة واقعيًا في الينط أو الخط المستخدم كثيرة ومتنوعة فقد قمنا باختيار الخط النسخ نظراً لشيوعه في الكتابة والطباعة. وقد اختير حجم الخلية الممثلة للحرف أو الرقم أكبر من الحد الأدنى الضروري لتمييز شكله مما يسمح بدراسة الفروق الدقيقة في بنية الأشكال.

ويمثل كل نمط من الأشكال بمتجه السمات؛ ش ١؛ المكون من عدد؛ ن؛ من السمات الأساسية؛ س ل؛ كالاتي:

$$ش ١ = (س١، س٢، س٣، س٤، س٥، س٦، س٧، س٨، س٩، س١٠، س١١، س١٢، س١٣، س١٤، س١٥، س١٦، س١٧، س١٨، س١٩، س٢٠، س٢١، س٢٢، س٢٣، س٢٤، س٢٥، س٢٦، س٢٧، س٢٨، س٢٩، س٣٠، س٣١، س٣٢، س٣٣، س٣٤، س٣٥، س٣٦، س٣٧، س٣٨، س٣٩، س٤٠، س٤١، س٤٢، س٤٣، س٤٤، س٤٥، س٤٦، س٤٧، س٤٨، س٤٩، س٥٠، س٥١، س٥٢، س٥٣، س٥٤، س٥٥، س٥٦، س٥٧، س٥٨، س٥٩، س٦٠، س٦١، س٦٢، س٦٣، س٦٤، س٦٥، س٦٦، س٦٧، س٦٨، س٦٩، س٧٠، س٧١، س٧٢، س٧٣، س٧٤، س٧٥، س٧٦، س٧٧، س٧٨، س٧٩، س٨٠، س٨١، س٨٢، س٨٣، س٨٤، س٨٥، س٨٦، س٨٧، س٨٨، س٨٩، س٩٠، س٩١، س٩٢، س٩٣، س٩٤، س٩٥، س٩٦، س٩٧، س٩٨، س٩٩، س١٠٠)$$

تعتبر مجموعة نقاط الصورة الثنائية هي السمات الذاتية لذلك الشكل. هذه السمات لا تصلح للمقارنة كما هي حيث أنها تتغير بالنقل أو الدوران أو الانكماش في فراغ الصورة. لذلك فإن شكل الحرف يجب معالجته أولاً باستخدام استحوالة خطية؛ لقياس بعد نقاط الشكل عن مركز ثقله ومعايرتها بالنسبة لمحاور الشكل الأساسية وحجمه<sup>٢</sup>. الشكل في الصورة الثنائية الناتجة يتميز بأنه لا يتغير مع أية تحويلات هندسية في فراغ الصورة ويمكن تمثيل كل نمط بمتجه السمات، ش، المكون من نفس العدد ن من السمات.

#### ٣-٢ قياس التشابه:

اخترنا في هذا البحث دالة الارتباط؛ ر؛ لقياس درجة التشابه. وتعرف دالة الارتباط؛ ر؛ بين شكلين ش ١، ش ٢ بالعلاقة الآتية:

$$R = \frac{\sum_{i=1}^n (ش١_i - \bar{ش١})(ش٢_i - \bar{ش٢})}{\sqrt{(\sum_{i=1}^n (ش١_i - \bar{ش١})^2)(\sum_{i=1}^n (ش٢_i - \bar{ش٢})^2)}}$$

$$R = (م، ل)$$

$$R = \frac{\sum_{i=1}^n (ش١_i - \bar{ش١})(ش٢_i - \bar{ش٢})}{\sqrt{(\sum_{i=1}^n (ش١_i - \bar{ش١})^2)(\sum_{i=1}^n (ش٢_i - \bar{ش٢})^2)}}$$

<sup>٢</sup> K. Fukunaga, "Statistical Pattern Recognition", Academic Press, inc., San Diego, 1990

المؤتمر الثاني لصحيفة اللغة، القاهرة، ١٨ أبريل ١٩٩٩م

حيث يتم حساب المجموع في التعريف السابق لكل النقاط (س ، ص ) داخل المنطقة المشتركة للشكلين؛ ق. هذه الدالة تمثل النسبة بين الطاقة الكامنة في النقاط المشتركة بين الشكلين والطاقة الكامنة في نقاط كل منهما. وكما هو واضح من التعريف فهذا القياس لا يعتمد على عدد النقاط الممثلة لكل شكل. ولكنه يعتمد على الوضع النسبي للشكلين معاً في فراغ الصورة. وهذه الدالة توافقية عادة ولها قيمة عظمى عامة عند أقصى درجة ارتباط وعدة قيم عظمى جانبية تنتشر حول القيمة العظمى العامة. ولزيادة مصداقية الاعتماد على نتائج قياس هذه الدالة يمكن قياس معامل الارتباط الأعظم؛ (م ع ، ل ع). أما في حالة معايرة الأشكال بالنسبة لنظام المحاور المستخدم في كتابة المنظومة العربية فإنه يمكن الاكتفاء بقياس معامل الارتباط عند مركز ثقل الشكل أي النقطة (م ع ، ل ع).

#### ٢-٤ تصميم شجرة التجمعات المتجانسة:

هناك عدة طرق لإيجاد التجمعات المتجانسة من الأنماط لخليط من الكائنات المجهولة. والطريقة التي نستخدمها في هذا البحث تدرج تحت طرق التجميع التراكمي الهرمي. في هذه الطريقة تبدأ تتابعية التجميع بتتابعية بها عدد من التجمعات مساو لعدد الكائنات ثم عند كل مستوى من مستويات التجميع تدمج كائنات التجمعات مع بعضها البعض، وحيث أن الكائنات المدمجة عند مستوى ما تبقى سوية في المستوى الأعلى فهذه التتابعية تسمى بالتجميع التراكمي الهرمي. وصحة شجرة التجمعات المتجانسة تعتمد بصورة عامة على الترتيب الابتدائي للتجمعات، وعلى ترتيب الكائنات داخل كل تجمع، وعلى اختيار معيار التشابه بين التجمعات. لذا فإن أي تصميم لهذه الإجرائية لا بد أن يأخذ في الاعتبار كل هذه العوامل. لذلك فقد قمنا بتطوير هذه الطريقة لتتلافى تأثير هذه العوامل بإضافة التالي:

١- عملية خلط عشوائي تام للكائنات أو التجمعات قبل مقارنتها مع بعضها البعض.

٢- استخدام معيار الزوج الأكثر تشابهاً لاستقطاب التجمعات ودمجها وهو في بحثنا هذا يمثل الوضع الأسوأ في الحكم على صحة النتائج.

سنستخدم هنا اللغة المعروفة بـ "شبه شفرة" لتمثيل إجرائية تصميم الشجرة. في ما نعرضه هنا تمثل الكلمات ذات المدلول الحاسوبي بالخط الكوفي العريض والكلمات ذات المدلول الرياضي بخط النسخ المائل العريض.

إجرائية التجميع التراكمي الهرمي:

مدخلات: فئة أشكال الأرقام والحروف؛ ش

مخرجات: شجرة التجمعات المتجانسة؛ ت

الخطوات:

خلق عدد ن من التجمعات يحتوي كل منها على شكل واحد فقط

أبدأ بعدد التجمعات الحالية ص = ن

كره ما يلي طالما عدد التجمعات الحالية ص < ١

بداية

اخلط التجمعات الحالية عشوائيا

ابحث عن أقرب زوج من التجمعات؛ (ج م ، ج ن)

أدمج زوج التجمعات (ج م ، ج ن)

احذف التجمع الأدنى ج ن

ارفع مستوى التجمع الأعلى ج م

أنقص عدد التجمعات الحالية واحدا

نهاية

### ٣ التطبيق والنتائج:

لإجراء التجارب المطلوبة في هذا البحث تم تصميم وتنفيذ تطبيق حاسوبي يعمل في بيئة النوافذ يتكون من الوحدات التالية:

- ١- واجهة التطبيق،
- ٢- مولد الأشكال،
- ٣- معالج الصور،
- ٤- مصنف الكائنات،
- ٥- عارض البيانات.

عند إجراء تجربة ما على أشكال الأرقام وحروف الخط العربي تتبع الخطوات التالية:

- ١- تولد أشكال منظومة كتابة الأرقام والحروف؛ ز؛ بواسطة برنامج مولد الأشكال.
- ٢- تعالج الأشكال الناشئة بحصرها داخل نافذة ثم تحذف هذه الأشكال وتحدد نافذة مكونة من  $20 \times 17$  نقطة.
- ٣- يحسب لكل شكل من الأشكال؛ ش<sub>١</sub>؛ متجه مركز الثقل؛ م<sub>١</sub>؛ ومصفوفة الانتشار؛ ك<sub>١</sub>؛ ودالة الاستحالة؛ د<sub>١</sub>.
- ٤- يعاد توقيع النقاط السوداء في كل شكل من الأشكال؛ ش<sub>١</sub>؛ باستخدام دالة الاستحالة الخطية؛ د<sub>١</sub>؛ لينشأ عن ذلك الشكل المعيار؛ ش<sub>١</sub>.
- ٥- تحسب دالة الارتباط بين جميع الثنائيات الممكنة للأشكال؛ (ش<sub>١</sub> ، ش<sub>٢</sub>)؛ في كل منظومة

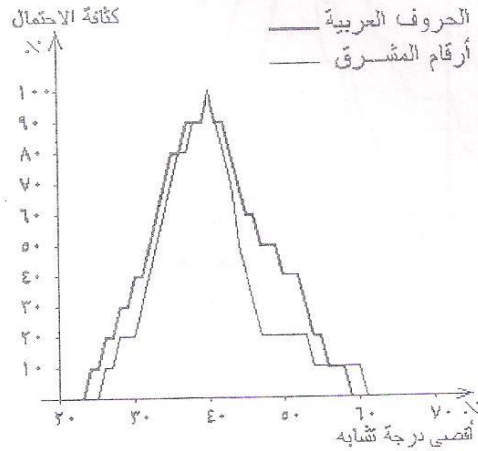


٦- منحني كثافة الاحتمال لفئة الحروف اليونانية.

ويرسم منحني كثافة الاحتمال لكل منظومة من منظومات الحروف مع احدى منظومتى الأرقام يمكننا معرفة إذا ما كانت فئة الأرقام المستخدمة تمثل نمطاً مستقلاً عن نمط بقية عناصر المنظومة أم أنها تمثل مجرد فئة جزئية من عناصر المنظومة. ويستدل على ذلك بظهور وسطين حسابيين في منحني كثافة الاحتمال لدرجات التشابه، وكلما زاد الفرق بين هذين الوسطين كلما قلت درجة انتماء كل من النمطين إلى الآخر، أي نمط الحروف ونمط الأرقام. ولقد تم حساب منحنيات كثافة الاحتمال السابقة ومثلت النتائج في الأشكال من رقم (٢) إلى رقم (١٠).

٣-١ نتائج القياسات على الأرقام العربية المشرقية والحروف العربية:

يمثل الشكل رقم (٢) منحني كثافة الاحتمال لفئة الحروف العربية مع منحني كثافة الاحتمال لفئة الأرقام العربية المشرقية.



شكل ٢ منحني كثافة الاحتمال لأقصى درجة تشابه لفئتي الحروف العربية والأرقام العربية المشرقية

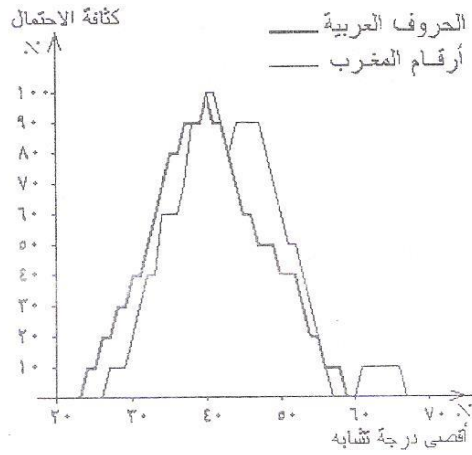
ويلاحظ في هذا الشكل أن الفئتين لهما نفس متوسط أقصى درجة تشابه، كذلك يلاحظ أن منحني كثافة الاحتمال لفئة الأرقام العربية المشرقية يقع بالكامل تحت منحني كثافة الاحتمال للحروف العربية. وهذا يعني في نظرية علم التعرف على الأنماط أن الفئتين يصعب التمييز بينهما. أي أن أشكال الأرقام العربية المشرقية تمثل فئة جزئية من فئة أشكال الحروف العربية مما يعنى

المؤتمر الثاني لمنظمة اللغة، القاهرة، ١٨ أبريل ١٩٩٩م

أن السمات الأساسية للأرقام العربية المشرقية والحروف العربية واحدة مما قد يدل على أن الثقافة التي أنتجت كلا من الفئتين واحدة. والجدير بالذكر أيضاً أن منحنى كثافة الاحتمال للأرقام العربية المشرقية باحتوانه على قمة واحدة يشير إلى أن هذه الأرقام متناسقة فيما بينها مما يعني أنها وليدة ثقافة واحدة.

### ٣-٢ نتائج القياسات على الأرقام الغبارية المغربية والحروف العربية:

لقد تم إجراء نفس القياسات بالنسبة للأرقام الغبارية المغربية حيث يمثل الشكل رقم (٣) منحنى كثافة الاحتمال لفئة الحروف العربية مع منحنى كثافة الاحتمال لفئة الأرقام الغبارية المغربية.



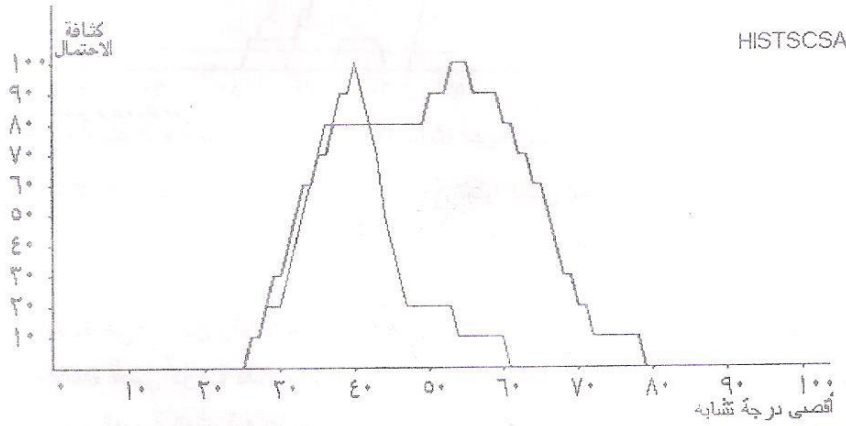
شكل ٣ منحنى كثافة الاحتمال لأقصى درجة تشابه لفئتي الحروف العربية والأرقام الغبارية المغربية

ويلاحظ في هذا الشكل أن متوسط أقصى درجة تشابه لفئة الأرقام الغبارية المغربية أكبر من متوسط أقصى درجة تشابه لفئة الحروف العربية. وكذلك يُلاحظ أن منحنى كثافة الاحتمال لفئة الأرقام الغبارية المغربية ممتد يمينا خارج منحنى كثافة الاحتمال للحروف وبه أكثر من قمة مما يعني احتواء مجموعة أرقام المغرب على فئات جزئية غير ظاهرة. وهذا يعني في نظرية علم التعرف على الأنماط أن الفئتين يسهل التمييز بينهما. أي أن أشكال الأرقام الغبارية المغربية تمثل فئة مستقلة عن فئة أشكال الحروف العربية، مما يعني اختلاف السمات الأساسية بين الأرقام الغبارية المغربية والحروف العربية. والجدير بالذكر أيضاً أن منحنى كثافة الاحتمال

للأرقام الغبارية المغربية باحتوائه على أكثر من قمة يشير إلى أن هذه الأرقام غير متناسقة فيما بينها مما يعنى أنها وليدة أكثر من ثقافة.

٣-٣ نتائج القياسات على الأرقام العربية المشرقية والحروف الهندية السنسكريتية:

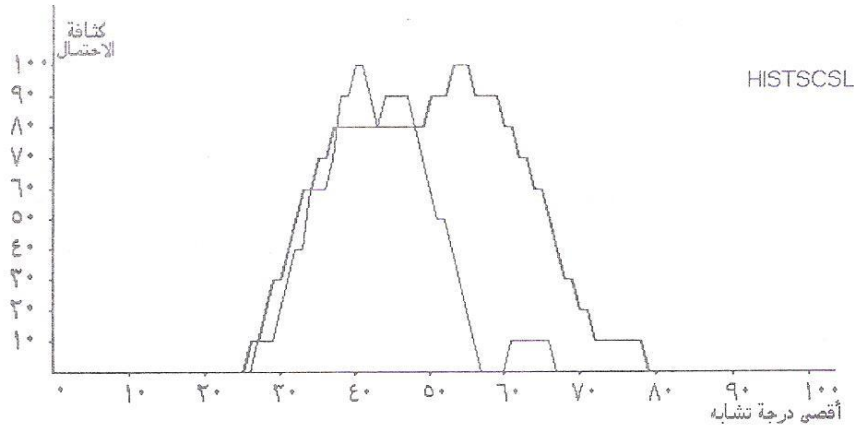
يمثل الشكل رقم ٤ منحنى كثافة الاحتمال لفئة الحروف الهندية السنسكريتية مع منحنى كثافة الاحتمال لفئة الأرقام العربية المشرقية.



شكل ٤ منحنى كثافة الاحتمال لأقصى درجة تشابه لفئتي الحروف الهندية والأرقام العربية المشرقية

٤-٣ نتائج القياسات على الأرقام الغبارية المغربية والحروف العربية:

لقد تم إجراء نفس القياسات بالنسبة للأرقام الغبارية المغربية حيث يمثل الشكل رقم ٥ منحنى كثافة الاحتمال لفئة الحروف العربية مع منحنى كثافة الاحتمال لفئة الأرقام الغبارية المغربية.

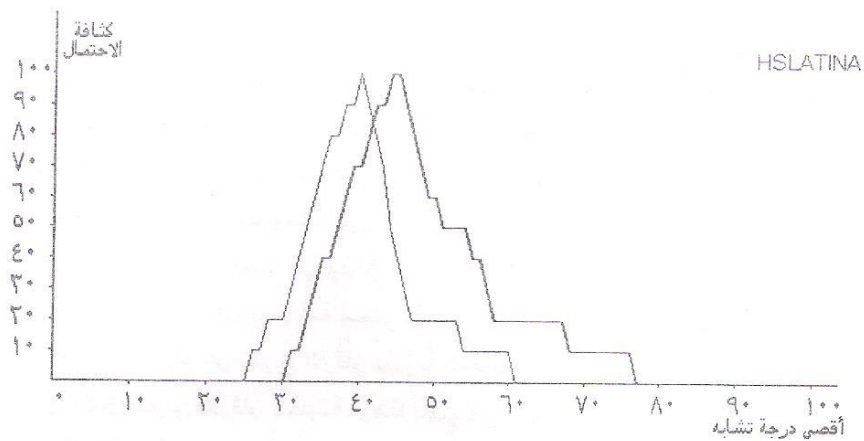


شكل ٥ منحنى كثافة الاحتمال لأقصى درجة تشابه لفتى الحروف الهندية والأرقام الغبارية المغربية

ويلاحظ في الشكلين رقم ٤،٥ أن متوسط أقصى درجة تشابه لكل فئة من فئتي الأرقام أقل من متوسط أقصى درجة تشابه لفئة الحروف الهندية، إلا أن متوسط أقصى درجة تشابه لفئة الأرقام العربية المشرقية أقل من متوسط أقصى درجة تشابه لفئة الأرقام الغبارية المغربية وبالتالي يبعد أكثر عن متوسط أقصى درجة تشابه لفئة الحروف الهندية. كما أن المساحة المشتركة بين منحنى الحروف الهندية والأرقام العربية المشرقية أقل منها في حالة الأرقام الغبارية المغربية. وهذا يعنى فى نظرية علم التعرف على الأنماط أن فئة الأرقام الغبارية المغربية أقرب إلى فئة الحروف الهندية من فئة الأرقام العربية المشرقية من ناحية السمات. كما يعنى اختلاف السمات الأساسية بين الأرقام العربية المشرقية والحروف الهندية بصورة أكبر عنها فى حالة الأرقام الغبارية المغربية والحروف الهندية.

### ٣-٥ نتائج القياسات على الأرقام العربية المشرقية والحروف اللاتينية:

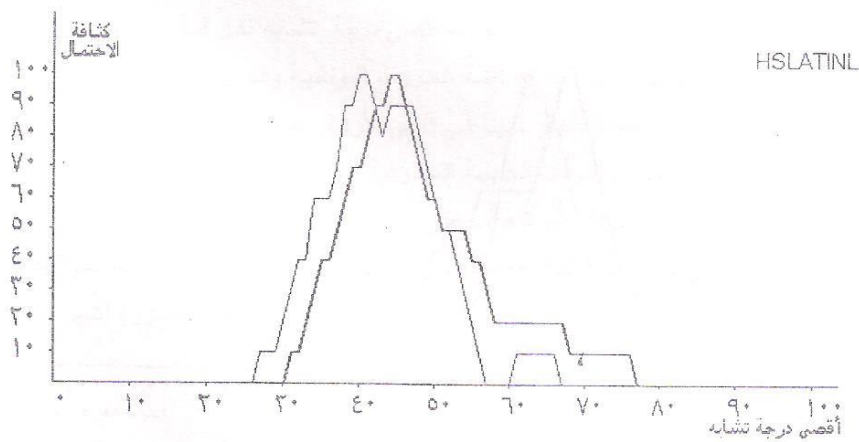
يمثل الشكل رقم (٦) منحنى كثافة الاحتمال لفئة الحروف اللاتينية مع منحنى كثافة الاحتمال لفئة الأرقام العربية المشرقية.



شكل ٦ منحني كثافة الاحتمال لأقصى درجة تشابه لفتى الحروف اللاتينية والأرقام العربية المشروقة

٦-٣ نتائج القياسات على الأرقام الغبارية المغربية والحروف اللاتينية:

لقد تم إجراء نفس القياسات بالنسبة للأرقام الغبارية المغربية حيث يمثل الشكل رقم (٧) منحني كثافة الاحتمال لفئة الحروف اللاتينية مع منحني كثافة الاحتمال لفئة الأرقام الغبارية المغربية.

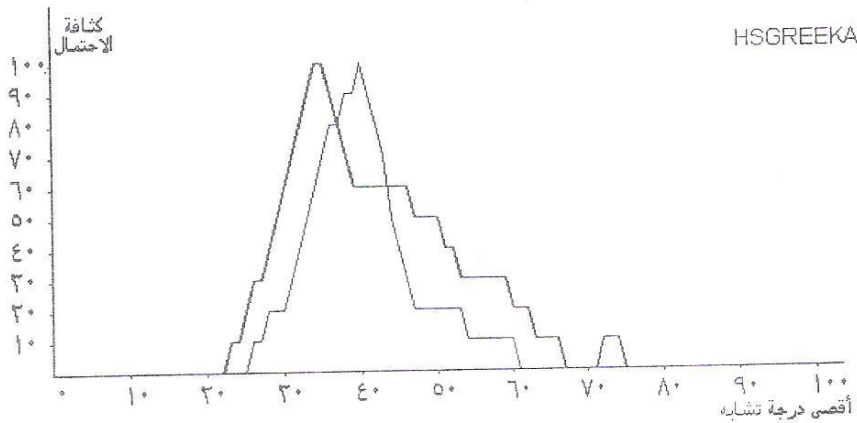


شكل ٧ منحني كثافة الاحتمال لأقصى درجة تشابه لفتى الحروف اللاتينية والأرقام الغبارية المغربية

ويلاحظ في الشكلين رقم ٧،٦ أن متوسط أقصى درجة تشابه لكل فئة من فئتي الأرقام أقل من متوسط أقصى درجة تشابه لفئة الحروف اللاتينية، إلا أن متوسط أقصى درجة تشابه لفئة الأرقام العربية المشرقية أقل من متوسط أقصى درجة تشابه لفئة الأرقام الغبارية المغربية وبالتالي يبعد أكثر عن متوسط أقصى درجة تشابه لفئة الحروف اللاتينية. كما أن المساحة المشتركة بين منحنى الحروف اللاتينية والأرقام العربية المشرقية أقل منها في حالة الأرقام الغبارية المغربية، بالإضافة إلى وقوع القمة الأصغر في منحنى كثافة الاحتمال لفئة الأرقام الغبارية المغربية بالقرب من قمة منحنى كثافة الاحتمال لفئة الحروف اللاتينية مما يثير بعض علامات الاستفهام عن تطويع الأرقام الغبارية المغربية لتصل لصورتها الحالية وإن لم تتخلص من الأصل العربي للأرقام الحديثة. وهذا يعنى فى نظرية علم التعرف على الأنماط أن فئة الأرقام الغبارية المغربية أقرب إلى فئة الحروف اللاتينية من فئة الأرقام العربية المشرقية من ناحية السمات. كما يعنى اختلاف السمات الأساسية بين الأرقام العربية المشرقية والحروف اللاتينية بصورة أكبر عنها فى حالة الأرقام الغبارية المغربية والحروف اللاتينية.

### ٣-٧ نتائج القياسات على الأرقام العربية المشرقية والحروف اليونانية:

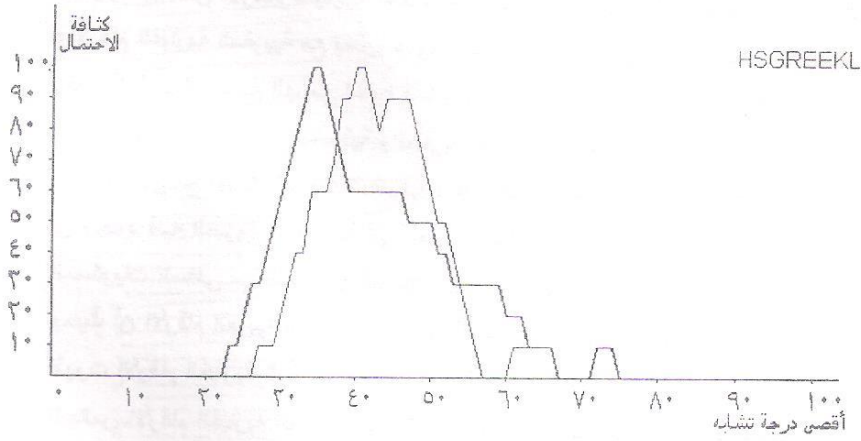
يمثل الشكل رقم (٨) منحنى كثافة الاحتمال لفئة الحروف اليونانية مع منحنى كثافة الاحتمال لفئة الأرقام العربية المشرقية.



شكل ٨ منحنى كثافة الاحتمال لأقصى درجة تشابه لفئتي الحروف اليونانية والأرقام العربية المشرقية

### ٣-٨ نتائج القياسات على الأرقام الغبارية المغربية والحروف اليونانية:

لقد تم إجراء نفس القياسات بالنسبة للأرقام الغبارية المغربية حيث يمثل الشكل رقم (٩) منحنى كثافة الاحتمال لفئة الحروف اليونانية مع منحنى كثافة الاحتمال لفئة الأرقام الغبارية المغربية.



شكل ٩ منحنى كثافة الاحتمال لأقصى درجة تشابه لفتى الحروف اليونانية والأرقام الغبارية المغربية

ويلاحظ كذلك في الشكلين رقم ٩،٨ أن متوسط أقصى درجة تشابه لكل فئة من فئتي الأرقام أكبر من متوسط أقصى درجة تشابه لفئة الحروف اليونانية وقد يشير هذا إلى تشابهات واضحة بين الحروف اليونانية أكثر منها في فئتي الأرقام. كما أن المساحة المشتركة بين منحنى الحروف اليونانية والأرقام العربية المشرقية أقل منها في حالة الأرقام الغبارية المغربية. وهذا يعنى فى نظرية علم التعرف على الأنماط أن فئة الأرقام الغبارية المغربية أقرب إلى فئة الحروف اليونانية من فئة الأرقام العربية المشرقية من ناحية السمات. كما يعنى اختلاف السمات الأساسية بين الأرقام العربية المشرقية والحروف اليونانية بصورة أكبر عنها فى حالة الأرقام الغبارية المغربية والحروف اليونانية.

### ٣-٩ نتائج القياسات على خليط أشكال حروف اللغة العربية وحروف اللغة اللاتينية والأرقام

العربية المشرقية والغبارية المغربية:

لحسم النتائج السابقة فقد كونا خليطاً من ١٠١ عنصراً تمثل الأرقام العربية المشرقية والأرقام

الغبارية المغربية مع حروف خط النسخ العربى وحروف اللغة اللاتينية الكبيرة والصغيرة. ولقد أجريت التجربة بنفس الشروط السابقة وكانت النتيجة هى ظهور تجمعين من الأشكال: التجمع الأول مكون من معظم حروف الخط العربى والأرقام العربية المشرقية مع بعض الحروف اللاتينية وبعض الأرقام الغبارية المغربية، والتجمع الثانى شمل معظم الحروف اللاتينية والأرقام الغبارية المغربية مع بعض حروف الخط العربى والأرقام المغربية. ويوضح الشكل رقم ١٠ شجرة التجمع الهرمى لخليط أشكال حروف اللغة العربية (خط النسخ) وحروف اللغة اللاتينية والأرقام العربية المشرقية والغبارية المغربية. والشكل يؤكد التجانس الواضح لحروف الخط العربى مع الأرقام العربية المشرقية وعدم تجانسها مع الأرقام الغبارية المغربية. ويلاحظ أن وجود فئة الحروف اللاتينية فى الخليط استقطب الأرقام الغبارية المغربية سريعاً عند المستويات السفلى من التجمعات الهرمية مؤكداً تجانسهما معاً.

وحيث أن الأرقام العربية المشرقية أقدم فى النشوء حيث ظهرت لأول مرة عام ٢٠٤ هـ؛ بينما ظهرت الأرقام الغبارية المغربية لأول مرة عام ٤٨٠ هـ بفارق حوالى ثلاثة قرون، وحيث تتجانس الأرقام الغبارية المغربية بصورة أكبر مع أشكال الحروف الهندية واللاتينية واليونانية لذا فمن المنطقى الحكم على تلك الأرقام الغبارية المغربية بأنها تطويع لمنظومة الأرقام العربية كى تلائم تلك اللغات الموجودة فى فترة نشوء تلك الأرقام. ومما يدعم ذلك تاريخياً أن تلك الأرقام الغبارية قد تطورت على ثلاثة مراحل منذ نشأتها حتى وصلت إلى صورتها الحالية. والجدير بالذكر أن كلمة غبارية التى اشتهرت بها تلك المنظومة فى الكتابات العربية أتت من طريقة كتابة الهنود حيث كانوا يثرون الغبار على الألواح ثم يكتبون عليها<sup>٦</sup>.

ونستنتج من المقارنات الإحصائية السابقة أن فئة الأرقام العربية المشرقية يمكن دمجها مع فئة الحروف العربية لتنتج فئة متجانسة فى خواصها الإحصائية. أما فئة الأرقام الغبارية المغربية عند دمجها مع فئة الحروف العربية تنتج فئة غير متجانسة، حيث يتميز منحنى كثافة الاحتمال الخاص بها بوجود قمتين تشيران إلى عدم تجانس تلك الفئة فى ذاتها وبالتالي مع فئة الحروف العربية. وبالتالي نستنتج أن فئة الأرقام العربية المشرقية أكثر تجانساً مع فئة الحروف العربية من تجانس فئة الأرقام الغبارية المغربية مع فئة الحروف العربية.

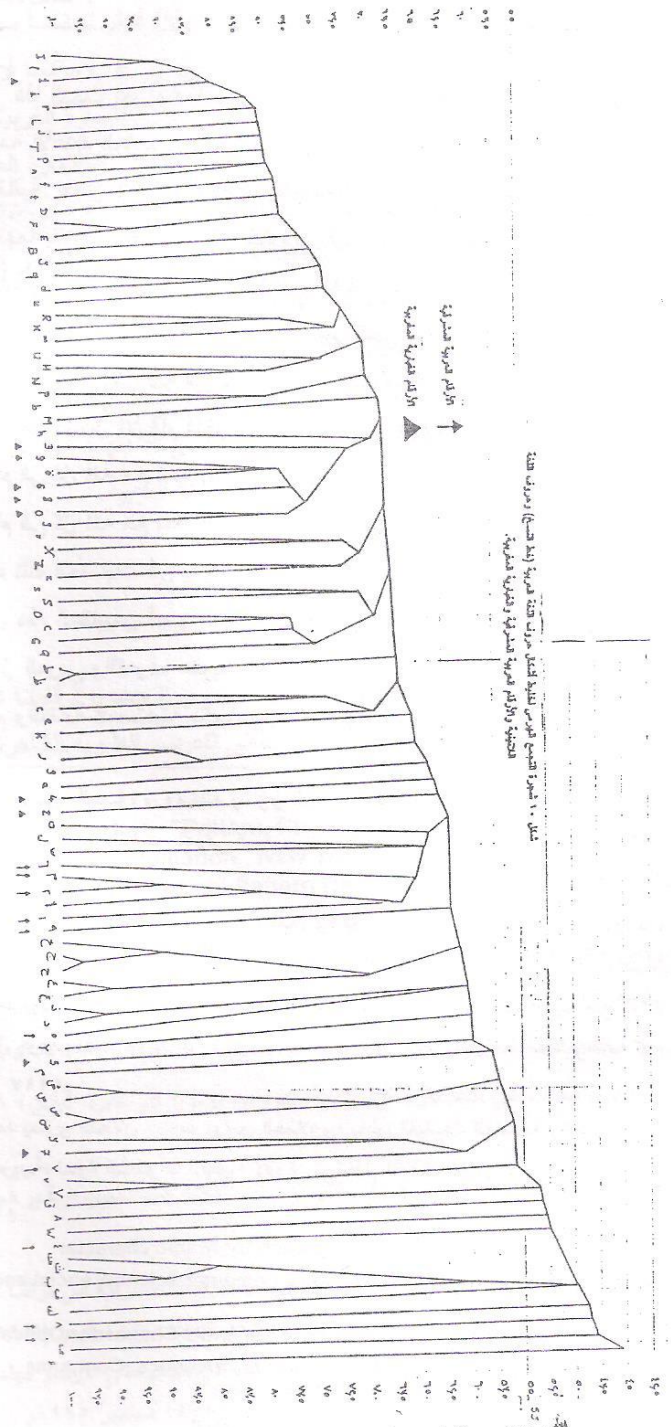
<sup>٦</sup> محمد بن موسى الخوارزمى؛ الجبر والمقابلة؛ دار الكاتب العربى للطباعة والنشر؛ القاهرة؛ ١٩٦٨م؛ صفحة ٢٤

<sup>٧</sup> لجنة الرياضة؛ مجمع اللغة العربية؛ الدورة الحادية والأربعون؛ القاهرة؛ ١٩٧٤-١٩٧٥م

موريس شريل؛ الرياضيات فى الحضارة الإسلامية؛ جروس برس؛ طرابلس، لبنان؛ ١٩٨٨م؛ صفحة ٩١

<sup>٨</sup> قاسم على سعد؛ الأرقام العربية: تاريخها وأصلاتها وما استعمله المحدثون وغيرهم منها، القسم الأول؛ مجلة الأحمديّة؛ العدد ٢؛ ١٩٩٨م؛ دبي

المؤتمر الثانى لمنظمة اللغة؛ القاهرة؛ ١٨ إبريل ١٩٩٩م



## ٤ الخلاصة:

تم في هذا البحث إجراء دراسة مقارنة بين فئتي منظومة الأرقام العربية المستخدمتين في منظومة الأعداد العشرية ومدى تناسقهما مع منظومة الكتابات العربية والهندية واللاتينية واليونانية. هاتان الفئتان تسميان خطأً فئة الأرقام الهندية وفئة الأرقام العربية. بينما الأصح تسميتهما بالأرقام العربية المشرقية والأرقام الغبارية المغربية.

ولقد أتاح المبحث الهندسي لنا استخدام القياس الكمي لسمات كل من الفئتين وإجراء المقارنات الموضوعية بينهما من خلال بعض المعايير الغير متحيزة، حيث تم حساب درجة انتماء نمط أشكال الأرقام في كل فئة مع نمط أشكال الحروف العربية. كما تم تغطية نتائج قياس درجة تشابه نمط أشكال الأرقام داخل كل فئة من فئتي الأرقام<sup>١٠</sup> ونتائج قياس درجة توافق نمط أشكال الأرقام في كل فئة مع نمط أشكال الحروف العربية<sup>١١</sup> وكذلك نتائج قياس درجة توافق نمط أشكال الأرقام في كل فئة مع نمط أشكال حروف بعض اللغات غير العربية<sup>١٢</sup>؛ حيث سبق للمباحثين إجراء تلك الحسابات في بحوث أخرى<sup>١٣</sup>.

وتمثل هذه المقاييس أهمية خاصة لأنها تؤثر بصورة مباشرة في حل مشاكل توصيف الرقم والرمز العربي والتعرف عليهما بالتقنيات الحديثة كالحاسوب<sup>١٤</sup> بالإضافة إلى تصميم منظومات الرسم والكتابة المستعينة بالحاسوب<sup>١٥</sup>.

<sup>١٠</sup> محمد يونس الحملاوي ومحمد يسرى النحاس؛ الأرقام الهندعربية، دراسة مقارنة: التاريخ، السمات، الاستخدام، التقنية؛ المؤتمر الدولي السادس عن الحاسب الآلي بين النظرية والتطبيق؛ الإسكندرية؛ ٥ - ٣ سبتمبر ١٩٩٦م

<sup>١١</sup> محمد يسرى النحاس ومحمد يونس الحملاوي؛ تجانس الأرقام الهندعربية مع أشكال حروف الخط العربي؛ المؤتمر الدولي السابع عن الحاسب الآلي بين النظرية والتطبيق؛ الإسكندرية؛ ٢ - ٤ سبتمبر ١٩٩٧م

<sup>١٢</sup> محمد يونس الحملاوي ومحمد يسرى النحاس؛ تجانس الأرقام الهندعربية مع أشكال الحروف العربية وأشكال حروف لغات أخرى؛ ندوة الأرقام ومكائنها في قضية التعريب؛ مجمع اللغة العربية؛ القاهرة؛ ٢٠ فبراير ١٩٩٧م

<sup>١٣</sup> محمد يسرى النحاس ومحمد يونس الحملاوي؛ بعض القياسات الهندسية والعلمية على مجموعتي الأرقام الهندعربية؛ ندوة الخطوات العملية لإقرار استخدام الأرقام العربية؛ جامعة الأزهر؛ القاهرة؛ ١ مارس ١٩٩٨م

<sup>١٤</sup> M.Y. Mahmoud (El Nahas) et al : "A statistical approach for Arabic character recognition", 12th international congress for statistics, computer science, social and demographic research, Cairo, 1987

<sup>١٥</sup> M.A. El Hamalaway & S.H. El Ramly , "A Novel Arabic Numerals Shapes for Optical Character Recognition", Thirteen International Congress for Statistics, Computer Science, Social and Demographic Research, 26-31 March 1988

المؤتمر الثاني لمنظمة اللغة، القاهرة، ١٨ إبريل، ١٩٩٩م

ولقد تبين أن الأرقام المشرقية (٠، ١، ٢، ٣، ٤، ٥، ٦، ٧، ٨، ٩) أكثر تجانساً مع أشكال الحروف العربية من الأرقام الغبارية المغربية (0, 1, 2, 3, 4, 5, 6, 7, 8, 9). كما أن الأرقام العربية المشرقية تنتمي بدرجة أكبر إلى الحضارة العربية منها إلى الحضارة الغربية؛ بينما الأرقام الغبارية المغربية تنتمي بدرجة أكبر إلى الحضارة الغربية منها إلى الحضارة العربية. وبناء على النتائج التي ظهرت من البحث فإن فئة الأرقام العربية المشرقية التي تنتمي إلى حضارة واحدة تبين أنها ليست الهندية وليست اللاتينية وليست اليونانية، بل هي الحضارة العربية وبالتالي لا يوجد سبب قوى أو ضعيف يبرر نبذ سلسلة الأرقام المستعملة في المشرق العربي واستعمال الأرقام المستعملة في المغرب العربي محلها، وذلك إضافة إلى عدة عوامل أهمها أن الأرقام المشرقية هي الأقدم في الاستعمال والشيوخ في مختلف أرجاء الأمة العربية والإسلامية واستعمالها في أكثر التراث العربي<sup>١</sup>، بالإضافة إلى كفاءتها عن المجموعة المغربية<sup>٢</sup>، زيادة على توافقها بصورة أكبر مع أشكال حروف اللغة العربية.

## ٥ المراجع:

١. محمود فهمي حجازي ومحمد يونس الحملأوى ومحمد يسرى النحاس؛ أرقامنا العربية: الأرقام المشرقية والأرقام المغربية؛ المؤتمر السنوى الثأنى لتعريب العلوم؛ القاهرة؛ ٢٠ - ٢١ مارس ١٩٩٦م
2. Duda & P. Hart, "Pattern Classification and Scene Analysis", A. Wiley Interscience Publication, New York, 1974
3. Fukunaga, "Statistical Pattern Recognition", Academic Press, inc., San Diego, 1990
٤. محمد بن موسى الخوارزمى؛ الجبر والمقابلة؛ دار الكأتاب العربى للطباعة والنشر؛ القاهرة؛ ١٩٦٨م؛ صفحة ٢٤
٥. لجنة الرياضة؛ مجمع اللغة العربية؛ الدورة الحادية والأربعون؛ القاهرة؛ ١٩٧٤-١٩٧٥م
٦. موريس شربل؛ الرياضيات فى الحضارة الإسلامية؛ جروس برس؛ طرابلس، لبنان؛ ١٩٨٨م  
صفحة ٩١
٧. قاسم على سعد؛ الأرقام العربية: تاريخها وأصالتها وما استعمله المحدثون وغيرهم منها،

<sup>١</sup> محمد يونس الحملأوى ومحمد يسرى النحاس؛ استعمال الأرقام العربية المشرقية فى تراثنا العلمى؛ المؤتمر السنوى الرابع لجمعية لسان العرب؛ جامعة الدول العربية؛ القاهرة؛ ١٥ - ١٦ نوفمبر ١٩٩٧م  
<sup>٢</sup> محمد يسرى النحاس ومحمد يونس الحملأوى؛ قياس درجة التشابه فى مجموعتى الأرقام الهندعربية؛ المؤتمر الخامس عن الحاسب الألى بين النظرية والتطبيق؛ الإسكندرية؛ ١٢-١٤ سبتمبر ١٩٩٥م  
المؤتمر الثأنى لخدمة اللغة، الحأامرة، ١٨ أبريل ١٩٩٩م

- القسم الأول؛ مجلة الأحمدية؛ العدد ٢؛ ١٩٩٨م؛ دبي
٨. محمد يونس الحملاوى ومحمد يسرى النحاس؛ الأرقام الهندعربية، دراسة مقارنة: التاريخ، السمات، الاستخدام، التقنية؛ المؤتمر الدولى السادس عن الحاسب الآلى بين النظرية والتطبيق؛ الإسكندرية؛ ٣-٥ سبتمبر ١٩٩٦م
٩. محمد يسرى النحاس ومحمد يونس الحملاوى؛ تجانس الأرقام الهندعربية مع أشكال حروف الخط العربى؛ المؤتمر الدولى السابع عن الحاسب الآلى بين النظرية والتطبيق؛ الإسكندرية؛ ٢-٤ سبتمبر ١٩٩٧م
١٠. محمد يونس الحملاوى ومحمد يسرى النحاس؛ تجانس الأرقام الهندعربية مع أشكال الحروف العربية وأشكال حروف لغات أخرى؛ ندوة الأرقام ومكانتها فى قضية التعريب؛ مجمع اللغة العربية؛ القاهرة؛ ٢٠ فبراير ١٩٩٧م
١١. محمد يسرى النحاس ومحمد يونس الحملاوى؛ بعض القياسات الهندسية والعلمية على مجموعتى الأرقام الهندعربية؛ ندوة الخطوات العملية لإقرار استخدام الأرقام العربية؛ جامعة الأزهر؛ القاهرة؛ ١ مارس ١٩٩٨م
12. Mahmoud (El Nahas) et al : "A statistical approach for Arabic character recognition", 12th international congress for statistics, computer science, social and demographic research, Cairo, 1987
13. El Hamalaway & S.H. El Ramly , "A Novel Arabic Numerals Shapes for Optical Character Recognition", Thirteen International Congress for Statistics, Computer Science, Social and Demographic Research, 26-31 March 1988
١٤. محمد يونس الحملاوى ومحمد يسرى النحاس؛ استعمال الأرقام العربية المشرقية فى تراثنا العلمى؛ المؤتمر السنوى الرابع لجمعية لسان العرب؛ جامعة الدول العربية؛ القاهرة؛ ١٥-١٦ نوفمبر ١٩٩٧م
١٥. محمد يسرى النحاس ومحمد يونس الحملاوى؛ قياس درجة التشابه فى مجموعتى الأرقام الهندعربية؛ المؤتمر الخامس عن الحاسب الآلى بين النظرية والتطبيق؛ الإسكندرية؛ ١٢-١٤ سبتمبر ١٩٩٥م

#### امتنان:

يشكر الباحثان الجمعية المصرية لتعريب العلوم؛ ص. ب ٥٣٠١ غرب مصر الجديدة، القاهرة ١١٧٧١، مصر على دعمها لجهود التعريب مما كان له أبلغ الأثر فى إنجاز هذا البحث.

LNGENG22.DOC