

Versatile Search of Scanned Arabic Handwriting

Sargur N. Srihari, Gregory R. Ball, and Harish Srinivasan

Center of Excellence for Document Analysis and Recognition
(CEDAR)

Department of Computer Science and Engineering
University at Buffalo, State University of New York

Email: srihari@cedar.buffalo.edu

Outline

- CEDAR Handwriting Analysis Systems
 - CEDAR-FOX and CEDARABIC
- Versatile Search
 - Query Types: Image and Text
- Word Spotting Algorithms
 - Word Segmentation Based
 - Holistic (Word Shape) Algorithm
 - Analytic (Character Shape) Algorithm
 - Word Segmentation Free
- Performance: Precision and Recall
- Conclusion

End-to-End Systems Developed (1983-Present)

Including corpora

1. Hand-Written Address Interpretation

HWAI

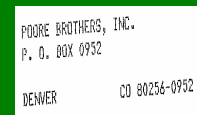
USPS, Australia Post, UK



2. Name and Address Block Reader

NABR

IRS



3. Handwriting Segmentation and Recognition

Penman

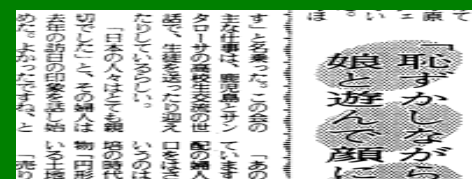
NSA



4. Japanese Character Recognition

Cherry Blossom

NSA

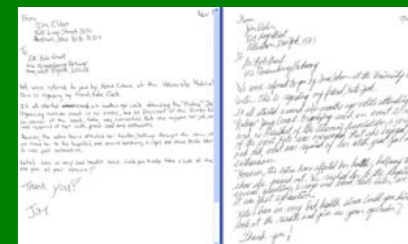


5. Writer Identification and Search

Cedar-Fox

English

NIJ



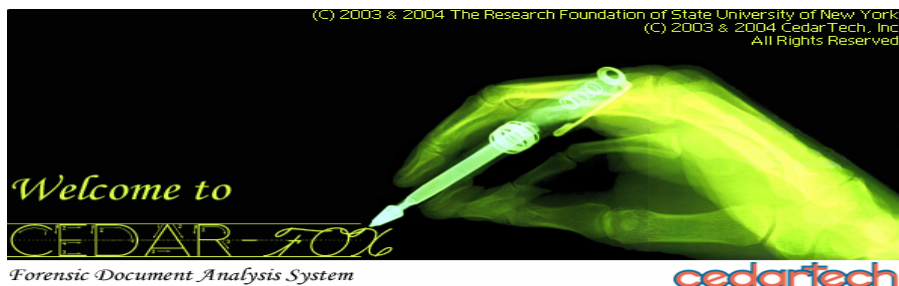
CEDARARABIC

Arabic

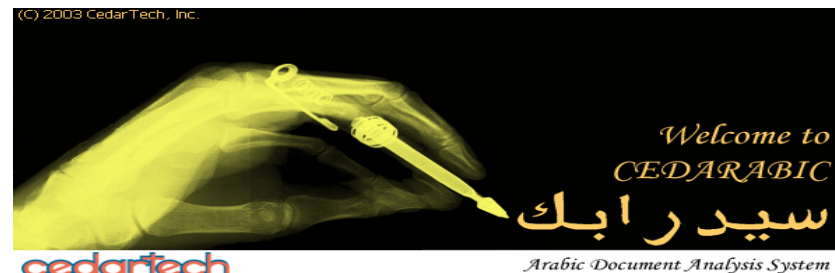


CEDAR-FOX vs. CEDARARABIC

Developed over 5 years in consultation with law enforcement agencies and professional QDEs



- English documents
- Forensic applications
- Search
 - Keyword, Image
- Recognition
 - Character, Word



- Arabic documents
- XML representation, Truthing Tools
 - Search
 - Keyword (English)
 - Image (Arabic)

- Verification and Identification
 - Writer
 - Signature
- Image Enhancement and Noise Removal
 - Two types of thresholding
 - Rule line and Underline removal
- Transcript Mapping for Creating Corpuses
- Database for Document Metadata

Search Problem

- Searching electronic documents for information related to a query is ubiquitous
- Searching scanned printed documents is a recent application
- Searching scanned handwritten images is a research frontier
 - Given a query and a repository of scanned handwritten documents, retrieve most relevant subset of documents

Approaches

- CBIR – Content-based information retrieval – broad topic in IR and data mining
 - Image based approaches—direct CBIR based on image retrieval (word spotting)
 - Text based approach—transcribe document to text and search electronic representation
- Both methods are error prone (grand challenge in computer vision)
- Combining both may achieve better performance than either alone

Versatile Search

1. Versatile query

1. Typed Arabic (eg. بمرکز)
UNICODE string of Arabic text
2. Typed English (eg. at, center)
corresponding to idea that should appear in Arabic document
3. Arabic Image (eg. بمرکز)
of Arabic word or words

2. Versatile search (combine search methods)

1. Image query (word spotting)
 - If query is image, preserve it throughout search
 - If query is text, extract or generate image query
2. Text query (needs recognition)
 - If query is image, convert to text
 - If query is text, preserve it throughout search

Image Query

ewi

Words Spotted

511

212

الحل

المالك

ewi

نوع

Database
(pre-segmented)

النقطة الرئيسية للخطاب الوطني للواء الفكري الشيخ صالح بن عبد الرحمن
الشيخ صالح بن عبد الرحمن بن فهد بن الشيخين والمختصين لتبادل الأدوار ووجهات
النظر حول اللواء الوطني، وهي قضية وطنية وذلك بمرور الأجيال.

وقد استعمل القاء بعوضها فيعلم من موطن الملك عبد العزيز الحواري الوطني
وعامر بك من موطن

بعد ذلك استعرض الشيخ عالم الوصي عبيدة مركز الملك عبدالعزيز
للدار الوطنية مؤكدا ان هذا المركز انشأ ليكون لبنة في بناء المجتمع من
علائك تدبج مفهوم الوار وثقافة الوار وعلمها فيهم عائلوا في نهضة
القرار الوطني وتكريس الوحدة الوطنية في الار القيد الانسانية رعاية
خيار اسلامي في ظل الولاية والانتقال وسامية جميع القضايا
الوطنية المختلفة في ظل الولاية والانتقال.

قدوت من معهم الزوار الوافدين في المملكة لاثنا الف إلى انه ينتقل من المناسم الاخرى وقال لاثنا الف من الاقربين وما معهموا التي في الوطية حيث انما في هذا الوطن كلنا صهرت وعلم انتموا ان يكون الاصلح منجى جاء لهم وهذه الدولة وهذا الوطن الذي تودعوا الفتيان ائمتنا قوتنا تودع في اماني الاصلح في اماني ان تكون الشريعة العرفية في الحكم الذي وان من الدومهي ان تكون حيت اعلامي.

Search Based on Text Query

CEDARABIC User Interface

Testing

Select Questioned

Select Knowns

king

Questioned

Knowns

Rank = 1, Probability 65.73 %
Meaning : king
Document : 006_1_1.arb

Rank = 2, Probability 63.79 %
Meaning : king
Document : 006_1_1.arb

Rank = 3, Probability 62.41 %
Meaning : king
Document : 008_1_1.arb

Rank = 4, Probability 61.92 %
Meaning : king
Document : 009_1_1.arb

Rank = 5, Probability 60.84 %
Meaning : king
Document : 009_1_1.arb

Page up

Page down

Page up

Page down

Tested 5 out of 5

Start Testing

English
Text
Query

Results

Word Search User Interface

Search by Text

Search String:

Search In: ☐ Current Document ☐ The Other Document ☒ Preprocessed .arb files

Use selected image for training

Go

Goto Top

Search Results

Rank	Word Image	Text	Distance	Document ID
No.1		king	0	002_1_1
No.2		king	0	002_1_1
No.3		king	0	
No.4		king	0	
No.5		king	0	
No.6		king	0	
No.7		kind	1	

1. Query (English Text)
2. Retrieved Style Choices
3. Chosen Styles
4. Results

Testing

Select Questioned

Select Knowns:

Questioned

	Rank = 1 . Probability 99.43 %
	Rank = 2 . Probability 97.91 %
	Rank = 3 . Probability 92.04 %
	Rank = 4 . Probability 86.07 %
	Rank = 5 . Probability 54.87 %

Page up Page down

Knowns

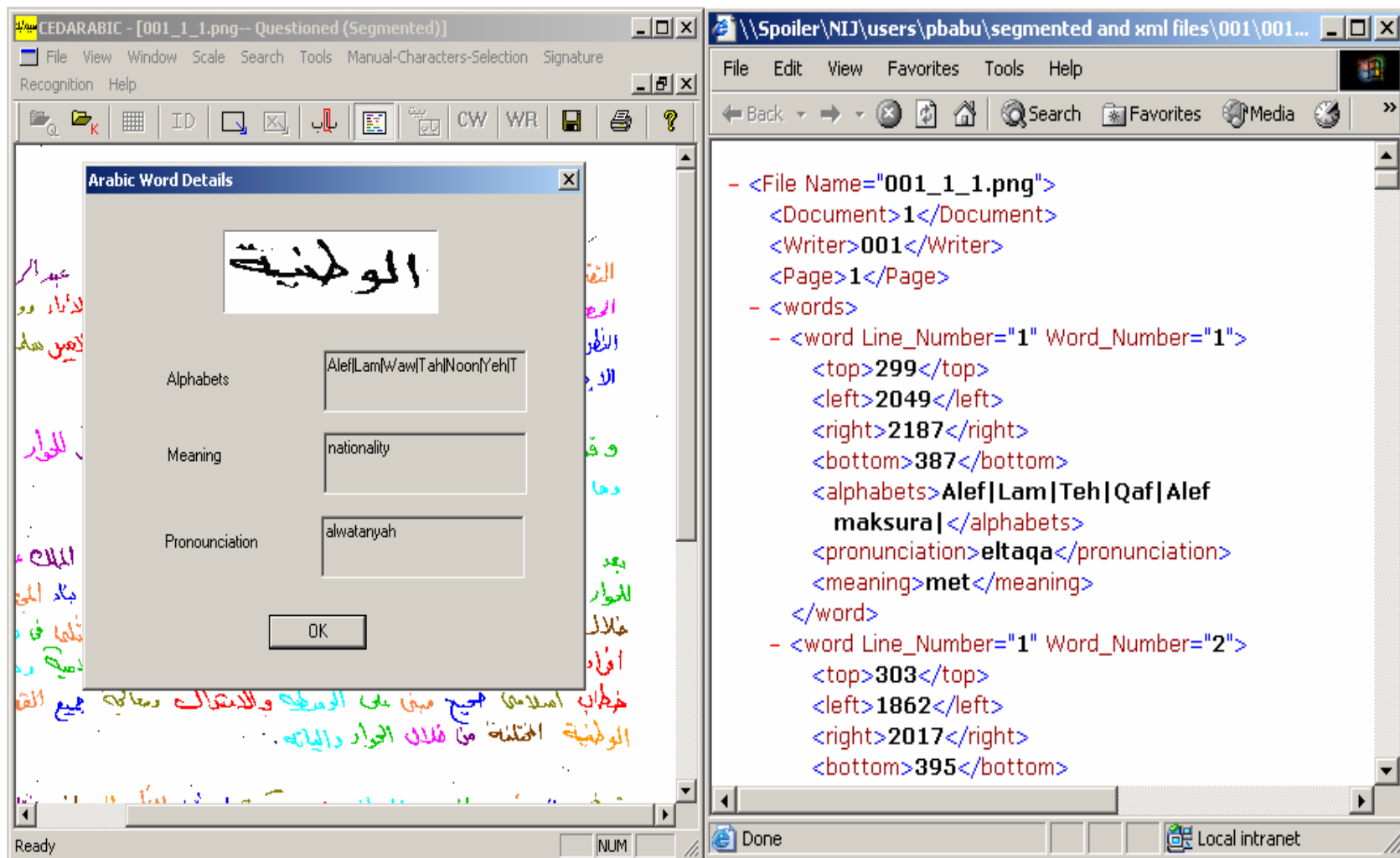
	AlefLamMeemLamKaf
	AlefLamMeemLamKaf
	AlefLamMeemLamKaf
	AlefLamMeemLamKaf
	AlefLamMeemLamKaf
	AlefLamMeemLamKaf

Page up Page down

Tested 7 out of 7

Start Testing

XML Representation



Detail of Recognition Module

Oversegment Words

Prototype Clusters

Dynamic Programming
(Maximization)

Find Nearest
Prototype

Character Based Word Recognition

Word Library

Combine Word
Library Features

Search for
Closest Match

Word Shape Recognition

Holistic Approach

Operates on Lines
(No Word Segmentation)

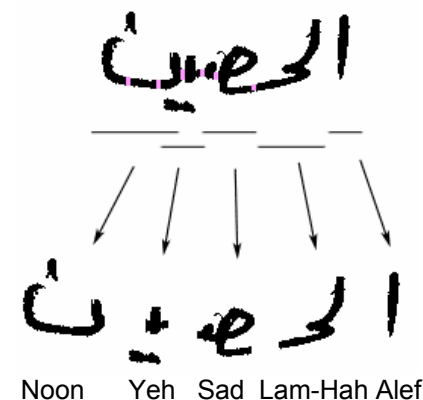
Maximize Word
Scores (for each line)

"Segmentation Free"

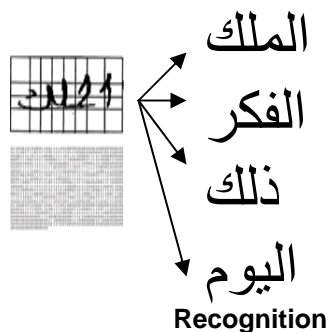
Holistic Line Recognition

Recognition

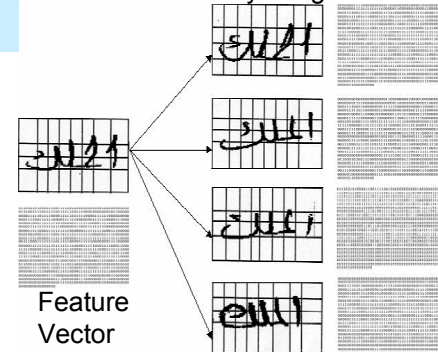
Character Based
Word Recognition



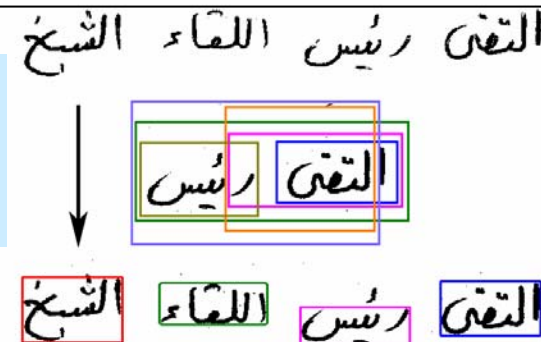
Word Shape Recognition

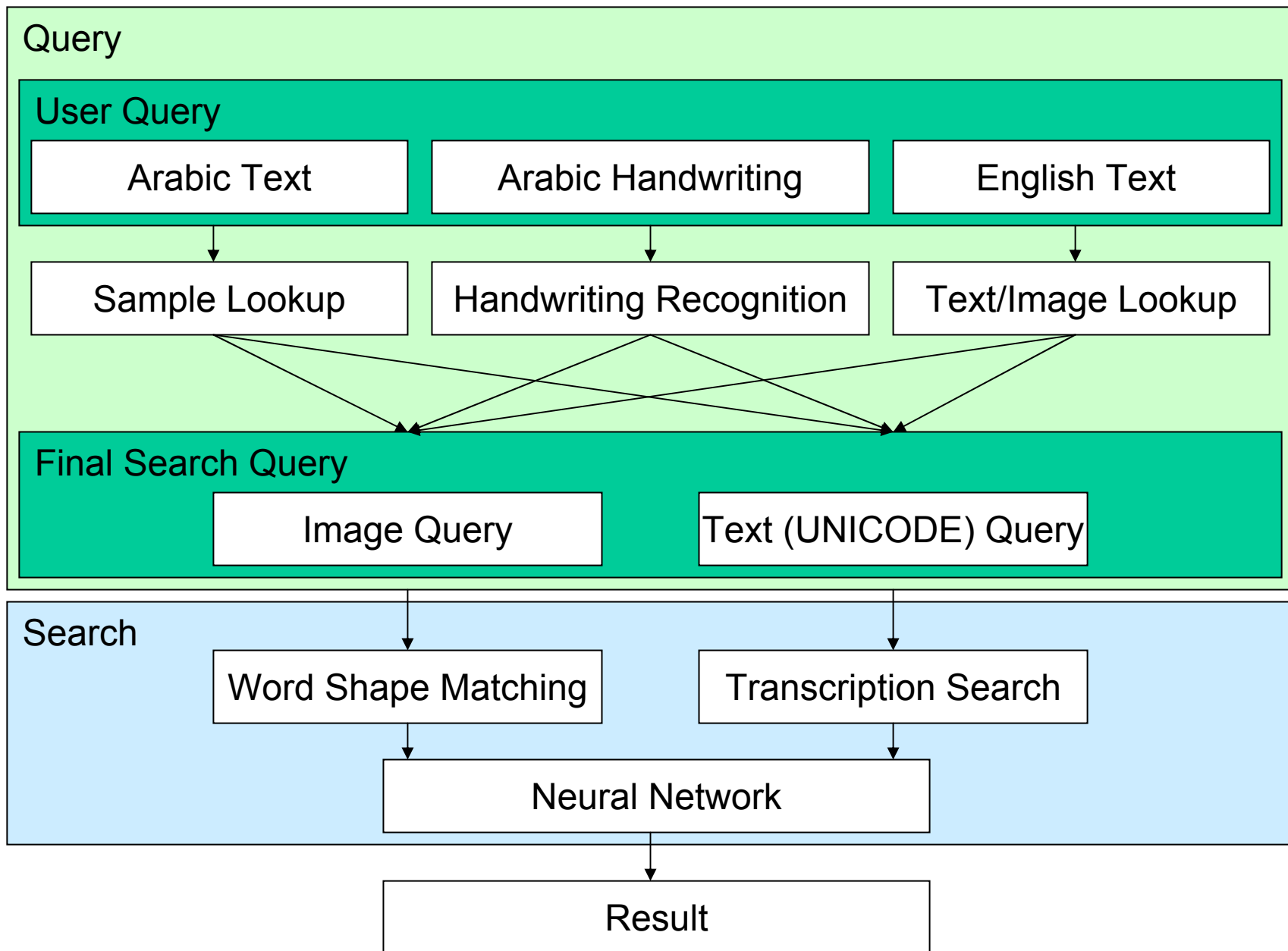


Library Images and Vectors



Holistic Line
Recognition
(Sliding Window)





Segmentation

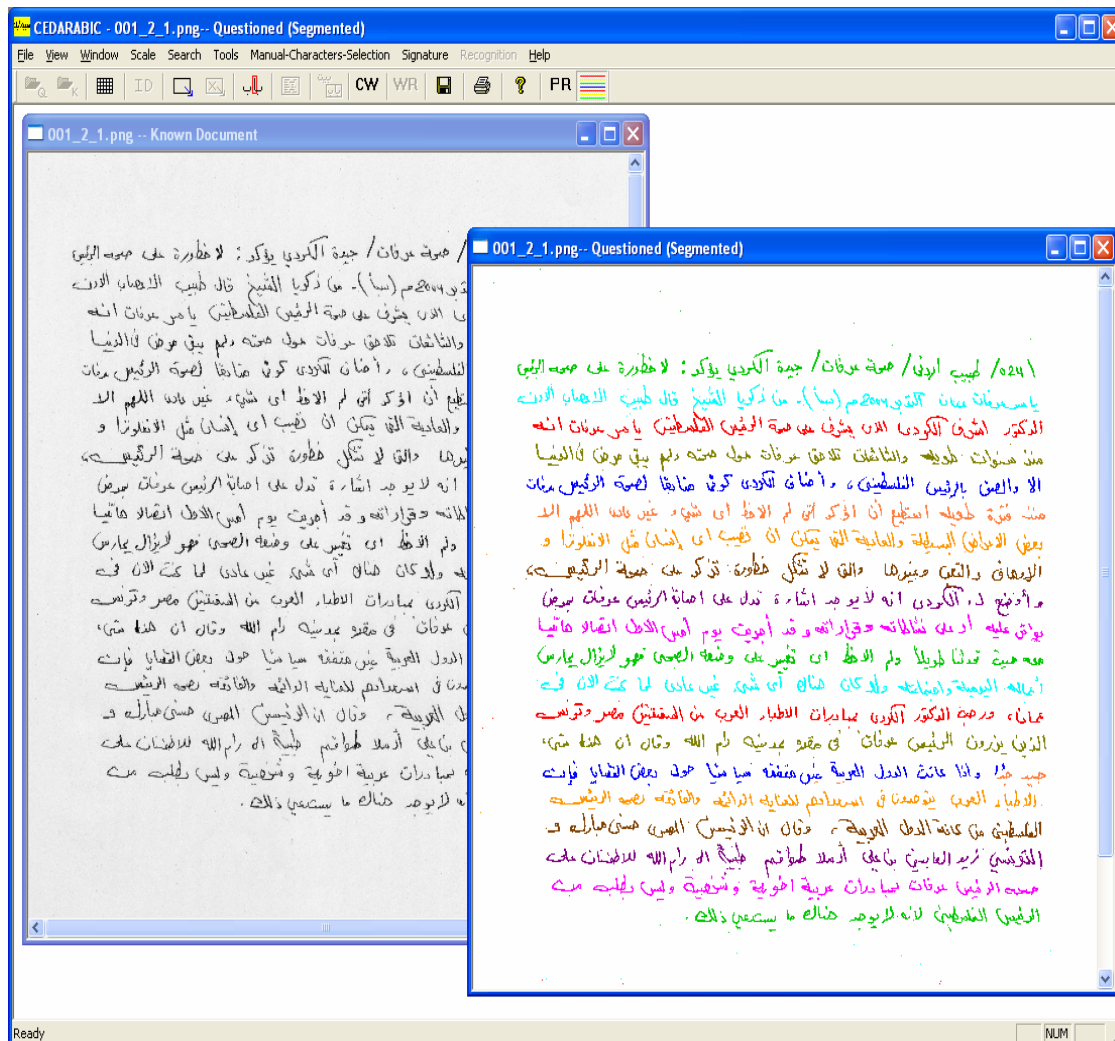
- Line – separating page into component lines
 - Most critical – new method achieves extremely successful line segmentation
- Word – separating line into component words
 - Developed automatic segmentation method;
 - Segmentation-free methods avoid need for word segmentation
- Character – separating word into component characters
 - Holistic approaches avoid character segmentation issues
 - Character based methods use prototypes to avoid need for complete character segmentation

Search depends on successful segmentation

Line Segmentation

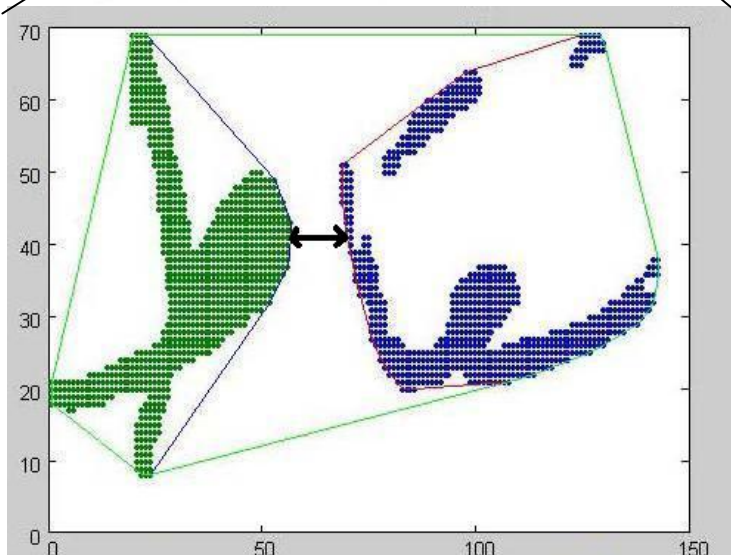
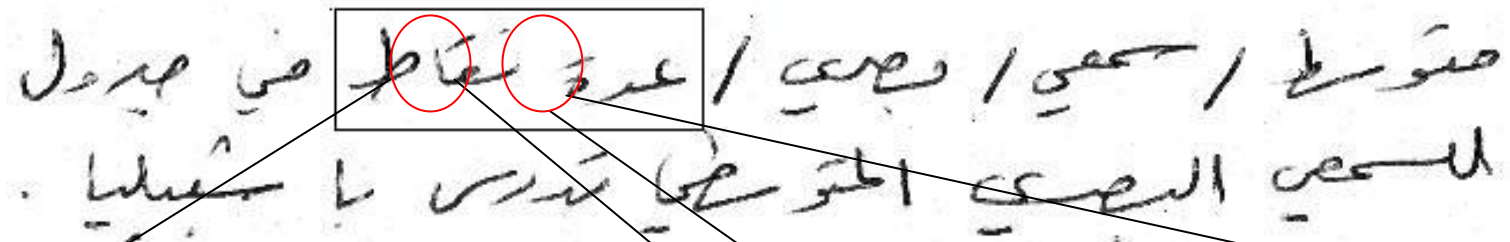
Algorithm

- Creates statistical models of adjacent lines
- In combination with top-down approaches
- To be presented at **SPIE, San Jose**
January 2006

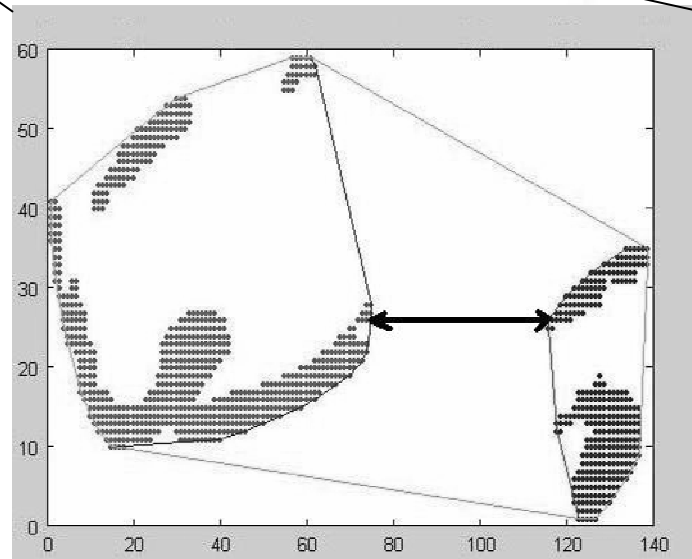


Word Segmentation

To determine whether a gap is a true word gap



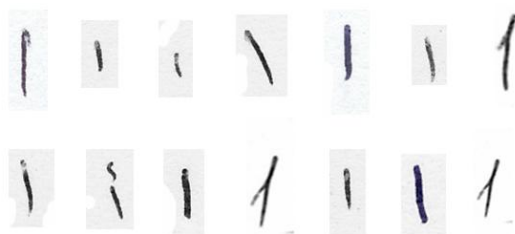
Not word gap



Word gap

Arabic Word Segmentation Algorithm

- Improved over method for Latin script segmentation
- Clustering of components
- Convex hulls of clusters
- Convex hull of pair of clusters
- Features(9)
 - Minimum distance between convex hulls
 - Ratio of area of pair to sum of individual areas
 - Heights of clusters
 - Alef Flag (words tend to begin with alef)

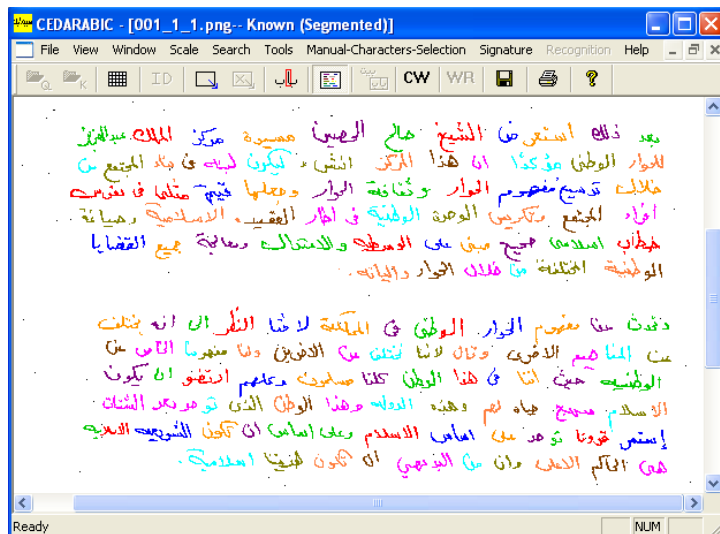


Height / width of
Components used

Word Segmentation Performance

Auto-segmentation

Truth



CEDARARABIC Word Segmentation

Automatic mode

الشيخ رئيس القمام الوطني للحوار الفكري الشيخ صالح بن عبد الرحمن
الشيخ حسام اليوم مع فيه من المتقنين والمفكرات لتبادل الآراء وجهات
النظر حول الحوار الوطني الوطفا وسيل تطويع وذلك يركز الاخص ههنا .
الاقتصاد في الرياض .

وقد استعمل القمام بعرضها فيلم عن مركز الملك عبد العزيز للحوار الوطني
وهو من مزايا .

بعد ذلك استعرض الشيخ صالح الشيخ حسام مركز الملك عبد العزيز
للحوار الوطني مؤكدا ان هذا المركز انشأ ليكون لبنة في بناء المجتمع من
خلال توسيع مفهوم الحوار وثقافة الحوار وجعلها قيم مشتركة في مجتمع
أفراد المجتمع وتكريس الوحدة الوطنية في إطار العقيدة الإسلامية ومبادئ
خطاب إسلامي صحيح مبني على الوسطية والاعتدال ومعالجة جميع القضايا
الوطنية المختلفة من خلال الحوار والحوار .

دعوت عن مفهوم الحوار الوطني في المملكة لآلها النظر الى انه يختلف
عن المناهج الفكرية وقال لآلها نشأت عن الاقوي ونا مشورا لآلها من
الوطنية حيث اننا في هذا الوطن كلنا مسلمون وعلمنا اننا نكون
الاسلام صحيح عيه لم وهذه الدولة وهذا الوطن الذي تودعنا الشان
استمع لقولنا تودعنا على اساس الاسلام وعلى اساس ان تكون الشخصية العربية
هي الحاكم الذي وان من الوديعي ان تكون لغتنا اسلامية .

Manual Mode

الشيخ رئيس القمام الوطني للحوار الفكري الشيخ صالح بن عبد الرحمن
الشيخ حسام اليوم مع فيه من المتقنين والمفكرات لتبادل الآراء وجهات
النظر حول الحوار الوطني الوطفا وسيل تطويع وذلك يركز الاخص ههنا .
الاقتصاد في الرياض .

وقد استعمل القمام بعرضها فيلم عن مركز الملك عبد العزيز للحوار الوطني
وهو من مزايا .

Useful for creating a corpus

$$score(w_i) = \frac{1}{n} \sum_{j=1}^n d(w_i, s_j)$$

$$d(X, Y) = \frac{1}{2} \left(1 - \frac{s_{11}s_{00} - s_{10}s_{01}}{[(s_{10} + s_{11})(s_{01} + s_{00})(s_{11} + s_{01})(s_{00} + s_{10})]^{1/2}} \right)$$

Chosen styles

Feature Vectors

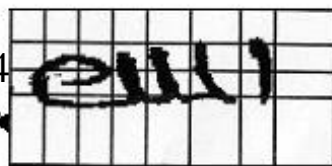
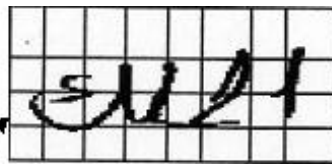
Candidate Word w_j in Database

2221

[illegible]
$$S_1$$
S₂

S₃

S


[illegible][illegible][illegible][illegible]

Spotting Based on Word Image Queries

User Interface

Latin Script-Handwriting

Search by Image


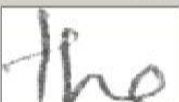
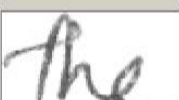
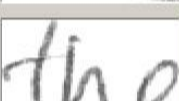
The Query Image: 

Search In:

- ☐ Current Document
- ☒ The Other Document
- ☐ Preprocessed Documents

☐ Search only in ROI

Searching results:

Rank	Word Image	Distance	Document ID	
No.1		0.212882	0001a	Highlight
No.2		0.269523	0001a	Highlight
No.3		0.298164	0001a	Highlight
No.4		0.309686	0001a	Highlight


Go

Goto Top

Go Up

Devanagari Script-Printed

Search by Image





The Query Image: 

Search In:

- ☐ Current Document
- ☐ The Other Document
- ☒ Preprocessed Documents

☐ Search only in ROI

Searching results:

Rank	Word Image	Distance	Document ID	
No.1		0	5_1	Highlight
No.2		0	5_1	Highlight
No.3		0	5_1	Highlight
No.4		0	6_2	Highlight

Go

Goto Top

Go Up

Word Image Query in English and Sanskrit

Analytic (Character Based):

Presegmentation using ligature points

- Query: UNICODE text of word
- UNICODE text mapped to positional variations of characters (initial(i), medial(m), final(f), separate positions)

Alef|Lam|Teh|Qaf|Alef maksura|

to

Alef_i|Lam_i|Teh_m|Qaf_m|Alef maksura_f|

- Candidate word is pre-segmented, based upon ligature points

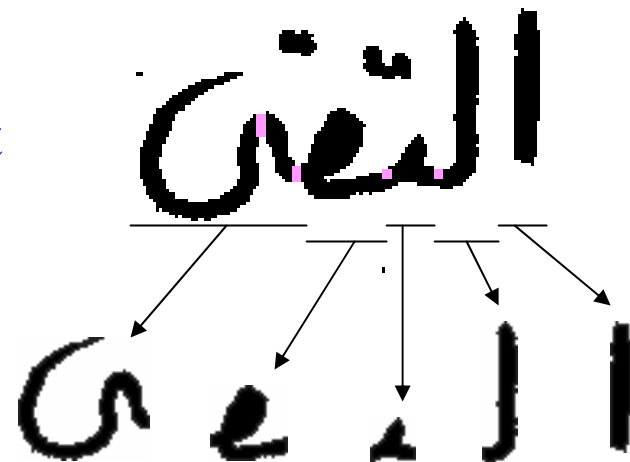
Pre-segmentation

Alef|Lam|Teh|Qaf|Alef maksura

Ligature based
segmentation of a
candidate word

Analytic (with char segmentation and recognition)

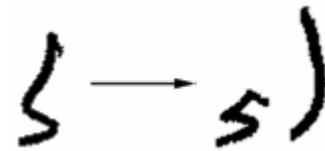
- Pre-segments reassembled into super-segments
- Candidate structures are measured against 2000 prototype chars (34 classes, 4 of each), WMR features, nearest-neighbor
- Scores of best candidate super-segments are combined into word-score
- Even with small prototype set, word to be spotted is in top 5 choices > 90% cases
- Advantage of not requiring any prototype word images



Best matching set of character super-segments

Character Based Spotting (with compound characters)

- Vertically oriented character combinations
 - Somewhat unique problem to Arabic
 - Dealt with by making compound character classes
 - Compound character classes dramatically improve recognition



Lam-ha

Ha Lam

Word-Segmentation Free Method

- Uses query to evaluate each potential word grouping
- Utilizes sliding window
 - Recognition and segmentation performed concurrently
 - Entire line acts as input
 - Splits line into connected component groups
 - Ligature based segmentation can further split components
 - Considers all realistic combinations of adjacent connected components

Candidate
Segmentations



Segmentation Free Method

الشيخ

اللقاء

رئيس

التقى

- Top 1 scoring regions for following text:
 - Alef|Lam|Teh|Qaf|Alef maksura|
 - Reh|Yeh+hamza|Yeh|Seen|
 - Alef|Lam|Lam|Qaf|Alef|Hamza|
 - Alef|Lam|Sheen|Yeh|Khah|

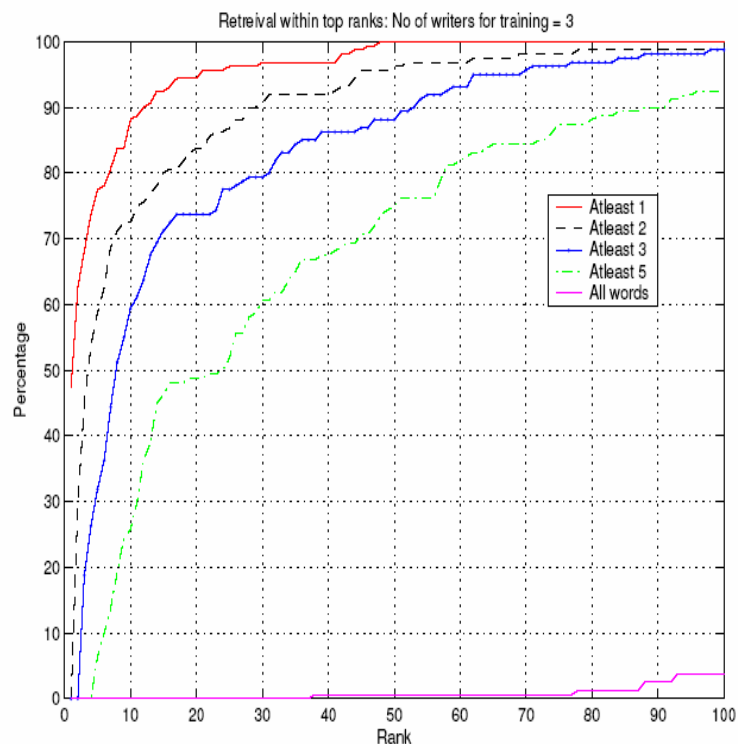
Combining Results

- After parallel image and text search, results combined with neural network
- Input: Output from each of the searches; optionally a set of features of the images
- Output: A combined score

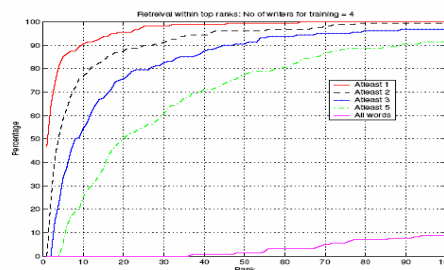
CEDARARABIC Word Spotting Performance

- Averaged over 150 Queries chosen randomly among: advancing, african, aims, algeria, algerian, allah-god ,am ,america, american ,ar, arabian, asian, atalanta, barcelona, because we, brescia, building,built, established, copeam, cagliari, cairo, chievo,country, department, developing, different views, european, existence, fiorentina,france, french, friday,gmt,gaza, germany, getting worse, gunmen,history ,influenced, intellectual, iran, iranian, iraq, islam, islamic, israel, italian, japanese, juventus, ke, khan younis (city), khartoum, lazio, lecce, legates, etc

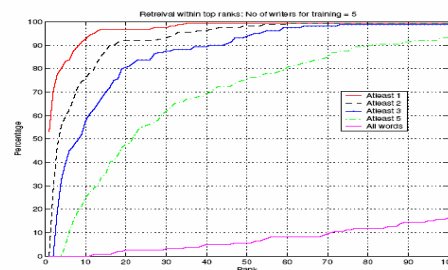
Styles = 3, Testing on 7 Writers



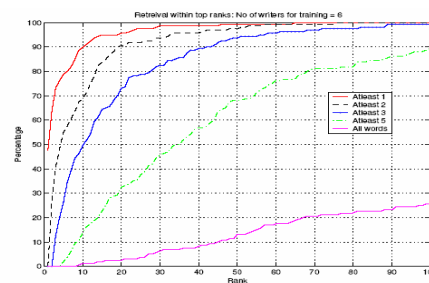
Performance increases with more styles



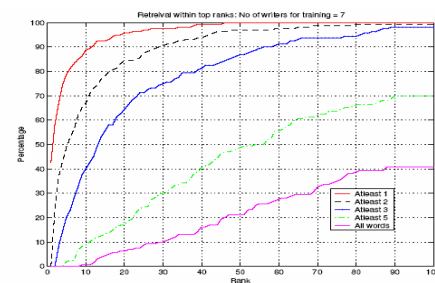
(a) 4 writers for training.



(b) 5 writers for training. vector



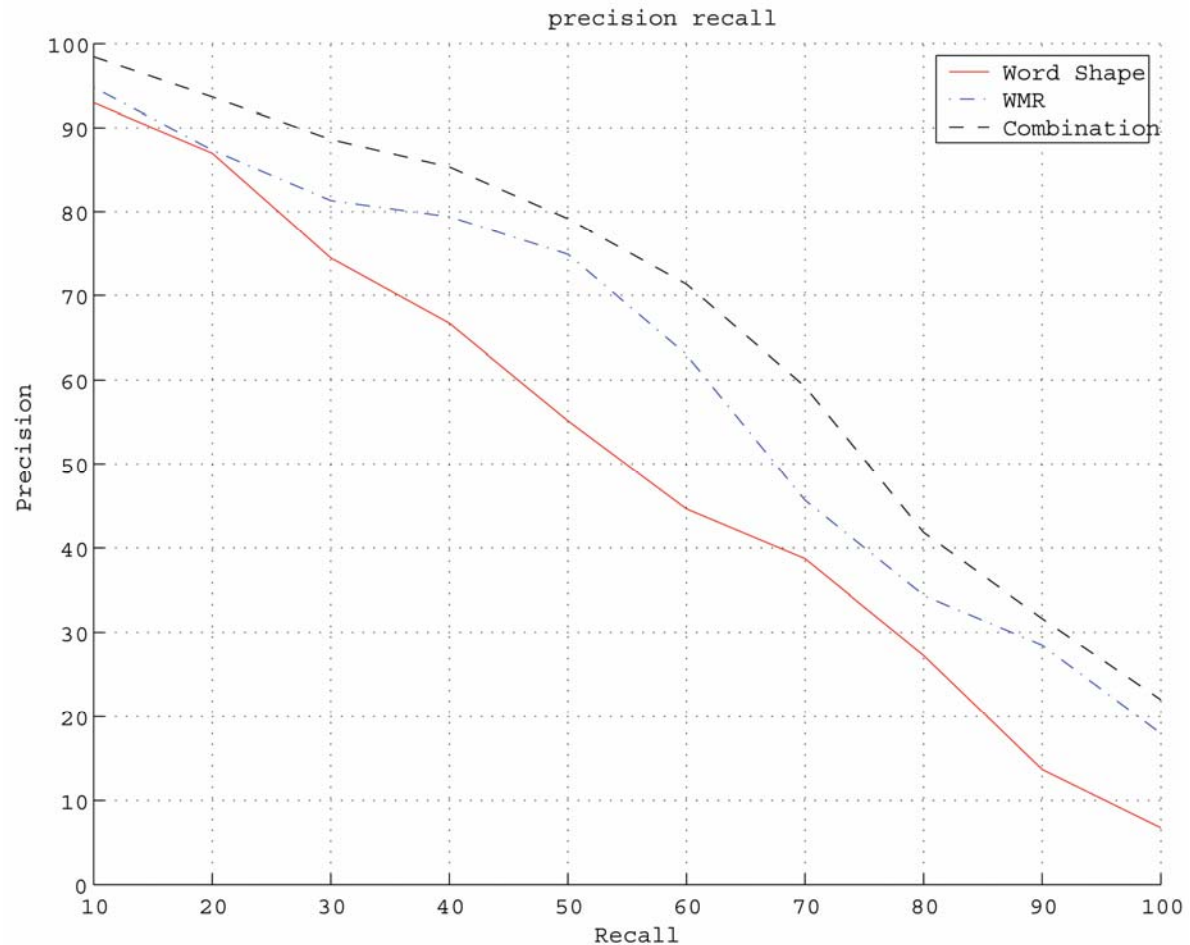
(a) 6 writers for training.



(b) 7 writers for training. vector

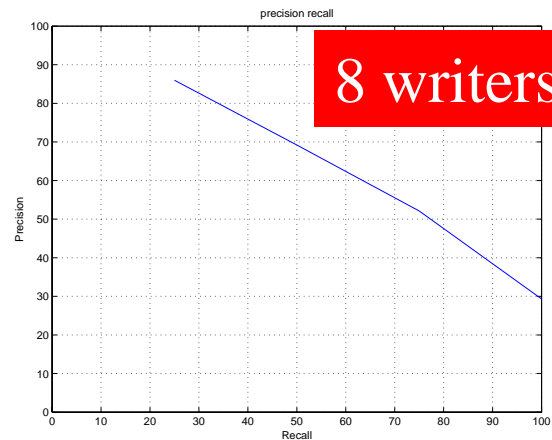
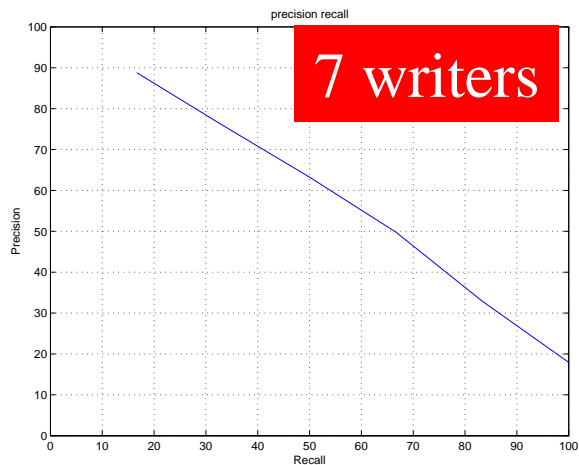
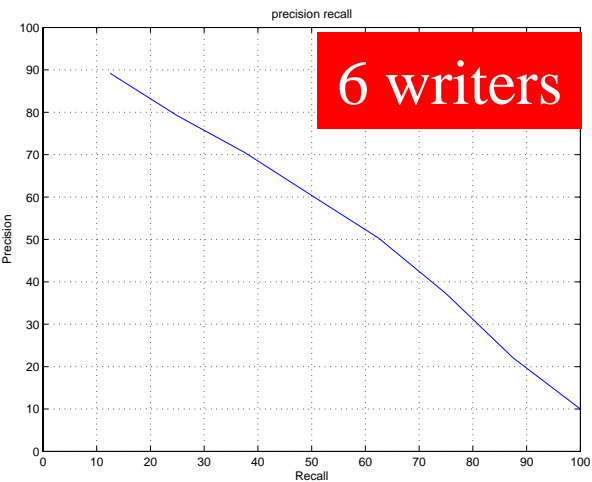
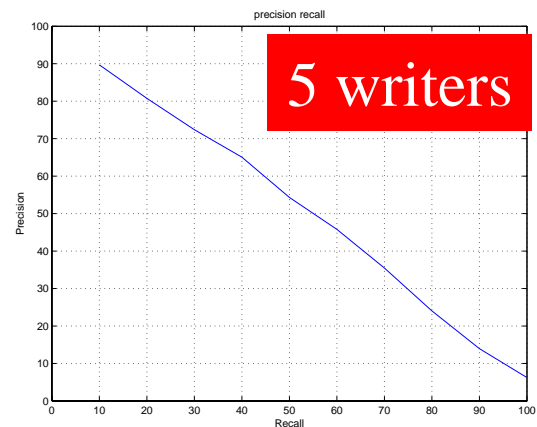
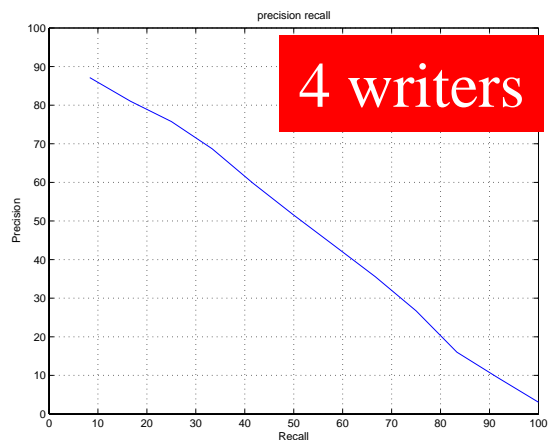
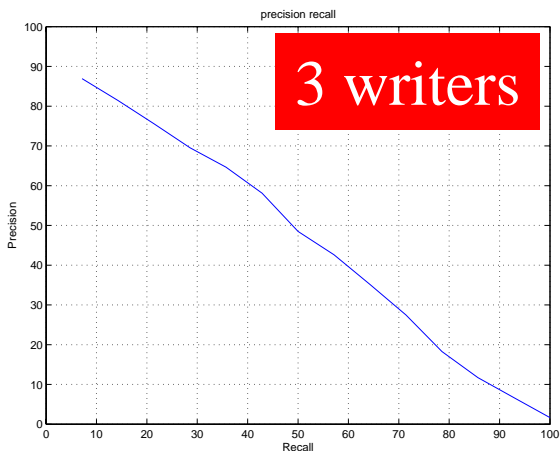
Results

- Higher performance than either method alone
- 91% raw classification accuracy
- At 50% recall, 55% precision was obtained in the word shape method, 75% precision for character based method
- Combined method about 80%

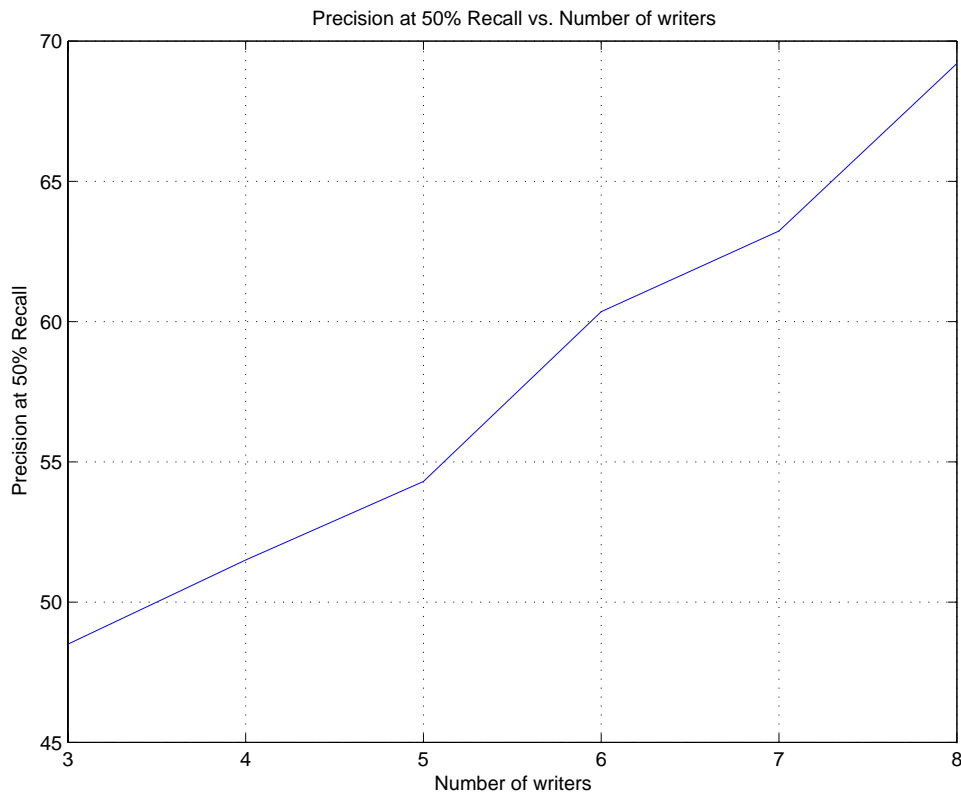


Word Spotting Precision-Recall

150 queries (king, nation, Friday,...)

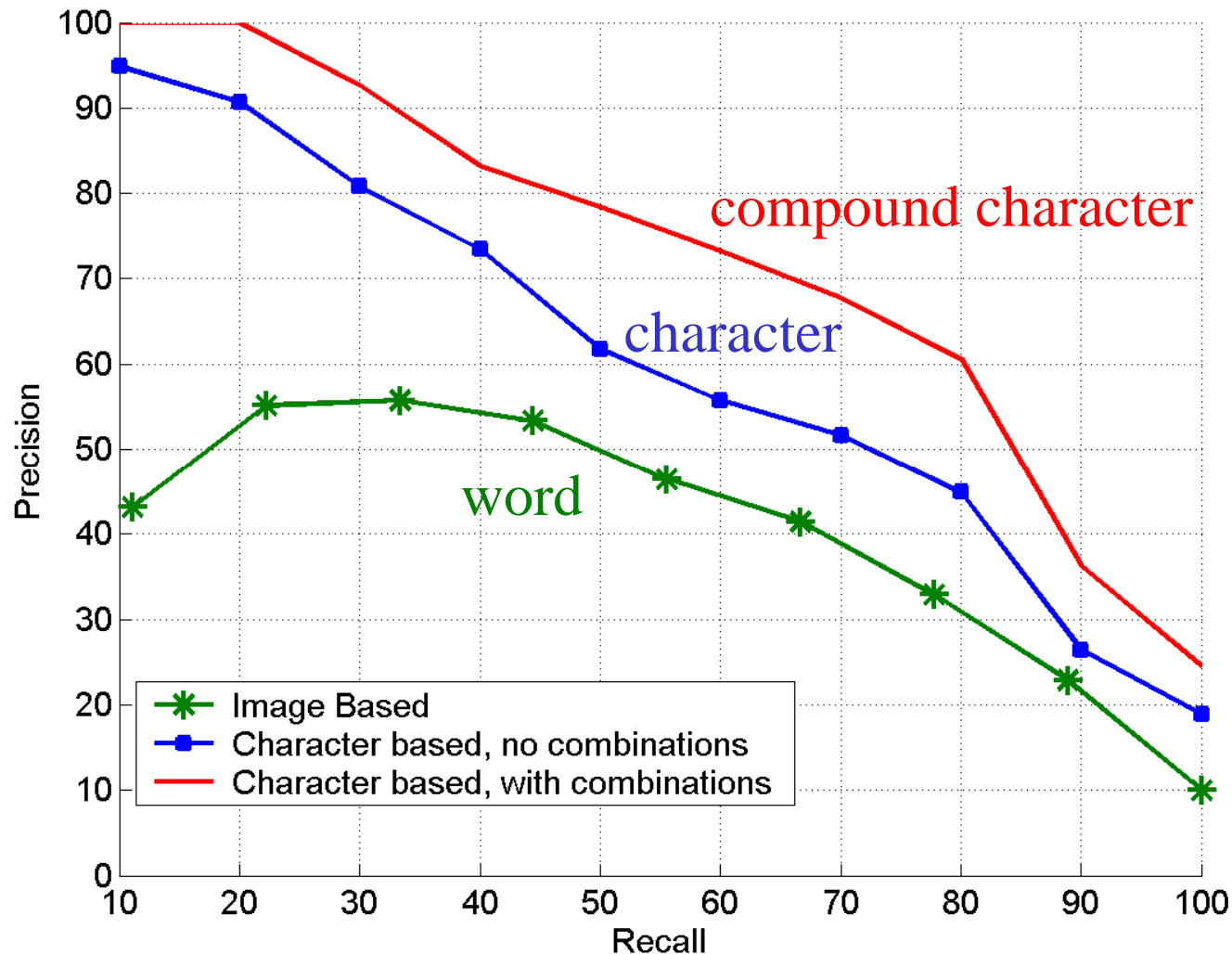


Performance as No of Styles Increase



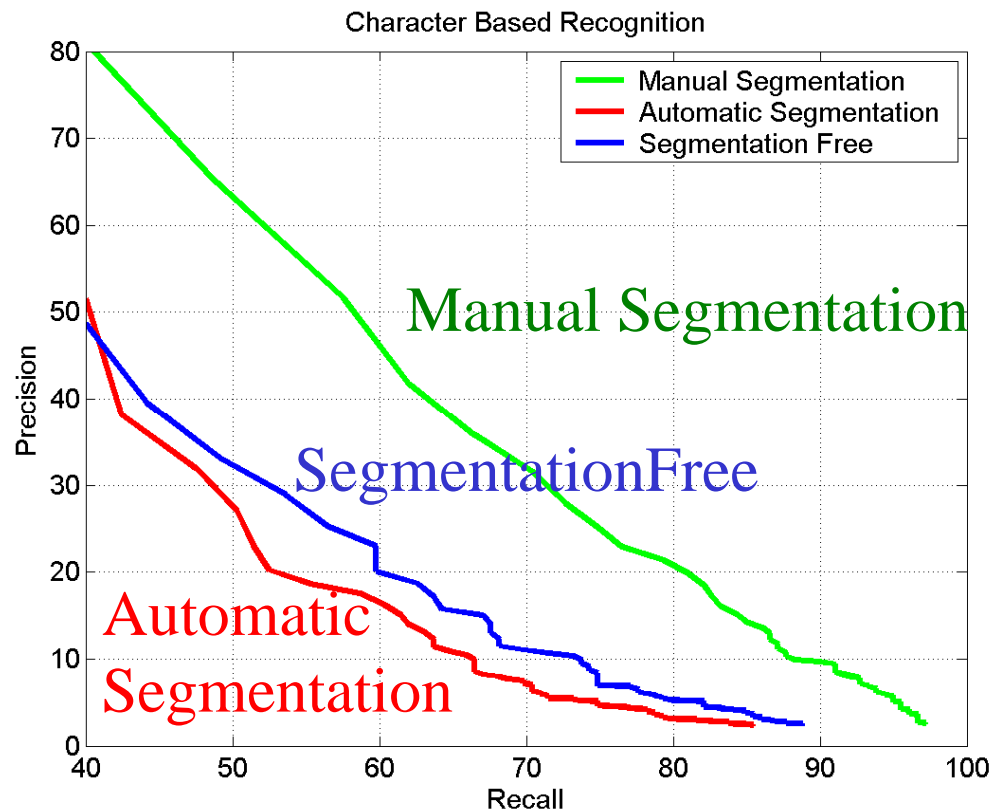
Precision at 50% Recall
vs Number of Writer Styles

Character Based versus Word Based



Performance of Segmentation Free – Character Based Method

- Comparison of manual, automatic, and segmentation free methods
- All use character based recognition; manual segmentation represents “ideal” recognition
- Segmentation free method offers significant performance increase over automatic segmentation
- Additional performance available by combining automatic/segmentation free method



Time comparison

- Methods compared on 200 word document, times in seconds on Pentium 4 (2.8 GHz)
- Overhead can be cached or preprocessed/stored before executing queries.

Method	Overhead	Per Query
Word Shape based	4	0.5
Character based	1	0.6
Word Segmentation Free	1	1.2 - 4

Summary

- CEDAR systems and corpuses
 - Developed over 25 years
 - Postal, IRS, Penman, Japanese, Indic, Forensic, Arabic
- CEDARABIC is an end-to-end system with user interfaces for:
 - Search based on keywords, writership, database functionality
 - Image enhancement, ROI selection, Transcript mapping

Summary

- Two methods for dealing with unsegmented lines
 - New method of automated word segmentation introduced for Arabic
 - Improved performance over Latin script segmentation
 - Segmentation free method
- Three methods of word spotting
 - Word based
 - Performance increases with no of styles chosen in search query
 - Character based
 - Character based with compound characters

Conclusions/Future Directions

- Processing image and text based queries in parallel can result in higher performance than either alone
- Versatile search framework can be applied to many search problems
- Using improved image or text-based search algorithms can push overall performance higher