

Writer Verification of Arabic Handwriting

Sargur N. Srihari and Gregory R. Ball

Center of Excellence for Document Analysis and Recognition (CEDAR)
University at Buffalo, The State University of New York, USA

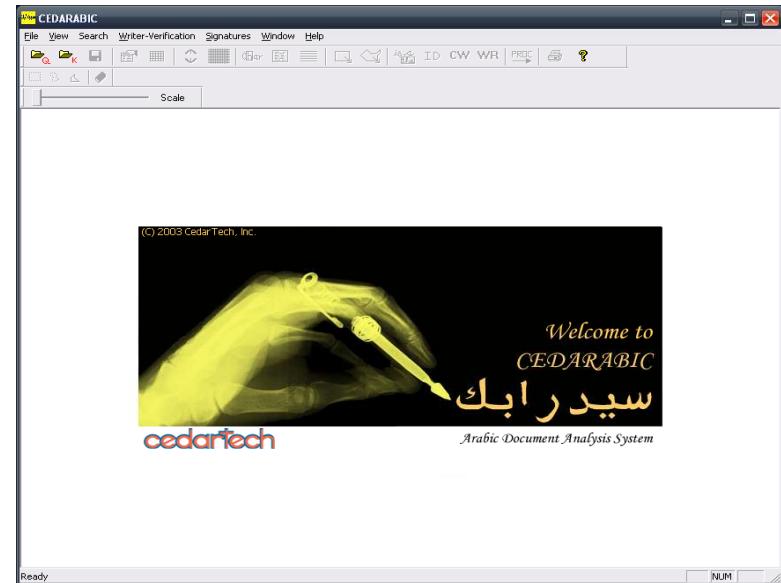
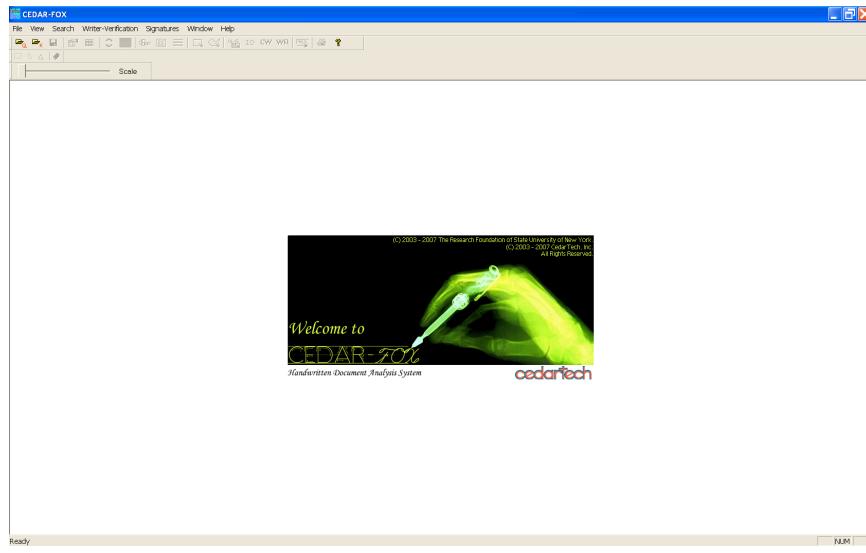
Document Analysis Systems (DAS), Nara, Japan, September 16-19, 2008

Overview

- Transition from Latin script to Arabic script
- Arabic Writer Verification Algorithm
 - Features and Similarity Measures
 - Statistical Model
 - Log Likelihood Ratio
- Performance
- User Interfaces
- Summary
- Future Plans

Transition from Latin to Arabic Script

Two systems: Cedar-Fox and Cedarabic



Functionalities: Writer Verification, Identification
Signature Verification
Word Spotting
Image Pre-processing (line/noise removal)

Latin vs Arabic Handwriting

English

It all started about six months ago while attending the "Rubeg" Jazz Concert. Organizing such an event is no picnic, and as President of the Alumni Association, a co-sponsor of the event, Kate was overworked. But she enjoyed her job, and did what was required of her with great zeal and enthusiasm.

Arabic

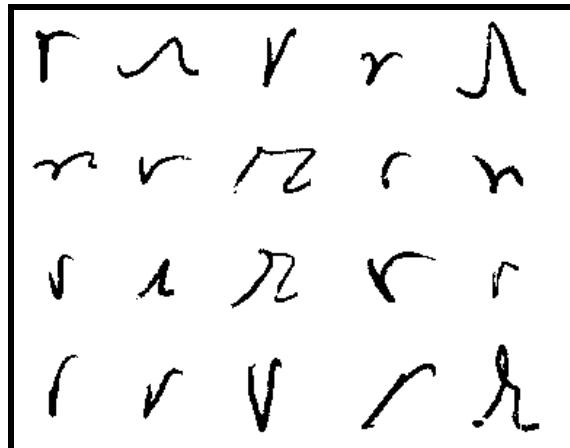
النقاش بين اللقائد الوطنيين والوزير الفكري الشيخ صالح بن عبد الله بن الحسين
صادر اليوم مع تجربة من الثقافيين والمتخصصين للتواصل والتبادل الحر حول الموارد
الوطنية وسبل تطويرها وذلك ضمن الامم المتحدة للتعاون الاجتماعي في
الرياض .
وقد استهل اللقاء ببرهان ضياع عاصم مركز الملك عبد العزيز للدراسات الوطنية
وطارئه من مراحل .

- Arabic is mostly cursive (running) writing
 - Need appropriate segmentation for characters
- Written right to left
 - Line/word segmentation, recognition algorithms need to handle reversed direction
- More character classes in Arabic
 - (approx. 200 vs. 52 character classes)

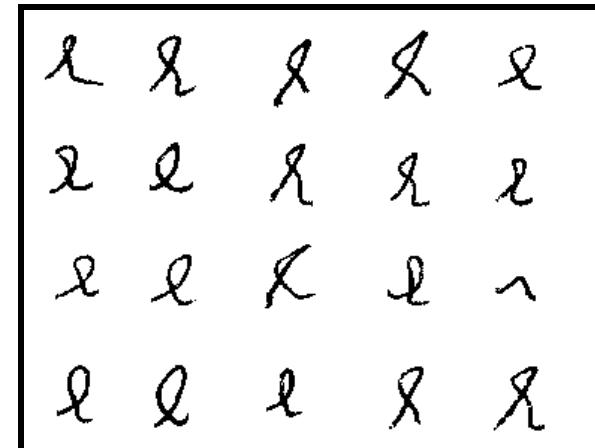
Writer Verification Principles

- “No two writers write the same and no writer writes the same twice”
- Variance of individual writers less than variance of general population
- Latin script example

General population: letter “r”



Single writer: “r”



Writer Verification Principle

- Arabic Samples of Two Writers

Writer A

السبعين (السبعين)
الفارس (الفارس)
الوطني (الوطني)
اللغة (اللغة)

Writer B

السبعين (السبعين)
الفارس (الفارس)
الوطني (الوطني)
اللغة (اللغة)

– *alef-lam* combinations have discriminatory power

- Difference can be used as feature for binary classification
 - Cumulative difference over all characters is an indicator of writership

Writer Verification Algorithm

- **Input:** Two scanned writing samples
 - same or different content
- **Output:** Probabilistic measure of same/different writer
 - Log-Likelihood Ratio (LLR)
 - Discretized Opinion on a scale
- **Process:**
 - 1 Extract macro (global) features
 - 2 Extract micro (character) features
 - 3 Compute similarities for macro and micro features
 - 4 Compute class probabilities from statistical model

1. Macro (Global) Features

- Global writing characteristics
 - Extracted from entire document
- Some features easily computed
 - 1 stroke width
 - 2 average word gap
 - 3 Entropy of gray values
 - 4 binarization threshold
 - 5 number of black pixels
 - 6 contour slope components
 - 1 horizontal (0 degree or flat stroke)
 - 2 positive (45 or 225 degree)
 - 3 vertical (90 or 270 degree)
 - 4 negative (135 or 315 degree)
 - 10 average height (per line)
 - 11 average slant (per line)
 - 12 number of interior contours
 - 13 number of exterior contours
- Features 3 and 4 not used in experiments

2. Micro Feature Extraction Algorithm

2.1 Word Segmentation

– from lines of text

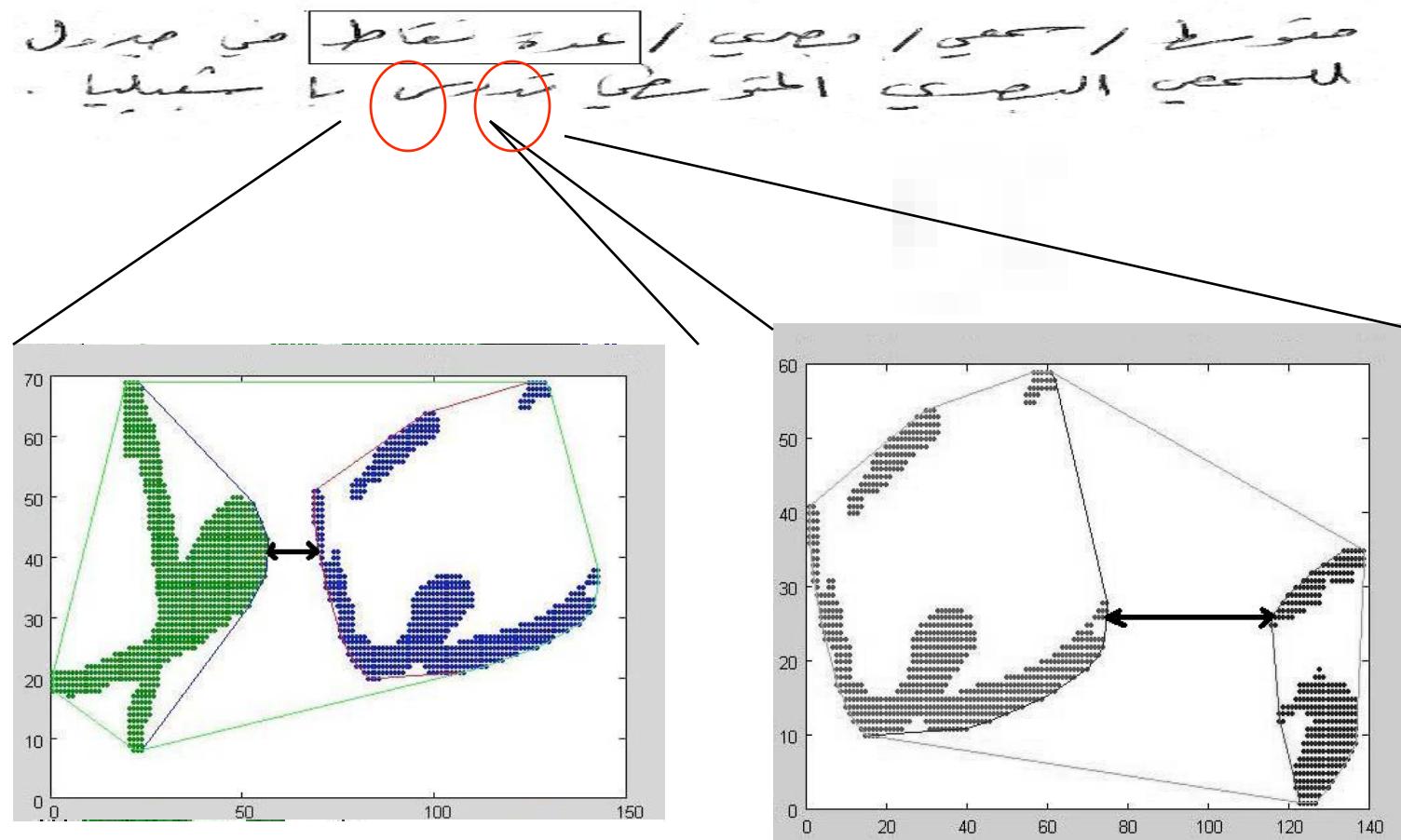
2.2 Character Segmentation

From words of text using word truth

2.3 Character Feature Measurement

2.1 Word Segmentation

- Determining whether a gap is a true word gap



Not word gap

Word gap

Features for Word Segmentation

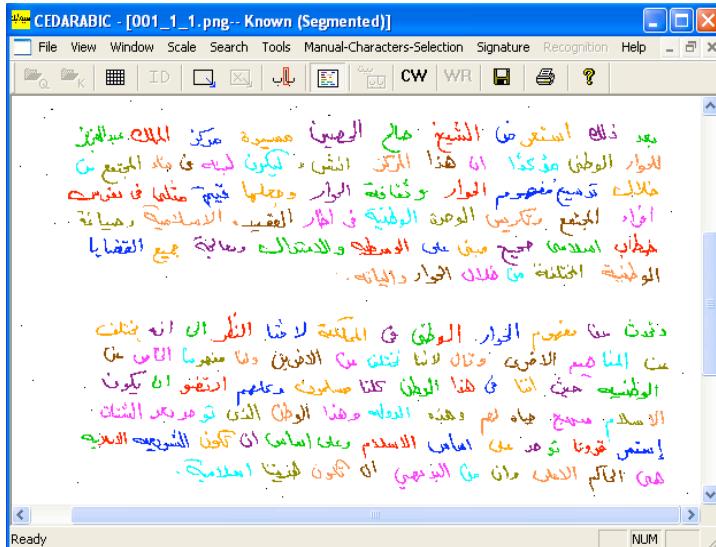
- Additional features from Latin script segmentation
- Clustering of components
- Convex hulls of clusters
- Convex hull of pair of clusters
- Features(9)
 - Minimum distance between convex hulls
 - Ratio of area of pair to sum of individual areas
 - Heights of clusters
 - Alef Flag (words tend to begin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



Height and width of components are used

Word Segmentation Result

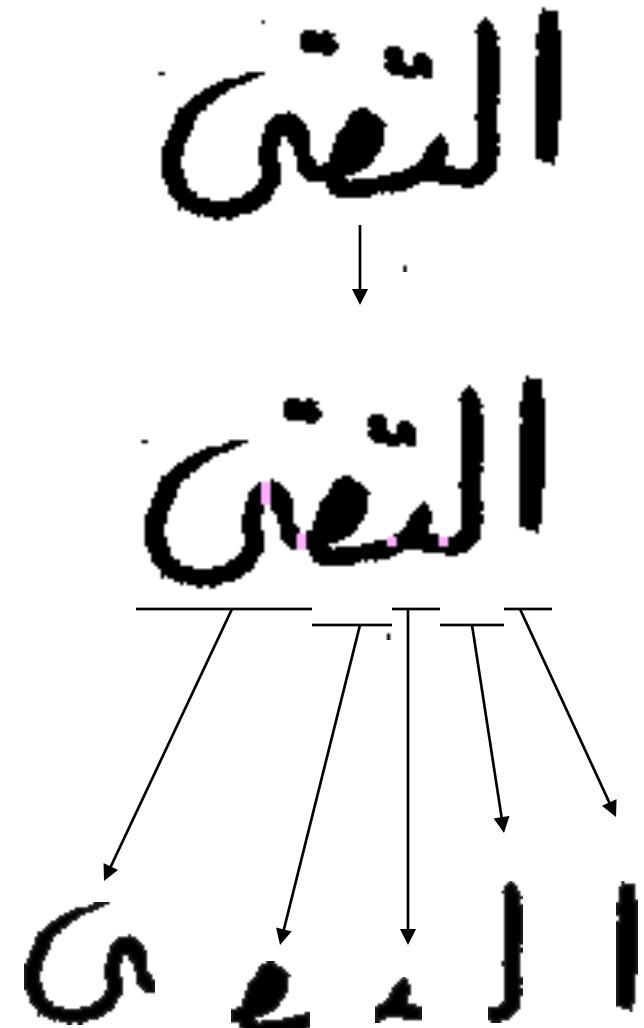
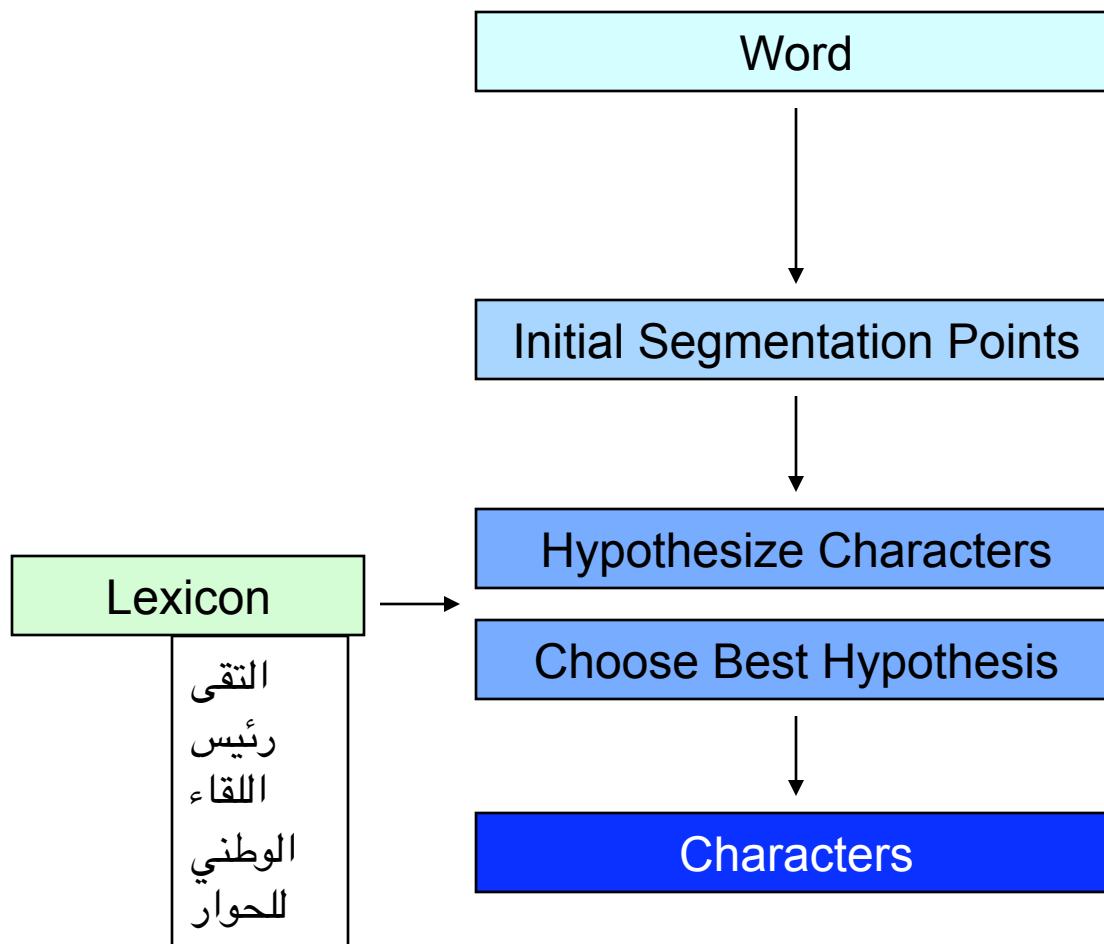
Auto-segmentation



Truth

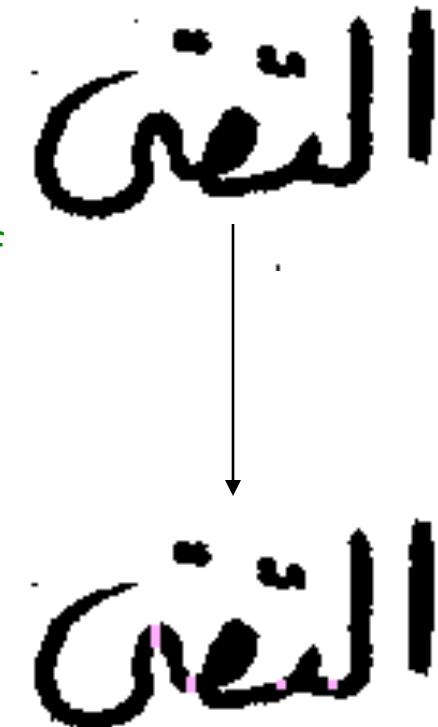


2.2 Character Extraction



Arabic Specific Character Segmentation

- Input:
 - candidate word image and UNICODE text of word
- Mapping:
 - UNICODE text is mapped to positional variations of characters (taking into consideration initial, medial, final, separate positions)
Ex) Alef|Lam|Teh|Qaf|Alef maksura|
to Alef_i|Lam_i|Teh_m|Qaf_m|Alef maksura_f|
- Candidate word is segmented, based upon ligature points
- Goal is to avoid undersegmentation
 - i.e. words will at least be segmented into component characters, with the possibility of component characters split into multiple segments

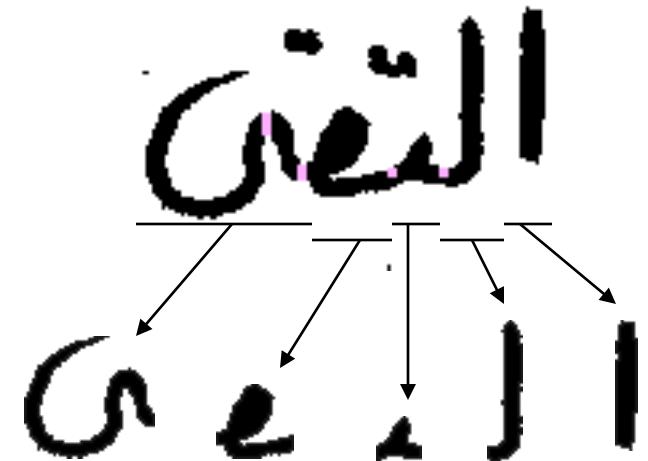


Ligature based
segmentation of a
candidate word

Summary of Character Segmentation

- Segment
 - word at ligatures
- Reassemble
 - Segments into super-segments
- Measure candidate structures
 - Against prototype chars, features, nearest-neighbor
- Obtain word score
 - Scores of best candidate super-segments are combined

Initial segmentation Points



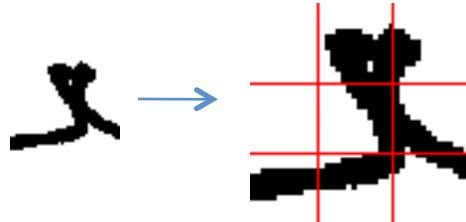
Best matching super-segments

2.3 Arabic Character Feature Extraction

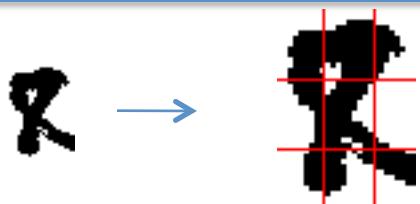
- Vertically oriented character combinations
 - Somewhat unique problem to Arabic
 - Dealt with by making compound character classes
 - Compound character classes dramatically improve recognition

$\{\rightarrow s\}$

Character Feature Set (Micro)



0.7773	0.8773	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
0.5800	0.2400	0.0960	0.1733	0.0000	0.3840	0.4000	1.0000	0.3718	0.4615		
0.7385	1.0000	0.2231	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.2417	1.0000	
1.0000	0.0000	0.2014	1.0000	1.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	
1.0000	0.0000	0.0000	0.2989	1.0000	0.0000	0.2759	0.1724	0.6265	0.0000		
0.0000	0.4815	0.9600	6667.0	0.1778	0.1481	0.7407	0.2589	0.2143	0.6000		
0.4643	0.6214	0.5143	0.0000	0.0000							



0.1968	0.3275	0.7738	0.2500	0.0000	0.0714	0.0000	0.2857	0.3482	1.0000		
1.0000	0.0000	0.0000	0.1538	0.0000	0.6154	0.1250	0.7179	0.3095	1.0000		
0.9400	2857.0	0.2857	0.0000	0.2857	0.0000	0.0000	0.0000	0.2188	0.8250		
0.0000	0.2750	1.0000	0.7109	0.0000	0.1444	0.7000	0.8800	0.2000	0.0000		
0.0000	0.2167	0.9333	0.0000	0.0000	0.0000	1.0000	0.0000	0.0000	0.0000		
0.0000	0.0000	0.0000	0.0000	0.0000	0.6769	0.3077	1.0000	0.7179	0.1083		
0.3500	0.8800	0.2500	0.4400	1.0000	0.0000	0.2333	0.3939	0.6364	1.0000		
0.0000	1.0000	0.7273	0.0000	0.0000							

2 Global Features: Aspect ratio
Stroke ratio

72 Local Features

$$F(i,j) = \frac{s(i,j)}{N(i)*S(j)} \quad i = 1,9 \\ \quad \quad \quad j = 0,7$$

s(i,j) - no. of components with slope j in sub image i

N(i) - no. of components in sub image i

$$S(j) = \max \frac{s(i,j)}{N(i)}$$

3. Similarity Measure for Characters

- Correlation measure as used to compute dissimilarity between two binary feature vectors.
 - Distance, $D(X, Y)$ between the two feature vectors X and Y is given by the formula:

$$\frac{1}{2} - \frac{S_{11}S_{00} - S_{10}S_{01}}{2\sqrt{(S_{10} + S_{11})(S_{01} + S_{00})(S_{11} + S_{01})(S_{00} + S_{10})}}$$

- Where S_{ij} is number of occurrences of matches with i in the first vector and j in the second vector at the corresponding positions

4. Statistical Model

- Two classes: same and different
- Class-conditional pdfs of similarities for feature x :

$p_s(x)$ and $p_d(x)$

– For macro features, Gamma for both classes

- Will not overlap negative part of real line

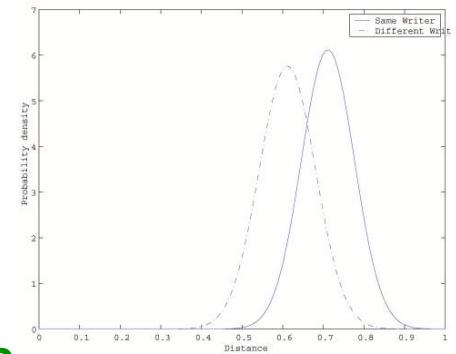
$$p_s(x) \sim \text{Gamma}(\alpha_s, \beta_s)$$

$$p_d(x) \sim \text{Gamma}(\alpha_d, \beta_d)$$

– For micro features, Gamma and Gaussian

$$p_s(x) \sim \text{Gamma}(\alpha_s, \beta_s)$$

$$p_d(x) \sim N(\mu_d, \sigma_d^2)$$



Maximum Likelihood Formulation

- Current system based on a maximum likelihood formulation
 - A Bayesian formulation is planned
- Assume statistical independence of features

$$LLR = \sum_{i=1} \sum_j \sum_k \ln p_s(d_i(j, k)) - \ln p_d(d_i(j, k))$$

Arabic Dataset

- Used for testing
 - Approximately 10 one page articles written by approximately 10 authors (100 full page documents)
 - Allows 5,000 pairwise comparisons
 - Realistic (page images)
 - Manually truthed to word level

العن رئيس القياده الوطنى للجوار العقاري الشيخ صالح بن عبد الرحمن
الصبيح مساداً اليوم مع فمه من المتفقين والمحفقات لتبادل الآراء وجهات
النظر حول الجوار الوطنى وسبل تطويره وذللك يرسّخ الاعتزاز بالعلماء
الباحثين على في الراحله .

وقد رستم القاء يعروها فيلم عن حكم الملك عبد العزيز للحوار الوطني
وهما يموي به من مهتماً

بعد ذلك استقرت الشيعة حال الهدى محسنة مركز الملك عبد العزيز لل WAR الوطنى مؤكداً أن هذا المركز انشئ ليكون لبيه هاد المجتمع من خلال توسيع مفهوم الwart وتنافه الwart وفعاليتها قيمـ مثلـها في نفس اداء المجتمع وتكتيس الوطـة الوطنـة في إطار القـيـدـ الاسلامـيـ وهـيـاـ طـرـدـ اـمـلاـمـهاـ حـيـحـ مـفـقـدـ علىـ الوـسـطـيـ وـالـمـنـدـالـ وـمـعـاـيـرـ جـمـيعـ الفـقـاهـيـاـ الوـطـنـيـةـ المـلـتـلـةـ منـ فـلـلـ الـwartـ والـمـائـةـ.

دوقت سدا عقورم الوار الوطن في المملكة لافتة الفار الى انه يختلف
عن المذاهب الافريقية وحال لاشا يختلف عن الاقريبيين ولما نفهم ما الايام عن
الوطنيه حيث اذنا في هذا الوطن كلنا مسلمون وعلمنا ارضضوا ان يكون
الاسلام هديج لهم وهذه الدواعيه وهذا الوطن الذي توعد به الشتات
يسعني قورنا توعد بـ اهانات الامم وعلي امامنا ان تكون المؤسسه الاسلاميه
هي الحاكم الالهي وان من الدهري ان تكون هرثها اهلاما .

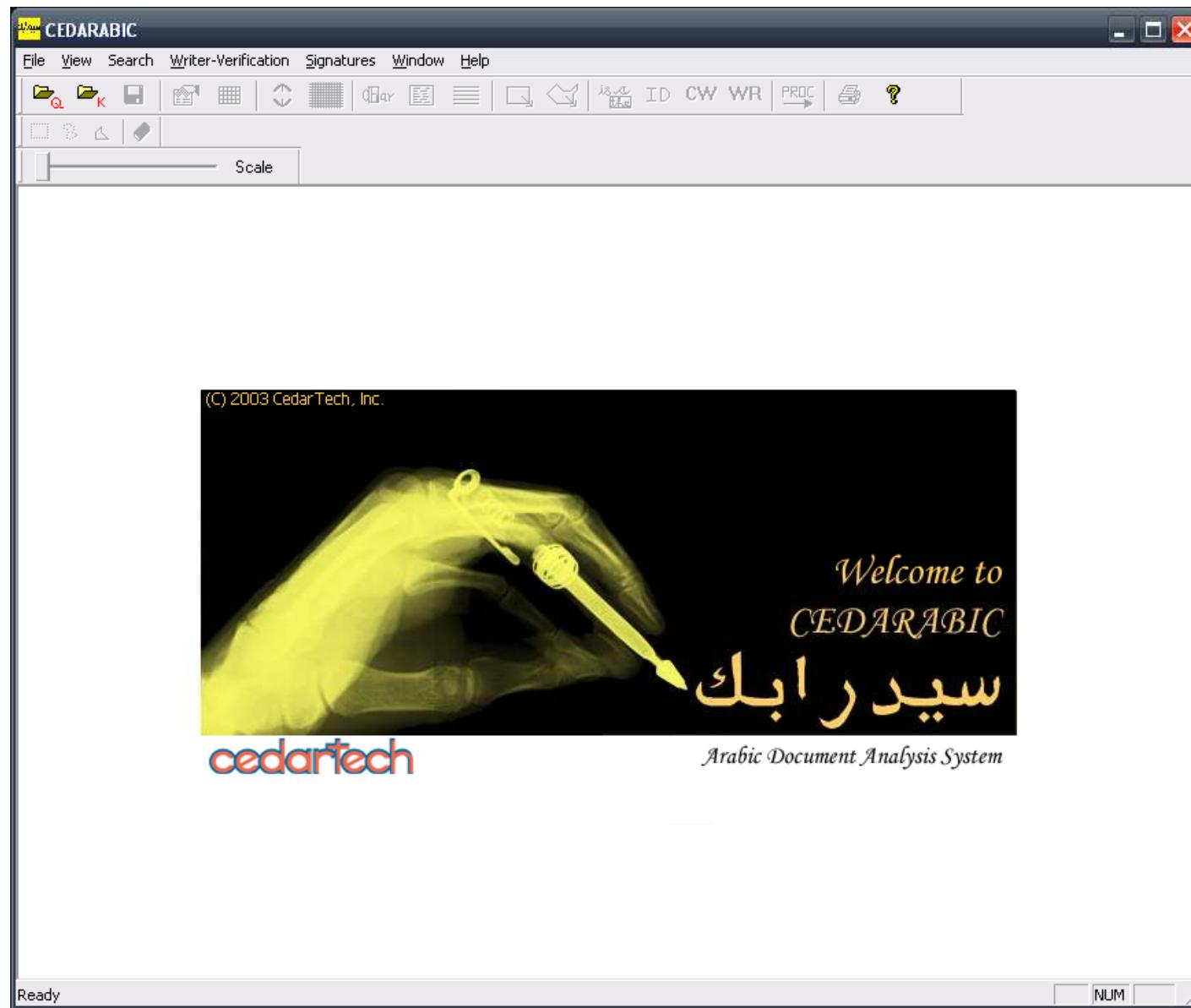
Results by thresholding LLR

- 144 comparisons for each test

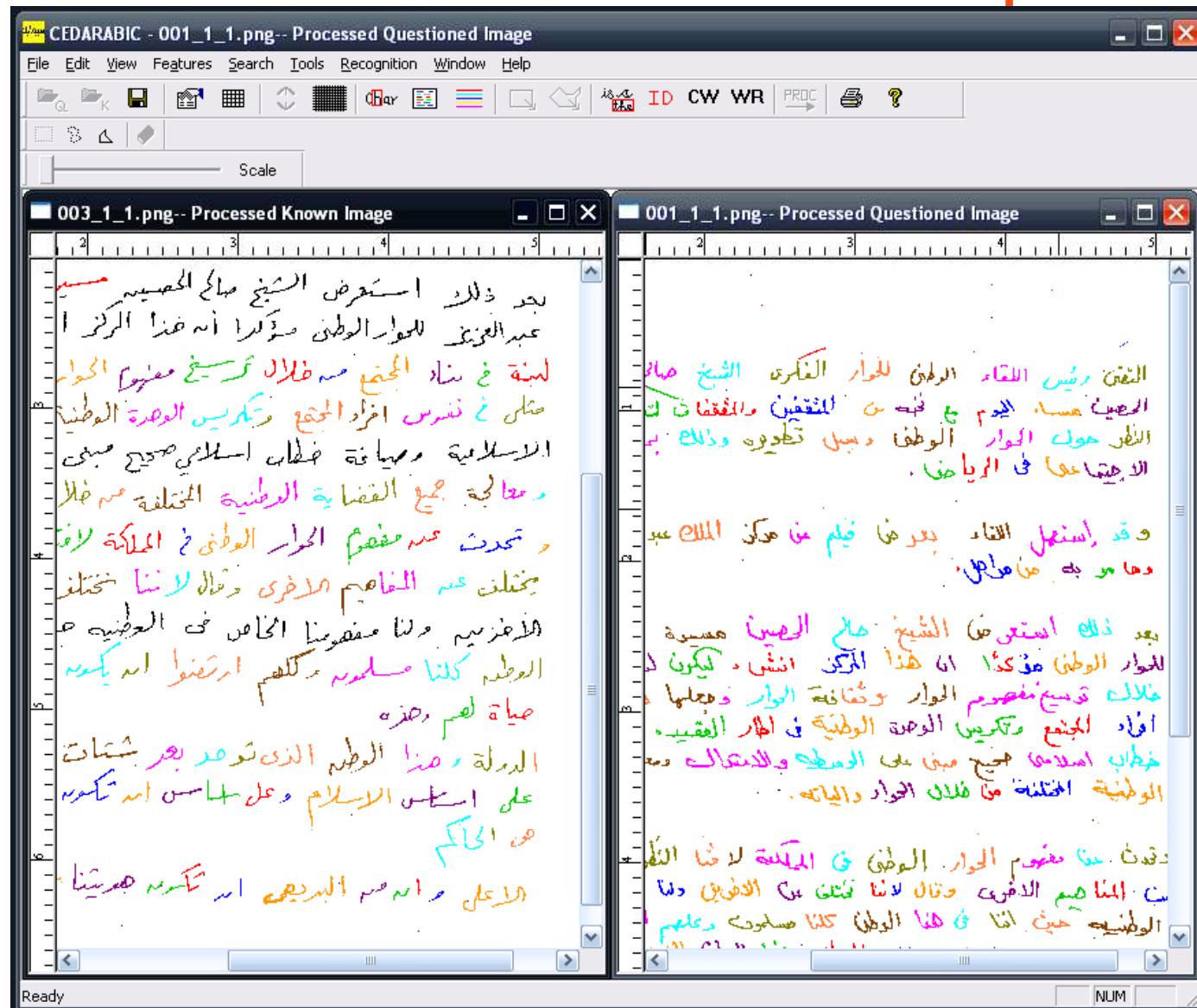
Discriminatory Method	Correct	Incorrect
Micro features, same writer	143 (99.3%)	1 (0.7%)
Micro features, different writer	105 (72.9%)	39 (27.1%)
Macro features, same writer	134 (93.1%)	10 (6.9%)
Macro features, different writer	99 (68.8%)	45 (31.2%)
Combination, same writer	143 (99.3%)	1 (0.7%)
Combination, different writer	104 (72.2%)	40 (27.8%)

- Overall accuracy: 86%

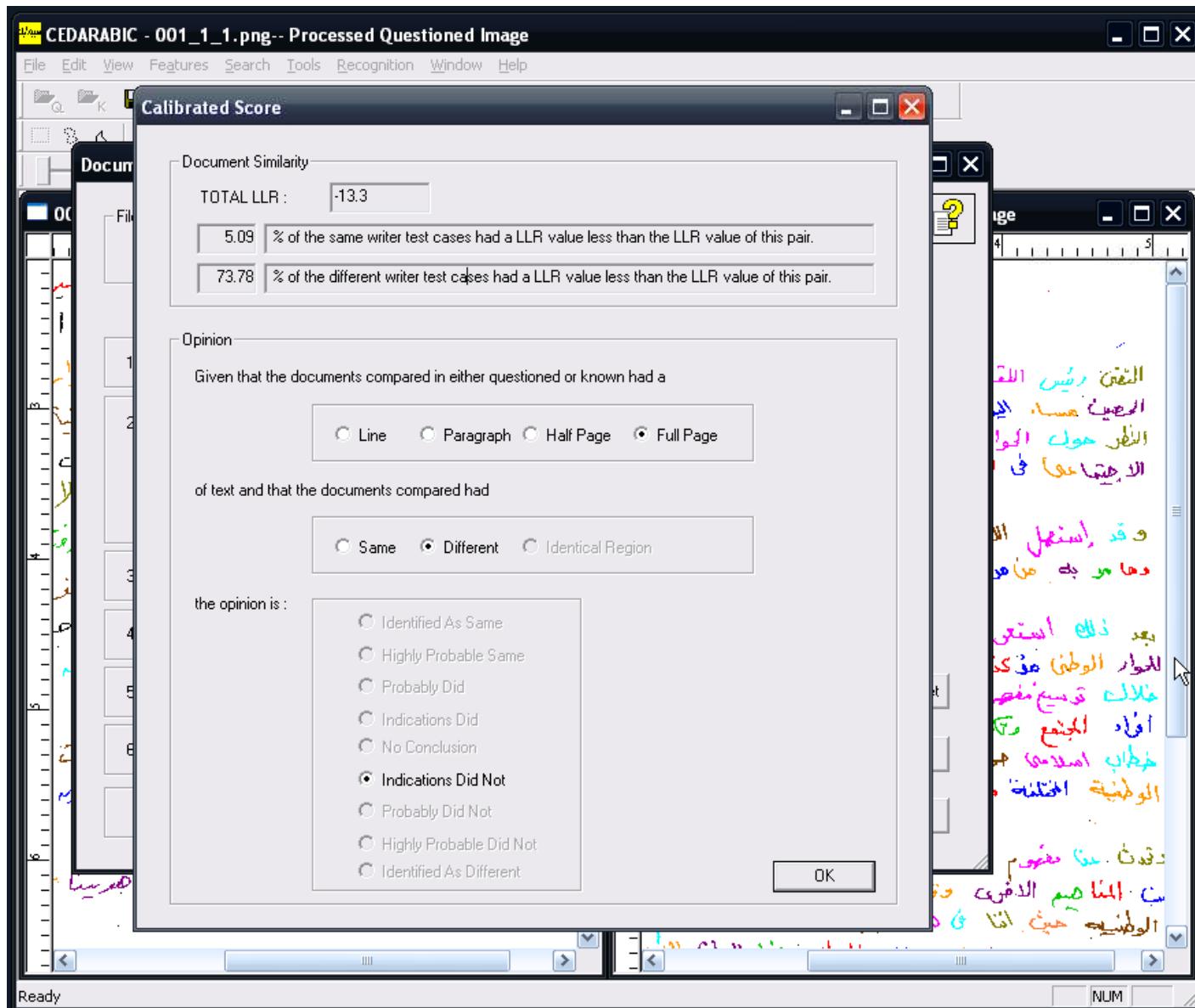
CEDARABIC User Interfaces



CEDARABIC – Document Comparison



CEDARABIC – Opinion, Strength of Evidence



Summary

- Transitioned writer verification system designed for English script to Arabic script
 - Macro features remain the same
 - Redesign of word segmentation algorithms
 - New character recognition module
- Writer Verification Performance
 - 86% average correct verification (10% lower than English)
 - Real-valued score gives relative confidence of answer
 - Performs better given longer documents

Future Plans

- Further Evaluations
 - Experiments with different data sets
 - Varying amounts of data (half-page and single line)
 - Writer identification performance
 - Several known authors, determine writer of test document
- Bayesian Formulation
 - Current system based on Maximum Likelihood formulation
 - Bayesian writer verification
 - To expect improved performance with small sample sizes
 - Already observed improvement in signature verification
 - Lesser effects of over-fitting
- Cedar-Fox and Cedarabic are now licensable
 - Looking for research partnerships