

Pattern Matching and Retrieval with Binary Features

Bin Zhang

Dept. of Human Genetics, David Geffen School of Medicine

Dept. of Biostatistics, School of Public Health

University of California, Los Angeles

April 28, 2004



Contents

- Background
- Dissimilarity Measures for Binary Vectors
- Fast k -Nearest Neighbor (k -nn) Classification
- Handwriting Identification with Binary Features
- Word Image Retrieval with Binary Features
- A Recognition-based System for Handwriting Matching and Retrieval
- Conclusions

Handwriting Matching and Retrieval

Three tasks of Handwriting Processing

- Handwriting Recognition
- Handwriting Identification
- Handwriting Retrieval

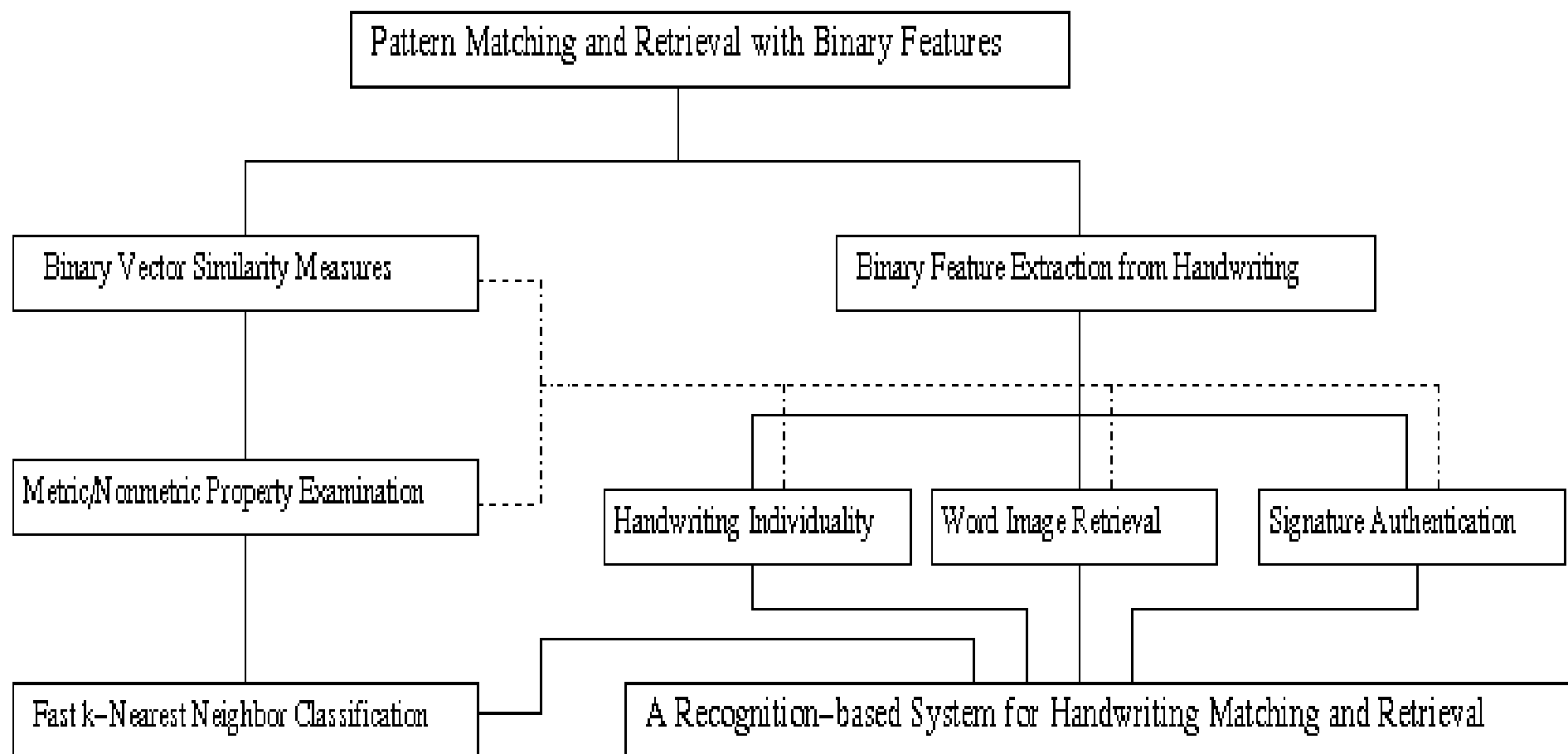
Handwriting recognition and identification are essentially two aspects of the handwriting pattern matching problem.

Motivation

Efficiency and effectiveness are two major concerns involved in handwriting matching and retrieval tasks:

- Effective Features
- Effective Recognizers
- Efficient Computation (for matching and searching)

Outline of the Research



Contents

- Background
- **Dissimilarity Measures for Binary Vectors**
- Fast k -Nearest Neighbor (k -nn) Classification
- Handwriting Identification with Binary Features
- Word Image Retrieval with Binary Features
- A Recognition-based System for Handwriting Matching and Retrieval
- Conclusions

Motivation

- Binary vector dissimilarity measures are widely used in PR and IR
- Properties of the measures have not been studied
- Metric/nonmetric examination is especially important as efficient pattern searching algorithms relying on metric properties usually fail for nonmetric dissimilarity measures
- Discovery of new properties
- Foundation for studying n-ary vector similarity measures

Metric Properties

Distance function d , defined in a set S , is a metric if it satisfies:

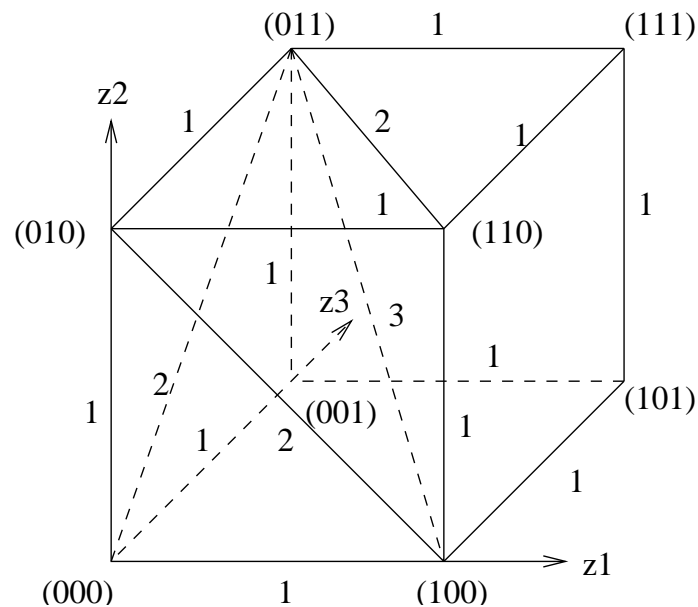
- Non-negativity: $d(x, y) \geq 0, \forall x, y \in S$;
- Commutativity: $d(x, y) = d(y, x), \forall x, y \in S$;
- Triangular Inequality:
 $d(x, z) \leq d(x, y) + d(y, z), \forall x, y, z \in S$;
- Reflexivity: $d(x, y) = 0 \iff x = y$.

Tri-Edge Inequality (TEI)

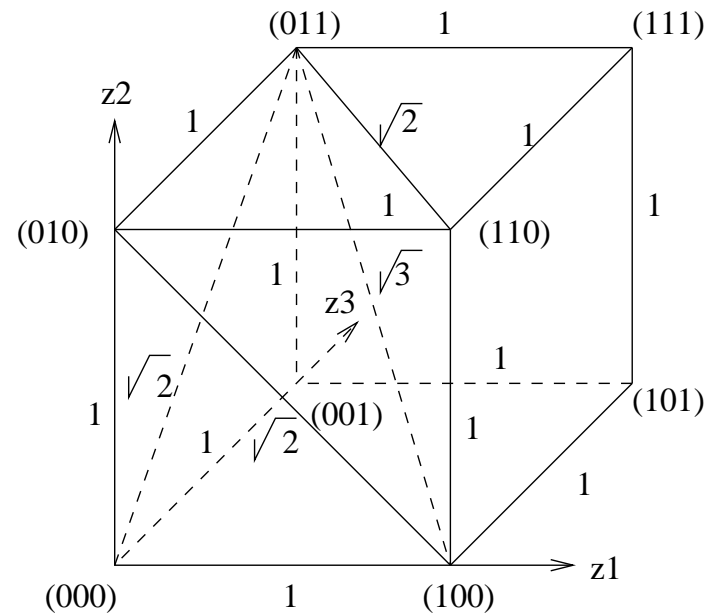
$\forall X, Y, U, V, P, Q \in \Theta$, the tri-edge inequality for a dissimilarity measure $D(\cdot, \cdot)$ is expressed as:

$$(1) \quad D(X, Y) + D(U, V) \geq D(P, Q)$$

Demonstration of TEI



Sokal-Michener Measure($\beta = 1$)



Euclidean Distance

TEI for Binary Vector Diss. Measur

Given two constants λ and M , Θ , as a subset of Ω (the set of all N -d binary vectors), is constrained by:

- $|X| \leq \lambda, \forall X \in \Theta$ and
- $X \cdot Y \leq M, \forall X, Y \in \Theta$

We discover and prove that several binary vector dissimilarity measures defined over Θ obey TEI.

Binary Vector Dissimilarity Measures

Measure	$D(\cdot, \cdot)$	metric	TEI
Jaccard-Needham	$\frac{S_{10}+S_{01}}{S_{11}+S_{10}+S_{01}}$	unknown	no
Dice	$\frac{S_{10}+S_{01}}{2S_{11}+S_{10}+S_{01}}$	no	no
Correlation	$\frac{1}{2} - \frac{S_{11}S_{00}-S_{10}S_{01}}{2\sigma}$	no	no
Yule	$\frac{S_{10}S_{01}}{S_{11}S_{00}+S_{10}S_{01}}$	no	no
Russell-Rao	$\frac{N-S_{11}}{N}$	no	$M \leq N/2$
Sokal-Michener	$\frac{2S_{10}+2S_{01}}{S_{11}+S_{00}+2S_{10}+2S_{01}}$	$\beta = 1$	$\beta \leq \frac{1}{1+\frac{2\lambda}{N-2M}}$
Rogers-Tanmoto	$\frac{2S_{10}+2S_{01}}{S_{11}+S_{00}+2S_{10}+2S_{01}}$	$\beta = 1$	no
Rogers-Tanmoto-a	$\frac{2(N-S_{11}-S_{00})}{2N-S_{11}-S_{00}}$	$\beta = 1$	yes
Kulzinsky	$\frac{S_{10}+S_{01}-S_{11}+N}{(S_{10}+S_{01}+N)}$	unknown	no

where, $\lambda = \max\{|X|\}$, $M = \max\{X \cdot Y\}$, and

$$\sigma = ((S_{10} + S_{11})(S_{01} + S_{00})(S_{11} + S_{01})(S_{00} + S_{10}))^{1/2}.$$

Rogers-Tanmoto-a measure

TEI holds under the following conditions:

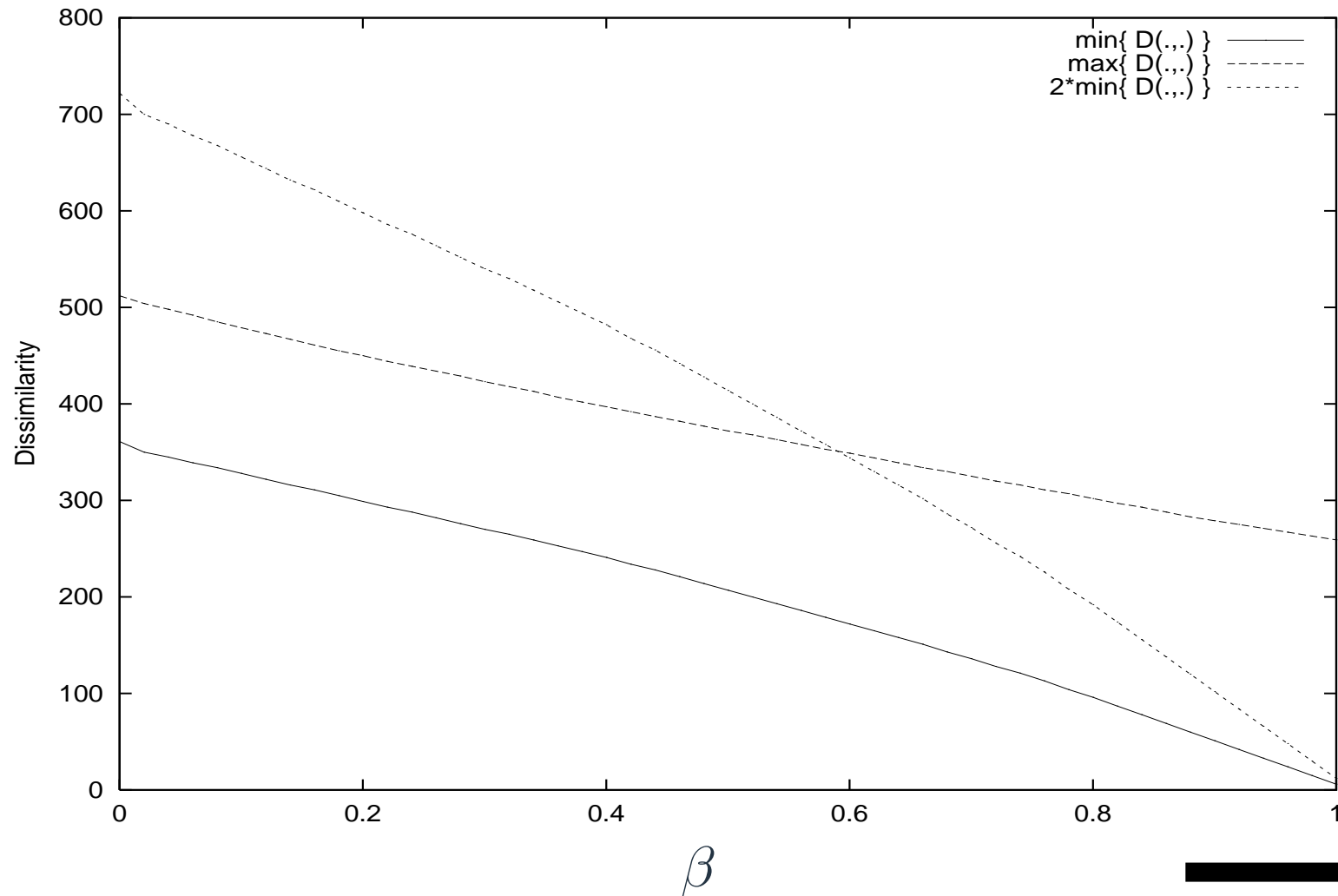
- (1) $\beta \leq 1 - \frac{N}{3(N-M)}$ if $\lambda = \frac{N}{2}$;
- (2) $\beta \leq \beta_{t2}$ if $\lambda < \frac{N}{2}$;
- (3) $\beta \in [0, 1] \cap [\beta_{t2}, \beta_{t1}]$ if $\lambda > \frac{N}{2}$. where β_{t1} and β_{t2} are given by

$$(2) \quad \beta_{t1} = \frac{a + \sqrt{\Delta}}{b}, \beta_{t2} = \frac{a - \sqrt{\Delta}}{b}$$

where, $a = 3N^2 + 2\lambda M - 4NM$, $b = 2(N - 2\lambda)(N - M)$,
and $\Delta = (2M(N + \lambda) - N^2)^2 + 8\lambda(2N - 3M)$.

Verification of TEI

$13565 * 14678 = 1.991070e + 08$ dissimilarities based on SM measure $D^{sm}(\cdot, \cdot)$.



Verification of TEI

- For $D^{sm}(\cdot, \cdot) | \beta = 0.5$, $\min\{D^{sm}(\cdot, \cdot)\} = 207$ and $\max\{D^{sm}(\cdot, \cdot)\} = 373$, so TEI holds.
- For $D^{sm}(\cdot, \cdot) | \beta \leq 0.58$,
 $2 * \min\{D^{sm}(\cdot, \cdot)\} \geq \max\{D^{sm}(\cdot, \cdot)\}$,
indicating TEI holds.
- For $D^{sm}(\cdot, \cdot)$, β_t is estimated as 0.43.

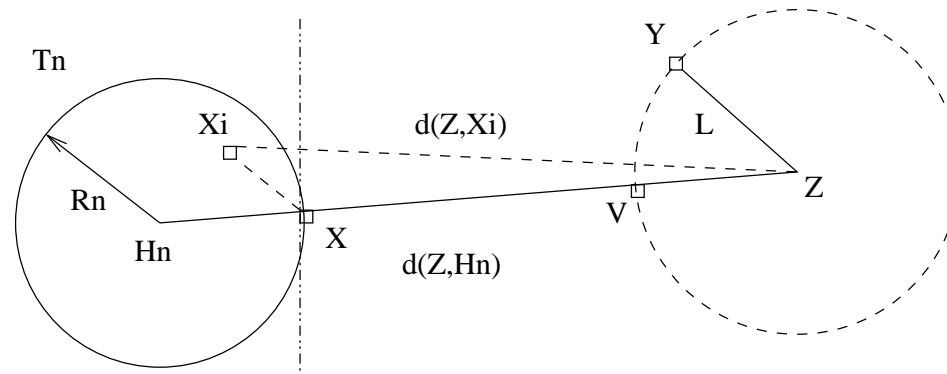
Contents

- Background
- Dissimilarity Measures for Binary Vectors
- **Fast k -Nearest Neighbor Classification**
- Handwriting Identification with Binary Features
- Word Image Retrieval with Binary Features
- A Recognition-based System for Handwriting Matching and Retrieval
- Conclusions

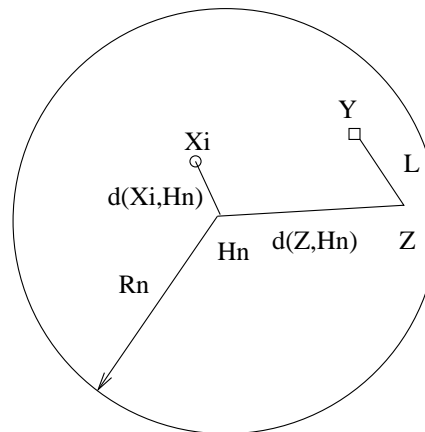
Motivation

- Exhaustive k -nearest neighbor search is slow
- Nonmetric measures are more robust in pattern matching than metrics
- The traditional reduction rules fail in some nonmetric spaces
- The condensation-based tree algorithm is not sufficiently efficient

Traditional Reduction Rules



Rule 1: no $X_i \in T_n$ can be the nearest neighbor to Z if $L + R_n < d(Z, H_n)$



Rule 2: $X_i \in T_n$ cannot be the nearest neighbor of Z if $L + d(X_i, H_n) < d(Z, H_n)$

Limitations of the Reduction Rules

- The rules utilize the fact that in $\Delta X Z X_i | \forall X_i \in T_n$, $Z X_i$ opposite an obtuse angle is the longest side, indicating $d(Z, X_i) > d(Z, X) > d(Z, Y) = L$.

The two rules fail in accelerating k -nn search using a nonmetric measure which

- obeys the Tri-Edge Inequality
since $L + R_n \geq d(Z, H_n)$
- OR violates the Triangle Inequality
if TI is violated, e.g. $d(Z, X_i) + d(X_i, X) < d(Z, X)$, without further information we cannot judge whether $d(Z, X_i)$ is bigger than $d(Z, Y)$ or not.

Condensation-based Tree Algorithm

- Developed by Randy L. Brown (IEEE TPAMI vol 15, no 3, 1995)
- Apply the condensation method to grow a template tree (top-down)
- Ensure that each training sample is correctly identified
- Require intense sorting operations in searching paths

The Cluster-based Tree Algorithm

- Use hierarchical clustering method (bottom-up)
- Introduce class-conditional clustering
- Establish two *decision levels* for early decision making

Local Properties of A Template

Given a template set \mathcal{T} . The local properties of $Z \in \mathcal{T}$ are:

- $\gamma(Z)$: the dissimilarity between Z and its nearest neighbor with a different class label.
- $\Psi(Z)$: a set of all neighbors which have the same class label as Z and are less than $\gamma(Z)$ far from Z .
- $\ell(Z)$: the size of the set $\Psi(Z)$.

Local Properties of A Cluster Center

Let \mathcal{T}^c be a set of cluster centers at a tree level.
The local properties of $Y \in \mathcal{T}^c$ are:

- $\Psi^c(Y)$, a set of all neighbors which are less than η far from Y
- $\ell(Y)$: the size of the set $\Psi^c(Y)$.

Growing A Cluster Tree

Step 0: Initialize the cluster tree to be a single level \mathcal{B} without node.

Step 1: Compute the localities of each template Z in \mathcal{T} , i.e., $\gamma(Z)$, $\Psi(Z)$, and $\ell(Z)$. Then rank all templates in \mathcal{T} in descendant order of $\ell(\cdot)$.

Step 2: Take the template with the biggest $\ell(\cdot)$, Z_1 , as a hyper node, and copy all templates of $\Psi(Z_1)$ as nodes at the bottom level of the tree, \mathcal{B} . Then remove all templates in $\Psi(Z_1)$ from \mathcal{T} , and set up a link between Z_1 and each pattern of $\Psi(Z_1)$ in \mathcal{B} .

Step 3: Repeat *Step 1* and *Step 2* till \mathcal{T} becomes empty. At this point, the cluster tree is configured with a hyper level, \mathcal{H} , and a bottom level, \mathcal{B} .

Step 4: Select a threshold η and cluster all templates in \mathcal{H} so that radius of each cluster is smaller or equal to η . All cluster centers form another level of the cluster tree, \mathcal{P} .

Step 5: Increase the threshold η and repeat *Step 4* for all nodes \mathcal{P} until only a single node is left in the resulting level.

Clustering above Hyper Level

- Relying only on nearness among a set of nodes.
- Iterate *Step 1* and *Step 2* based on *the localities of the cluster centers* at the current level.
- All cluster centers are actual data points.
- The same dissimilarity measure can be used in both tree generation and classification.

Parameter Choices

- η should be an increasing function of the number of iterations at *Step 4*
- Let $\eta(i)$ be the threshold for the i th iteration, a simple solution is

$$\eta(i) = \mu_i - \alpha \frac{\sigma_i}{1+i},$$

where α is a constant and μ_i and σ_i represent the mean and the standard deviation of the dissimilarities between the nodes at the current level respectively.

Classification Through Cluster Tree

Classification of an unseen template, \mathbf{X} :

Step 0: Compute dissimilarity between \mathbf{X} and each node at the top level of the cluster tree, and choose the ζ nearest nodes as a node set \mathcal{L}_x .

Step 1: Compute dissimilarity between \mathbf{X} and each subnode linked to the nodes in \mathcal{L}_x , and again choose the ζ nearest nodes, which are used to update the node set \mathcal{L}_x .

Step 2: repeat *Step 1* till the *hyper level* in the tree. As searching stops at the *hyper level*, \mathcal{L}_x consists of ζ hyper-nodes.

Step 3: Search \mathcal{L}_x for the hyper-nodes, each of which has a dissimilarity to \mathbf{X} less than its cluster radius. Let \mathcal{L}_h be a set consisting of these hyper-nodes. If all nodes in \mathcal{L}_h have the same class label, then this class is associated to \mathbf{X} and the classification process stops; otherwise, go to *Step 4*.

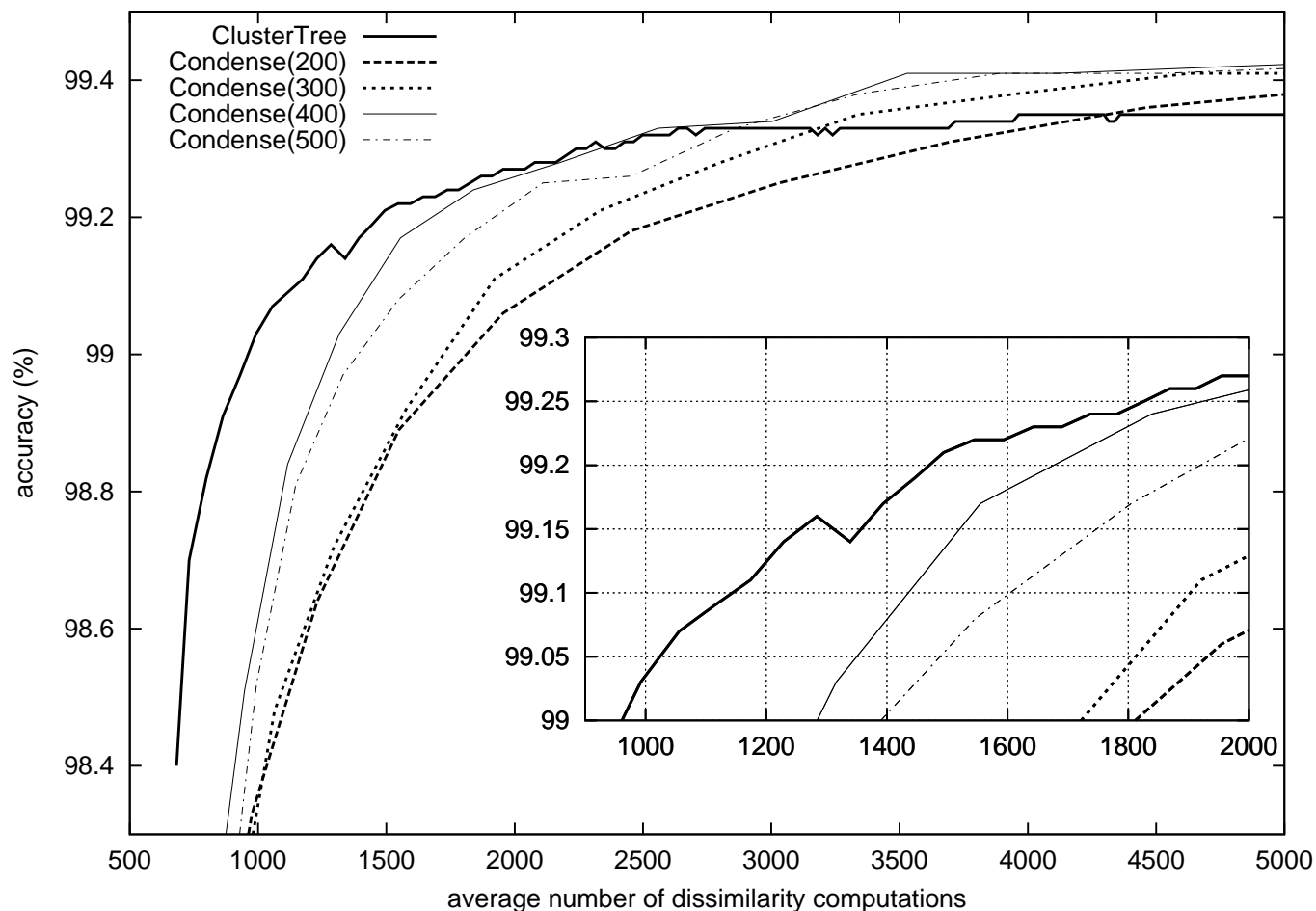
Step 4: Compute the dissimilarity between \mathbf{X} and every subnode linked to the nodes in \mathcal{L}_x , and choose the k nearest templates. Then, take a majority voting among the k nearest templates to decide the class label of \mathbf{X} .

Experimental Settings

- NIST database: a training set of 53,449 handwritten numeral images and a testing set of 5,3313 numeral images
- MNIST database: a training set of 60,000 handwritten numeral images and a testing set of 10,000 numeral images
- Feature extraction: 512-dimensional binary feature vector corresponding to gradient (192 bits), structural (192 bits), and concavity (128 bits) features.
- k -Nearest neighbor classification using a Sokal-Michener nonmetric
- UltraSPARC workstation: an 800MHz CPU, 512MB RAM, and SunOS 5.8

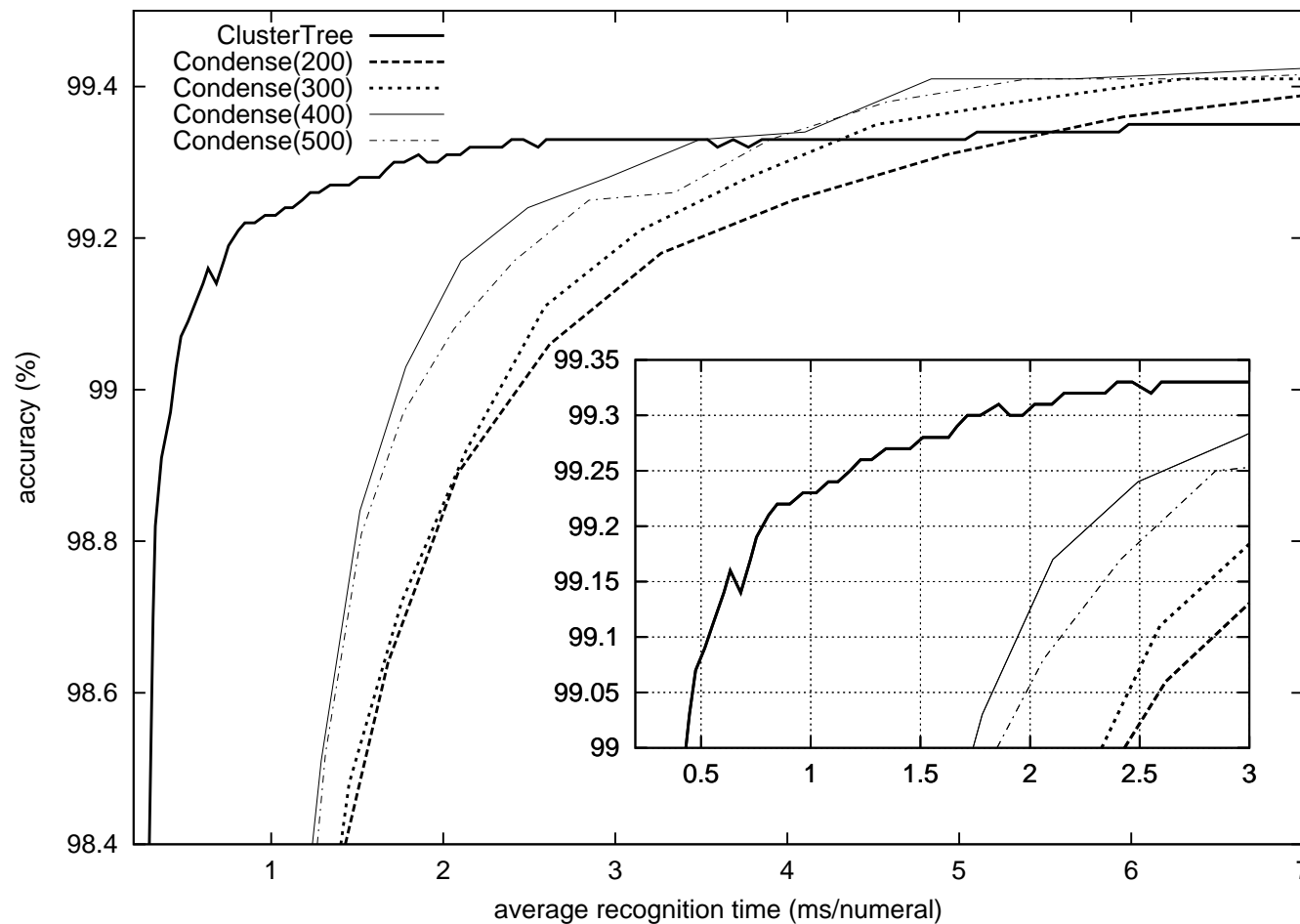
Dissimilarity Computations

Number of Dissimilarity Computations vs Accuracy



Recognition Time

Average Recognition Time vs Classification Accuracy



Detailed Comparison (1)

Comparison of a cluster tree and a condensation tree

	Recognition Performance on NIST Database							
	accur.=99.03%		accur.=99.17%		accur.=99.24%		accur.=99.28%	
	comput.	time	comput.	time	comput.	time	comput.	time
Cluster	992	0.45	1394	0.73	1737	1.08	2079	1.51
Conden.	1316	1.78	1555	2.10	1839	2.49	2178	2.96

Note: the exhaustive k -nn achieves an accuracy 99.25% with 20.68 ms per digit.

Detailed Comparison (2)

Comparison of a cluster tree and a condensation tree

	Recognition Performance on MNIST Database							
	accur.=97.25%		accur.=97.70%		accur.=97.86%		accur.=98.00%	
	comput.	time	comput.	time	comput.	time	comput.	time
Cluster	1115	0.51	1919	1.07	2329	1.48	3840	4.01
Conden.	1580	2.00	2675	3.42	3509	4.52	4052	5.23

Note: the exhaustive k -nn achieves an accuracy 98.24% with 23.26 ms per digit.

Discussions

Advantages of the cluster tree algorithm:

- Early decision making
- Minimal side-operations for choosing searching paths
- Tunable to fit any tradeoff between accuracy and speed
- Several times faster than a condensation tree
- Applicable to k -nn search using any dissimilarity measure

Contents

- Background
- Dissimilarity Measures for Binary Vectors
- Fast k -Nearest Neighbor (k -nn) Classification
- **Handwriting Identification with Binary Features**
- Word Image Retrieval with Binary Features
- A Recognition-based System for Handwriting Matching and Retrieval
- Conclusions

Motivation

- Measuring discriminability of each of sixty-two alphanumeric characters.
- Providing scientific guidance for choosing most-discriminative characters in examining forensic documents.
- Seeking new features from handwriting
- Developing high-performance classifiers for handwriting identification

Binary Features for Characters

Given a handwritten character image I , 512-bit gradient, structural, and concavity features are extracted as follows:

Step 0: Preprocessing: binarization and slant normalization.

Step 1: Equi-mass sampling: dividing I into 4 x 4 grids with equal mass for each of 4 rows and equal mass for each of 4 columns.

Step 2: Gradient features are computed by convolving 3*3 Sobel operators on I and capture the low-level gradient-direction frequency.

Step 3: Structural features are computed by applying 12 rules to each pixel in the gradient map to capture middle-level geometric characteristics including the presence of corners and lines at several directions.

Step 4: Concavity features are computed from the binary image to capture high-level topological and geometrical features including direction of bays, presence of holes, and large vertical and horizontal strokes.

Advantages of GSC Features

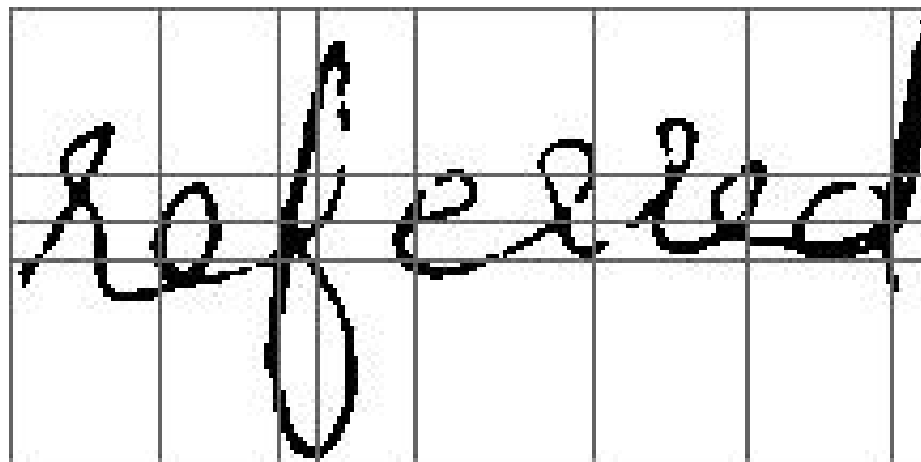
- **Quasi-multiresolution**: several distinct feature types are generated at different scales in image
- **Efficiency**: binary features lead to very efficient dissimilarity computation using bit-based operations
- **Effectiveness**: GSC features, coupled with k -nn classification, result in one of the best handwritten character recognizer.

Binary Feature Extraction for Word

Adapt the character-oriented GSC algorithm to handwritten word images:

- Equi-mass sampling: dividing I into $4 \times N$ grids with equal mass for each of 4 rows and equal mass for each of N columns.
- $N > 4$ depending on the width of a word image
- Resulting in $4 \times N \times 32$ binary features

Example of GSC Features



```
11100011100000000000000000111000111000111011110010011100000010111111100100000000000000011100011
1001111001111001111111111111111111111111111111111111111111111111111111111111111111111111111111111111111111
011111011110011110011110011110001111111111111111111111111111111111111111111111111111111111111111111111111111
000000111000101100111101111001001101000001111111111111111111111111111111111111111111111111111111111111111111
010100000110100100011000000000000000000000000000000000000000000000000000000000000000000000000000000000000000
00000000010000101100001011111100000010111100011001111000100111000000011011000110101111100
1111101111000110111111001000011111000001100110000000000110000011001100000011000000001100
000000001100000000001100110000000000000000000000000000000000000000000000000000000000000000000000000000000000
000000000000000000000000000000000000000000000000000000000000000000000000000000000000000000000000000000000000
000000100001101011111111000000010000000000000000000000000000000000000000000000000000000000000000000000000000
000000000011111111111111111111111111111111111111111111111111111111111111111111111111111111111111111111111111
000000000000
```

Classifier Design

Two models for Handwriting
Identification: Identification and Verification.
Identification Model

- Nearest neighbor classification

Verification Model: working in a distance space
converted from the feature space

- k -nearest neighbor (k -nn) classification
- Naive Bayes (NB) classification

Naive Bayes Classification

$$(3) \quad \omega_{NB} = \underset{\omega_i \in \{\omega_1, \dots, \omega_m\}}{\operatorname{argmax}} P(\omega_i) \prod_{j=1}^n p(f_j | \omega_i)$$

- Parametric Naive Bayes (P-NB) classification
- Non-Parametric Naive Bayes (NB) classification

Handwriting Sample Collection

3081 documents written by 1027 writers in US. Each writer copied three times **the CEDAR Letter** including:

- 156 words (concise)
- all alphanumeric characters (complete)
- each character in different position of a word
- certain character combinations of interest

A Handwriting Sample

From

Jim Elder

829 Loop Street, Apt 300
Allentown, New York 14707

To

Dr. Bob Grant

602 Queensberry Parkway
Omar, West Virginia 25638

We were referred to you by Xena Cohen at the University Medical Center. This is regarding my friend, Kate Zack.

It all started around six months ago while attending the "Rubeq" 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.

However, the extra hours affected her health; halfway through the show she passed out. We rushed her to the hospital, and several questions, x-rays and blood tests later, were told it was just exhaustion.

Kate's been in very bad health since. Could you take a look at the results and give us your opinion?

Thank you!

Jim

Nov 10, 1999

From

Jim Elder

829 Loop Street, Apt 300
Allentown, New York 14707

To

Dr. Bob Grant

602 Queensberry Parkway
Omar, West Virginia 25638

We were referred to you by Xena Cohen at the University Medical Center. This is regarding my friend, Kate Zack.

It all started around six months ago while attending the "Rubeq" 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.

However, the extra hours affected her health; halfway through the show she passed out. We rushed her to the hospital, and several questions, x-rays and blood tests later, were told it was just exhaustion.

Kate's been in very bad health since. Could you kindly take a look at the results and give us your opinion?

Thank you!

Jim

Nov 10, 1999

Character- and Word- Images

From each handwriting sample, we manually extracted

- 62 alphanumeric character images
- 4 characteristic word images
("been", "Cohen", "Medical" and "referred")

Features for Handwriting Sample

Features characterizing a handwriting sample include

- 11 real-valued global features from the entire document, consisting of darkness, contour and averaged line-level features.
- GSC binary features from 62 character images
- GSC binary features from 4 word images

Experimental Settings for Identification

- Query set: 875 randomly selected documents by 875 writers randomly chosen from 1027 writers.
- Exemplar set: the remaining 2206 documents.
- Each sample in the query set has two samples by the same writer

Experimental Settings for Verification

- Consider 3000 documents written by 1000 writers.
- Partition 1000 writers into two groups, each with 500 writers.
- Obtain two sets of document images (by two groups of writers), each with 1500 images.
- For each set, there are 1,500 document pairs, each by the same writer and $C_{500}^1 C_3^1 C_{499}^1 C_3^1 / 2 = 1,122,750$ document pairs, each by two different writers

Experimental Settings for Verification

- A test set consists of 1,500 document pairs, each by the same writer and 1,500 document pairs, each by two different writer.
- A training set consists of 1,500 document pairs, each by the same writer and 6,000 document pairs, each by two different writer.

Tasks of Experiments

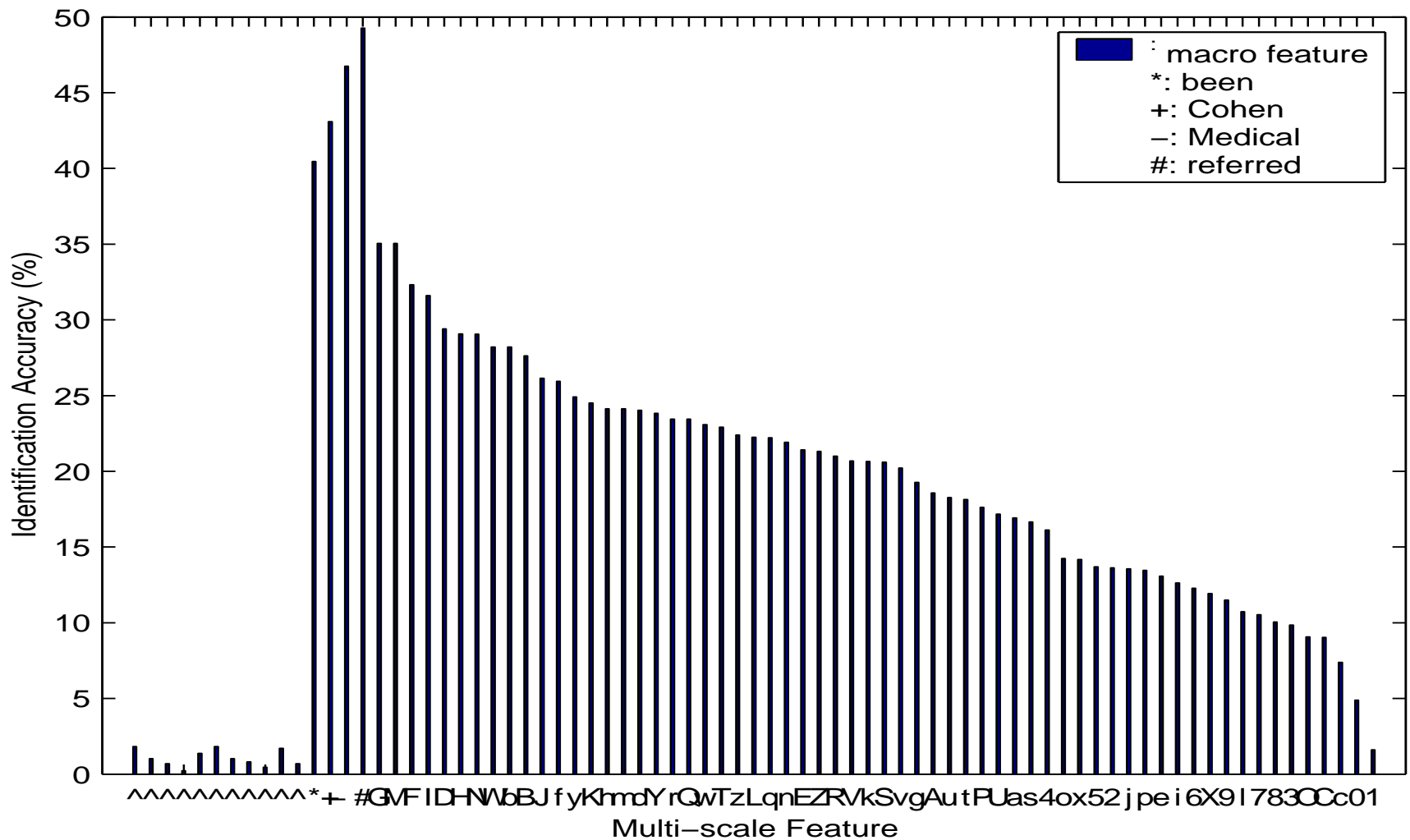
Identification using word-level GSC features under different divisions, and identification and verification performance under the following scenarios:

- Identification using word-level GSC features under different divisions
- Individual features: each macro feature, each word, and each character.
- Eleven macro features.
- Sixty-two characters.
- Four characteristic words.
- Ten numerals.
- Fifty-two alphabetic characters.
- All features

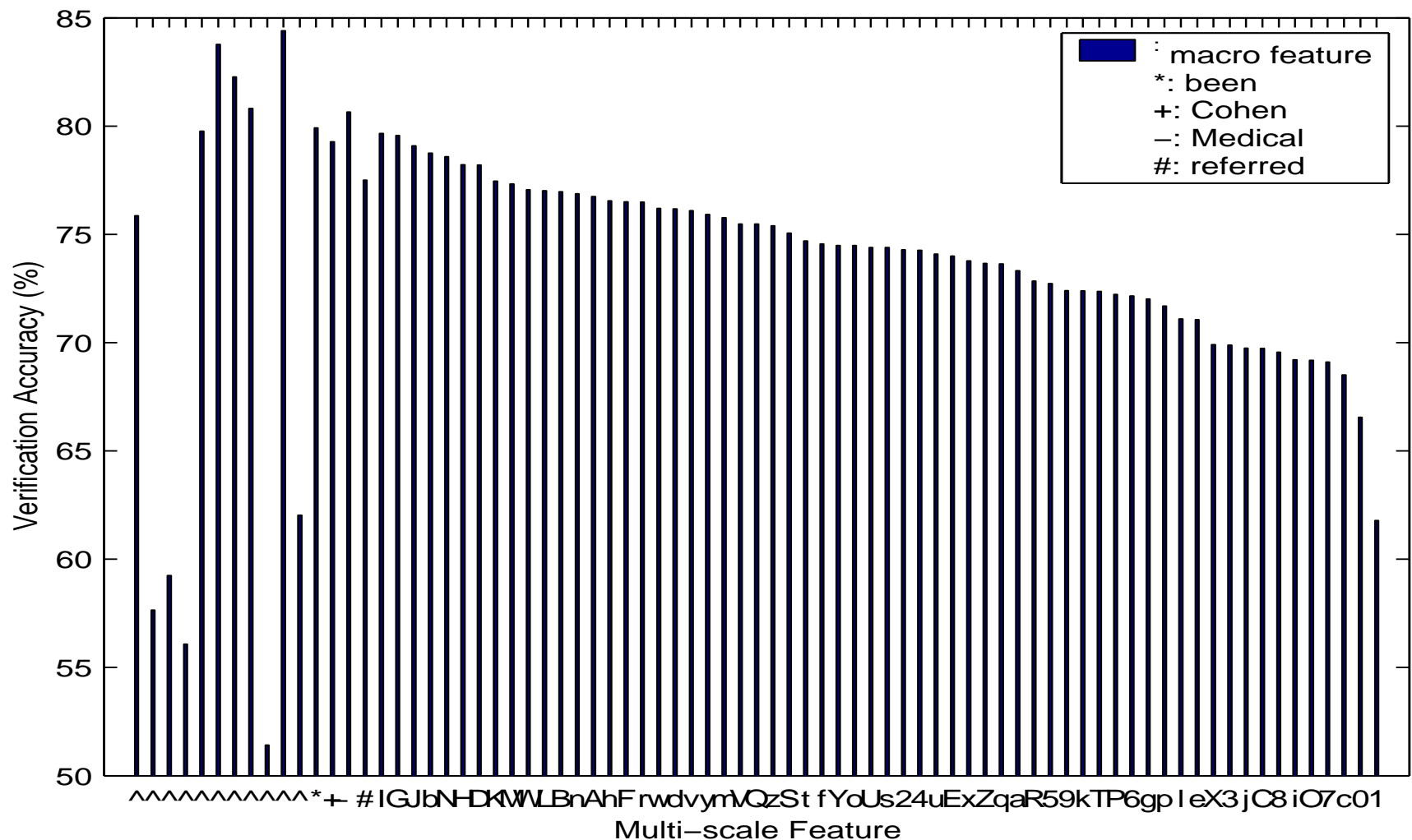
Identification with Word Images

Division	"been"	"Cohen"	"Medical"	"referred"
4 x 4	35.89	33.71	37.03	37.26
4 x 6	45.35	41.35	43.79	42.91
4 x 8	40.46	43.09	46.74	49.26
4 x 10	40.23	43.09	44.57	47.31
4 x 12	39.66	41.14	43.09	44.46

Identification by Individual Feature



Pattern Matching and Retrieval with Binary Features – p.52/87



Effect of Combined Features

Accuracy	Macro	Numerals	Alphabet	Alpha-Numerals
Identification	65.94%	74.86%	97.71%	97.83%
Verification	94.25%	89.46%	96.10%	96.40%

Accuracy	Lower-case Alphabet	Upper-case Alphabet	Words	Multi-scale
Identification	96.11%	95.77%	82.97%	98.06%
Verification	93.80%	94.52%	90.94%	97.07%

Comparison of k -nn and NB

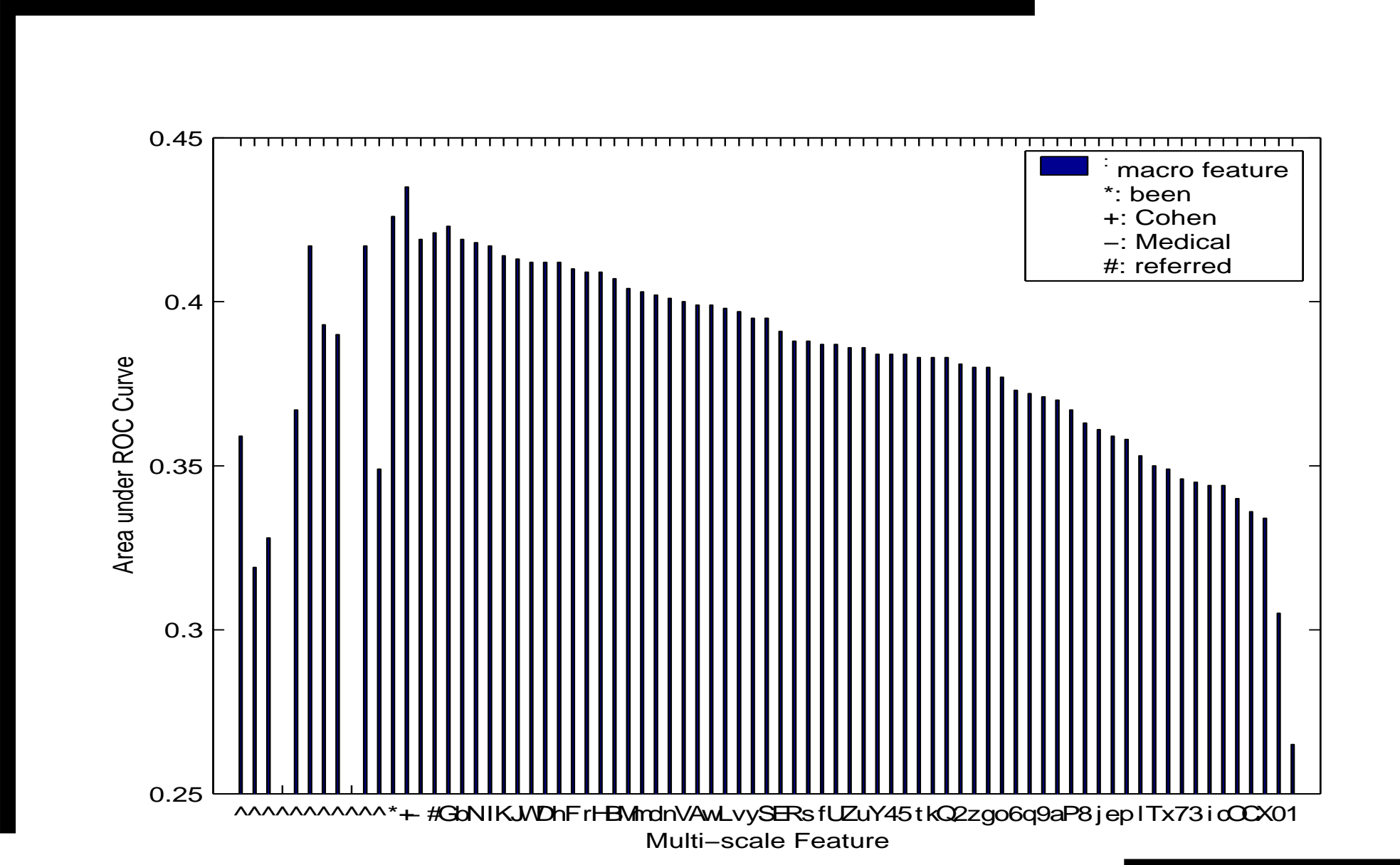
Accuracy	Macro	Numerals	Alphabet	Alpha-Numerals	Words	Multi-scale
k -nn	94.25%	89.46%	96.10%	96.40%	90.94%	97.07%
P -NB	91.67%	89.10%	93.63%	93.77%	91.00%	95.37%
nP -NB	92.83%	58.50%	79.50%	93.73%	83.57%	91.87%

Discriminability Measurement

Measure the discriminability of each character by using the associated receiver operating characteristic (ROC) curve.

- Compute the probability distributions $p_1(x)$ and $p_2(x)$ of the distance between two character samples conditioned on whether the samples belong to the same individual (intra-class) or to different individuals (inter-class)
- Construct the ROC curve for each character from the two distributions
- Determine the area under each ROC curve for each character

Ranking Order in Discriminability



Discussions

- Word-features usually bear high discriminability.
- Different characters have different discriminability.
- Numerals '0' and '1' have lowest discriminability.
- Ranking order of sixty-two alphanumeric characters:
'G'; 'b'; 'N'; 'l'; 'K'; 'J'; 'W'; 'D'; 'h'; 'F'; 'r'; 'H'; 'B'; 'M'; 'm'; 'd'; 'n'; 'V'; 'A'; 'w'; 'L'; 'v';
'y'; 'S'; 'E'; 'R'; 's'; 'f'; 'U'; 'Z'; 'u'; 'Y'; '4'; '5'; 't'; 'k'; 'Q'; '2'; 'z'; 'g'; 'o'; '6'; 'q'; '9'; 'a';
'P'; '8'; 'j'; 'e'; 'p'; 'l'; 'T'; 'x'; '7'; '3'; 'i'; 'c'; 'O'; 'C'; 'X'; '0'; '1'.
- Ranking order provides a scientific guidance for selecting most-informative characters in examining forensic documents.

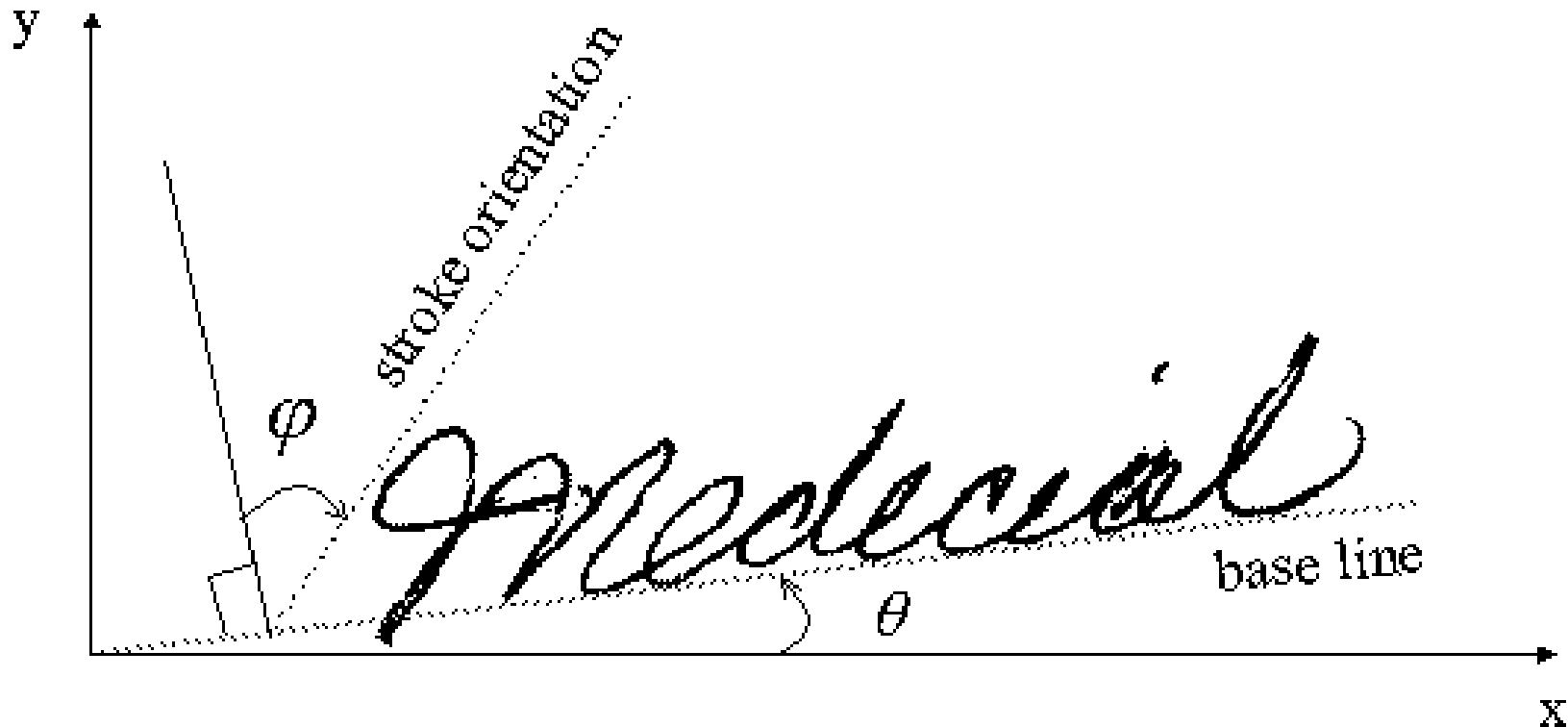
Contents

- Background
- Dissimilarity Measures for Binary Vectors
- Fast k -Nearest Neighbor (k -nn) Classification
- Handwriting Identification with Binary Features
- **Word Image Retrieval with Binary Features**
- A Recognition-based System for Handwriting Matching and Retrieval
- Conclusions

Motivation

- Difficulty of handwriting recognition in unconstrained domains
- Searching a repository of handwritten documents for those most similar to a given query
- Serving as a new human computer interface (HCI) and for forensic document examination
- Existing word image retrieval algorithms suffer from either low retrieval performance or high computation complexity.
- Apply word-level GSC features for retrieval purpose

Word Image Preprocessing



Skew and Slant Corrections

Skew and Slant Detections








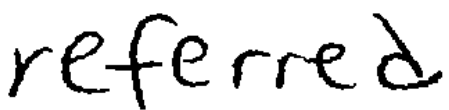
● Skew Detection

The skew angle is estimated as the orientation of the least eigenvector of the 2×2 matrix consisting of the 2-nd order moments.

● Slant Detection

The slant angle is computed as the average of all angles of vertical or near-vertical strokes, weighted by each stroke length

Examples of Preprocessing

before	after
	
	
	
	

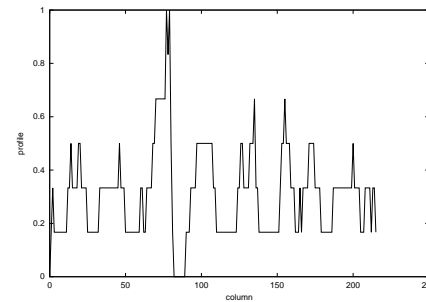
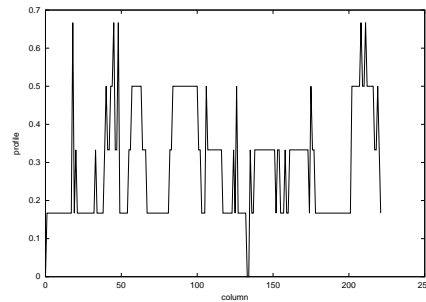
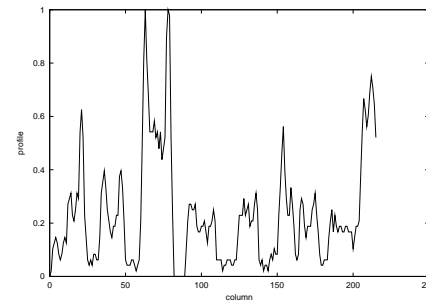
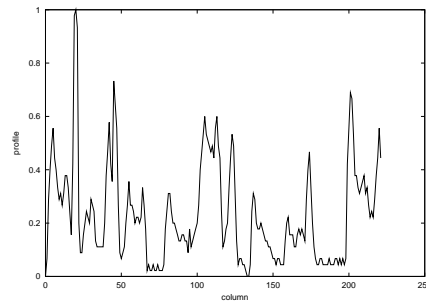
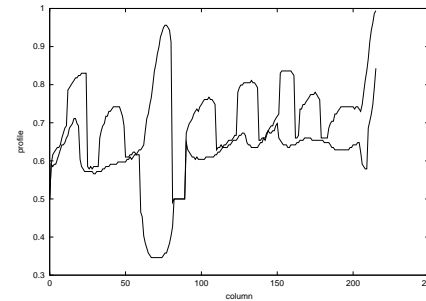
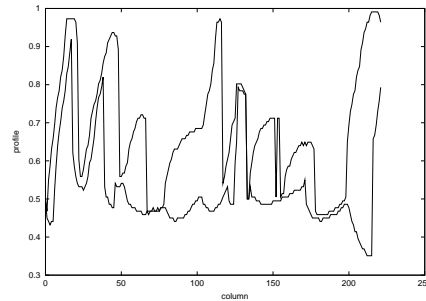
Existing Algorithms

- **XOR**, **SSD** (Sum of squared differences), **SLH** (by Scott and Longuet-Higgins), **SC** (Shape Context), **EDM** (Euclidean distance mapping), **DTW** (dynamic time warping matching) and **CORR** (recovered correspondences)
- **DTW** with profile-based shape features have been found the best among the 7 algorithms [R. Manmatha and Toni M. Rath, SDIUT 2003]

Profile-based Shape Features

Medical

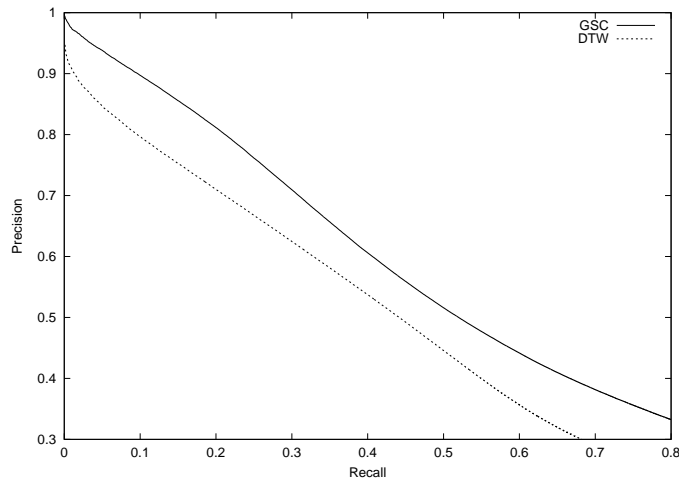
referred



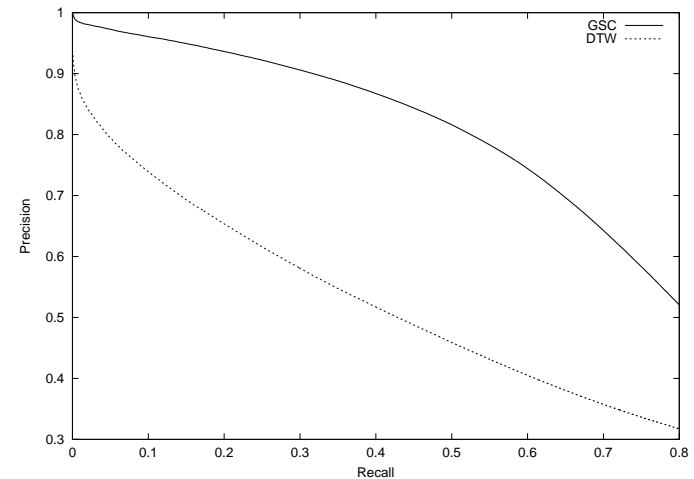
Experimental Settings

- four characteristic words, “been”, “Cohen”, “Medical”, and “referred”
- 9,312 word images (2328 images for each word) from 2328 handwriting samples by 776 individuals, each with three samples
- A query set includes $776 \times 4 = 3104$ word images
- A exemplar set includes $776 \times 4 \times 2 = 6208$ word images
- For each query image, there are $776 \times 2 = 1552$ relevant images in the exemplar set

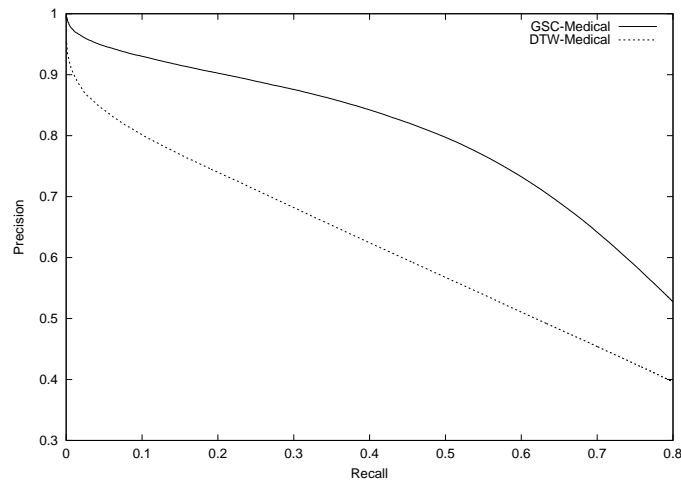
Comparison DTW and GSC (1)



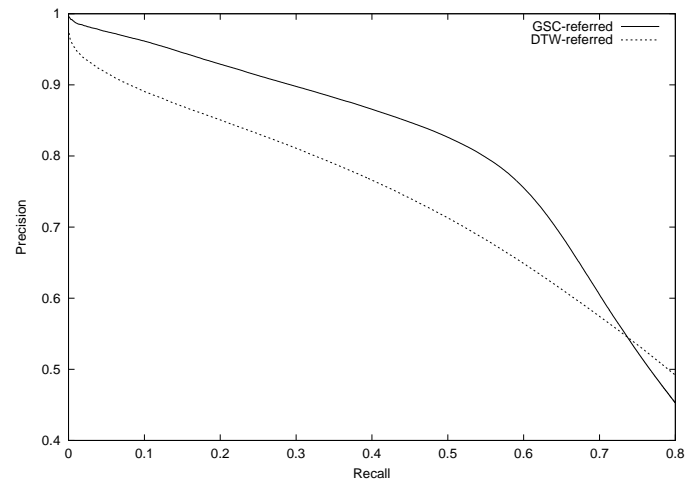
(a) word "been"



(b) word "Cohen"



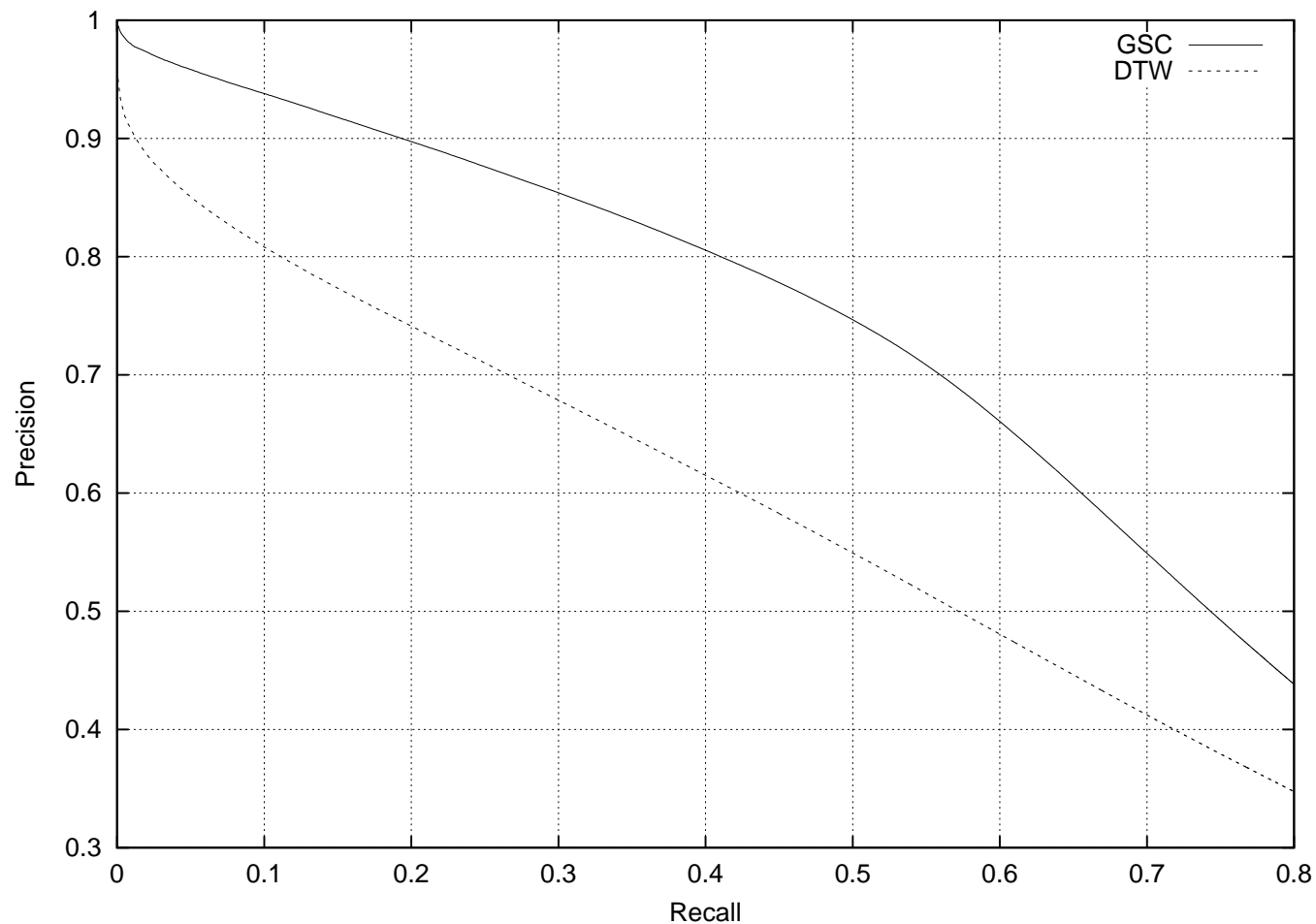
(c) word "Medical"



(d) word "referred"

Comparison DTW and GSC (2)

Average performance of four words



Comparison DTW and GSC (2)

N	Precision		Recall	
	GSC	DTW	GSC	DTW
1	0.9958	0.9507	0.00064	0.00061
5	0.9883	0.9305	0.00318	0.00300
10	0.9833	0.9174	0.00634	0.00591
20	0.9773	0.9023	0.01259	0.01163
50	0.9670	0.8756	0.03115	0.02821
100	0.9534	0.8462	0.06143	0.05452
500	0.8636	0.7211	0.27823	0.23231
1000	0.7554	0.6167	0.48675	0.39734
1552	0.6293	0.5294	0.62928	0.52941
2000	0.5456	0.4735	0.70311	0.61025
3000	0.4219	0.3841	0.81559	0.74255
4000	0.3461	0.3248	0.89199	0.83699

Discussions

- Running on an UltraSPARC workstation with an 800MHz CPU, 512MB RAM, and SunOS 5.8.
- DTW requires about seven days to finish 3104 queries.
- GSC takes only 11 minutes, 893 times faster than DTW.
- GSC also significantly outperforms DTW in retrieval precision and recall
- Quasi-multiresolution approach is better than profile-based shape features at low resolution
- The similarity between binary vectors is calculated using extremely fast bit-based operations while DTW suffers from the very high computation complexity.

Contents

- Background
- Dissimilarity Measures for Binary Vectors
- Fast k -Nearest Neighbor (k -nn) Classification
- Handwriting Identification with Binary Features
- Word Image Retrieval with Binary Features
- **A Recognition-based System for Handwriting Matching and Retrieval**
- Conclusions

Motivation

- Developing an automated system for handwriting identification and retrieval
- Increasing the efficiency and effectiveness of handwriting matching and retrieval
- Enabling massive handwritten data collection
- Searching large databases of handwritten documents
- Manual extraction of handwriting elements requires a large amount of effort and time, *e.g., it took a half year for 4 students to manually extract 186,000 alphanumeric images from about 3000 handwritten documents by 1000 writers.*

Handwriting Sample with Transcript

From

Jim Elder

829 Loop Street, Apt 300
Allentown, New York 14707

To

Dr. Bob Grant

602 Queensberry Parkway
Omar, West Virginia 25638

We were referred to you by Xena Cohen at the University Medical Center. This is regarding my friend, Kate Zack.

It all started around six months ago while attending the "Rubeq" 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.

However, the extra hours affected her health; halfway through the show she passed out. We rushed her to the hospital, and several questions, x-rays and blood tests later, were told it was just exhaustion.

Kate's been in very bad health since. Could you take a look at the results and give us your opinion?

Thank you!

Jim

Nov 10, 1999

From

Jim Elder

829 Loop Street, Apt 300
Allentown, New York 14707

To

Dr. Bob Grant

602 Queensberry Parkway
Omar, West Virginia 25638

We were referred to you by Xena Cohen at the University Medical Center. This is regarding my friend, Kate Zack.

It all started around six months ago while attending the "Rubeq" 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.

However, the extra hours affected her health; halfway through the show she passed out. We rushed her to the hospital, and several questions, x-rays and blood tests later, were told it was just exhaustion.

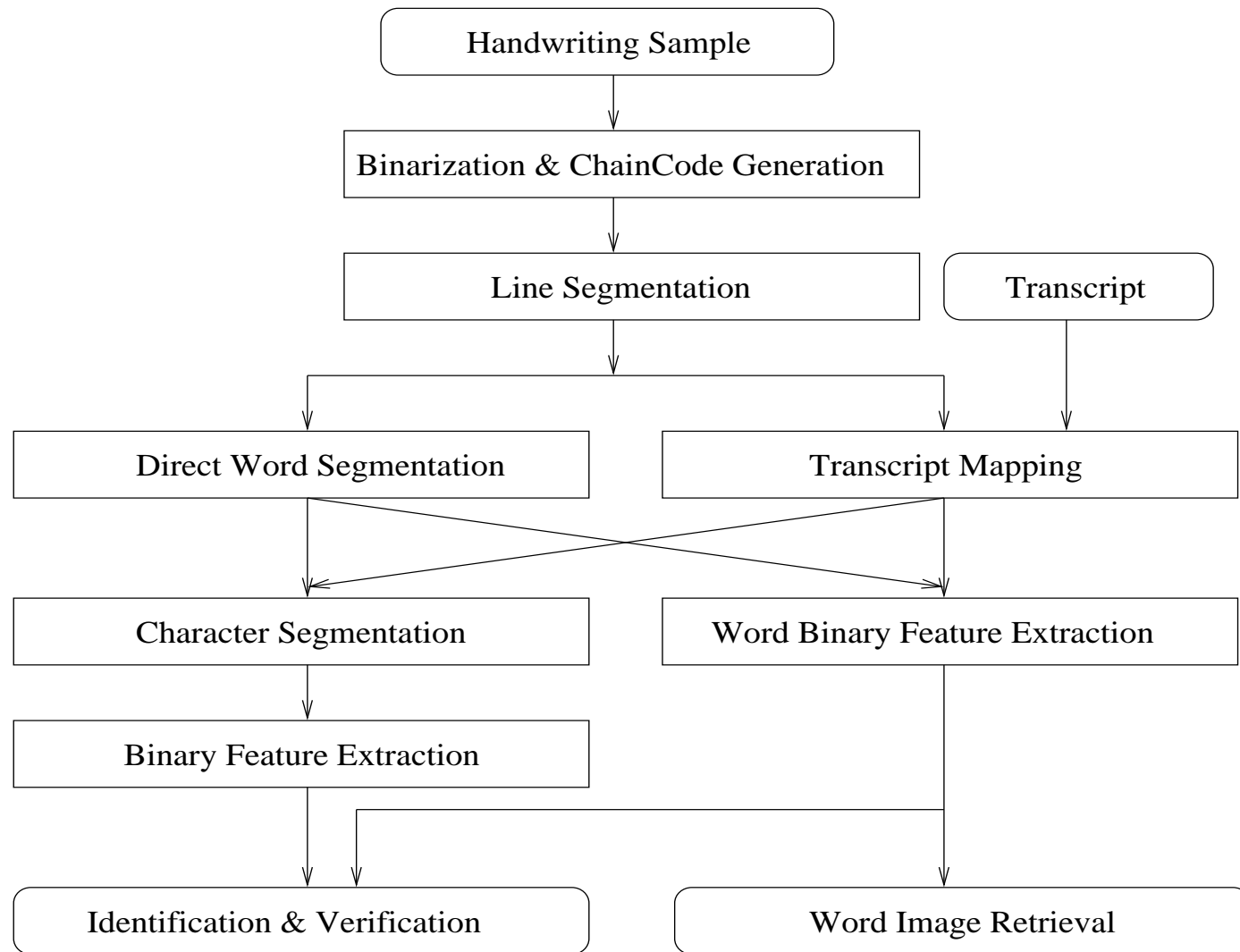
Kate's been in very bad health since. Could you kindly take a look at the results and give us your opinion?

Thank you!

Jim

Nov 10, 1999

System Architecture



Multiple Word-break Hypotheses

medical Center, This is regarding my friend, Kate Zack,

A line image

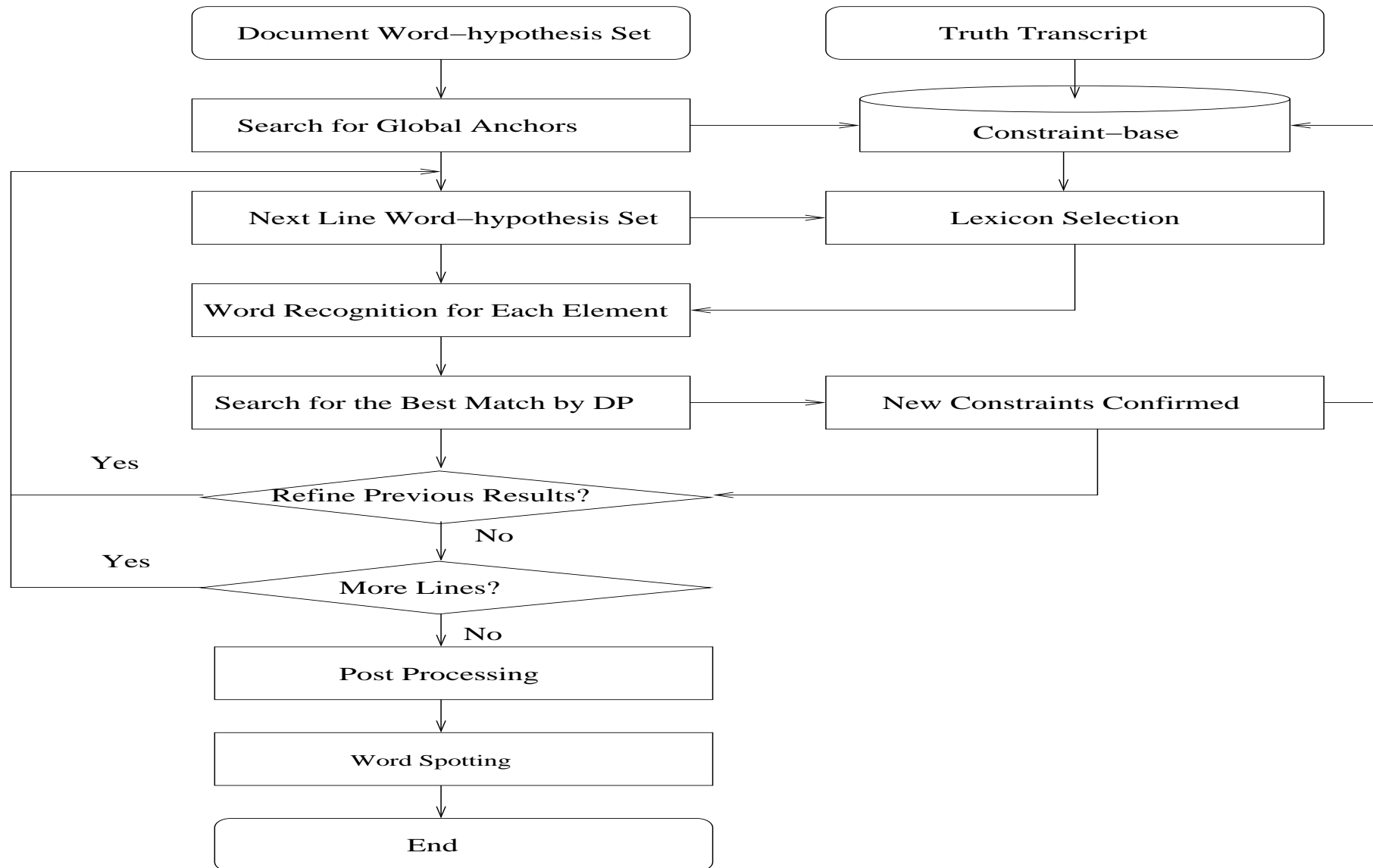
medical Center, This is regarding my friend, Kate Zack,

A word-segmentation hypothesis (big gap threshold)

medical Center, This is regarding my friend, Kate Zack,

Another word-segmentation hypothesis (small gap threshold)

Transcript-based Mapping



Lexicon Reduction & Optimization

A Line Image

- Line lexicon generation

For each word-break hypothesis

- Coarse word mapping
- Word-model Word Recognition
- Word alignment by a dynamic programming algorithm to Longest Common Subsequence (LCS)

Finalizing word mapping in the line

- Each word w in the transcript will be associated with several (or none for some cases) word images (hypotheses), and the word image h with the highest confidence is chosen as a mapping to w .

A Mapping Sample

From
Jim Elder
829 Loop Street Apt 300
Allentown New York 14707

Dr. Bob Grant
602 Queensberry Parkway
Omni West Virginia

We were referred to you by Xena Cohen at the University Medical Center. This is regarding my friend, Kate Zack.

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

However, the extra hours her health halfway through the show she passed out. We rushed her to the hospital and several x-rays and blood tests.

Kate been in very bad health since. Could you take a look at the results and give us your?

Thank you

From Nov 10 1999
Jim Elder
829 Loop Street Apt 300
Allentown New York 14707

??
Dr Bob Grant
602 Queensberry Parkway
Omni West Virginia

We were referred to you by Xena Cohen at the University Medical Center. This is regarding my friend Kate Zack.

?? all started around six months ago while attending the Rubeq Jazz Concert. Organizing such an event is no picnic and as President of the Alumni Association a cosponsor of the event, Kate was overworked. But she enjoyed her job and did what was required of her with great zeal and enthusiasm.

However the extra hours ??????? her health halfway through the show she passed out. We rushed her to the hospital and several ??????? x-rays and blood ????? were it was just ???????????

Kate been in very bad health since. Could you kindly take a look at the results and give us your ??????

Thank you
Jim

Experimental Results

- Experimented on 2997 documents written by 999 individuals
- Requiring 12 hours on P4 with 2.4 GHz CPU, 1 GB RAM, Windows 2000 OS
- 373,825 word images, 80% of the total 467,532, were extracted and recognized
- 1,331,253 character images were extracted from 373,825 word images

Evaluation of Word Images

With a verification-based evaluation method

- 18% (67,288) of all extracted word images are accepted with about 100% reliability.
- 37% (138,315) of all extracted word images are accepted with 99.2% reliability.
- 52% (194,389) of all extracted word images are accepted with about 95.2% reliability.
- 70% (261,677) of all extracted word images are accepted with about 82% reliability.

Evaluation of Character Images

- Manual extraction of 1,331,253 character images will take about 7 years for a full-time student.
- 56,179 numeral images (2.7% of the total) verified with an accuracy 93%
- Overall verification accuracy is 85.13%
- 74.6% of the 1500 document pairs written by the same writers and 95.6% of the 1500 pairs by the different writers can be correctly verified.

Analysis

The reason for the drop of the verification rate when multiple instances of each character from more than one word are used:

- a writer inscribes a character statistically differently with respect to its different positions in words
- matching two character images of the same content but at different relative positions in two words is inappropriate
- A training set with characters at various relative positions in words will greatly boost the performance of writer verification.

Contents

- Background
- Dissimilarity Measures for Binary Vectors
- Fast k -Nearest Neighbor (k -nn) Classification
- Handwriting Identification with Binary Features
- Word Image Retrieval with Binary Features
- A Recognition-based System for Handwriting Matching and Retrieval
- **Conclusions**

Contributions

- Property Examination of Binary Vector Dissimilarity Measures
- Discovery and Proof of a new property, the Tri-Edge Inequality
- Fast k -nn Classification Algorithm
- Individuality of Handwritten Characters
- Effective and Efficient Word Image Retrieval Algorithm
- A Recognition-based System for Handwriting Matching and Retrieval

Publications

19 conference and journal papers:

- "Fast k-Nearest Neighbor Classification Using Cluster-based Trees", **IEEE TPAMI 26(4), 2004.**
- "Individuality of Handwritten Characters", **the Best Paper Award of ICDAR 2003**, Edinburgh, Scotland.

Papers submitted or in review:

- IJDAR
- IEEE TPAMI
- Pattern Recognition.
- Journal of Forensic Document Examination
- Journal of Visual Communication and Image Representation

Thank you!

