

Segmentation and labeling of documents using Conditional Random Fields

S. Shetty, H. Srinivasan, S.N.Srihari and
M.Beal

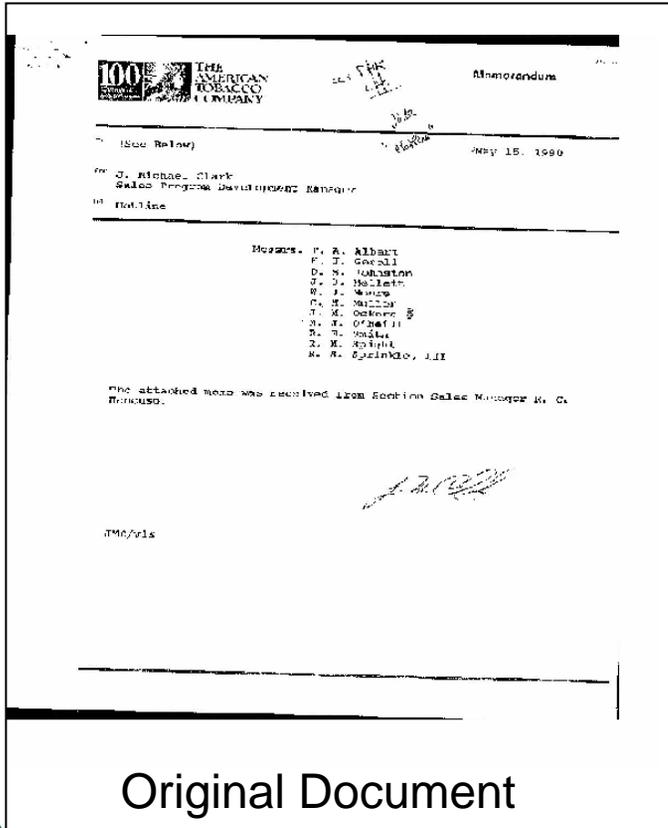
Center of Excellence for Document
Analysis and Recognition (CEDAR)

University at Buffalo
State University of New York

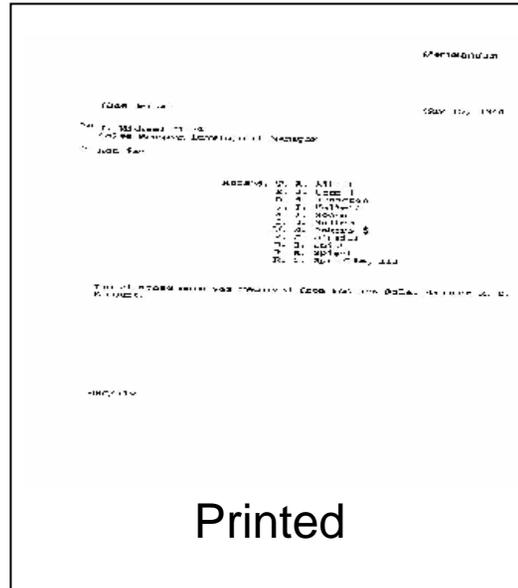
Overview

- Objective
- Segmentation and Neighbor Detection
- Classification using CRF
- Results
- Application – Ex: Signature Retrieval
- Conclusion

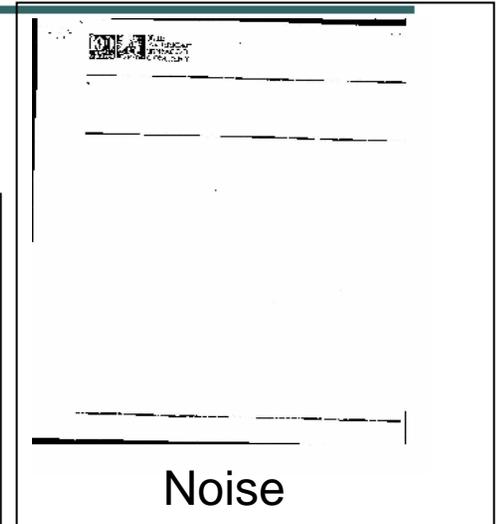
Objective



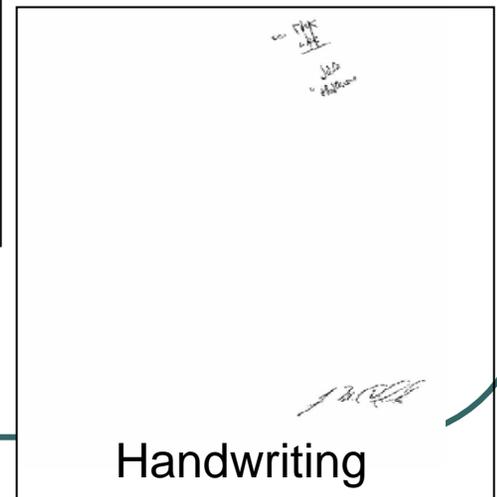
Original Document



Printed

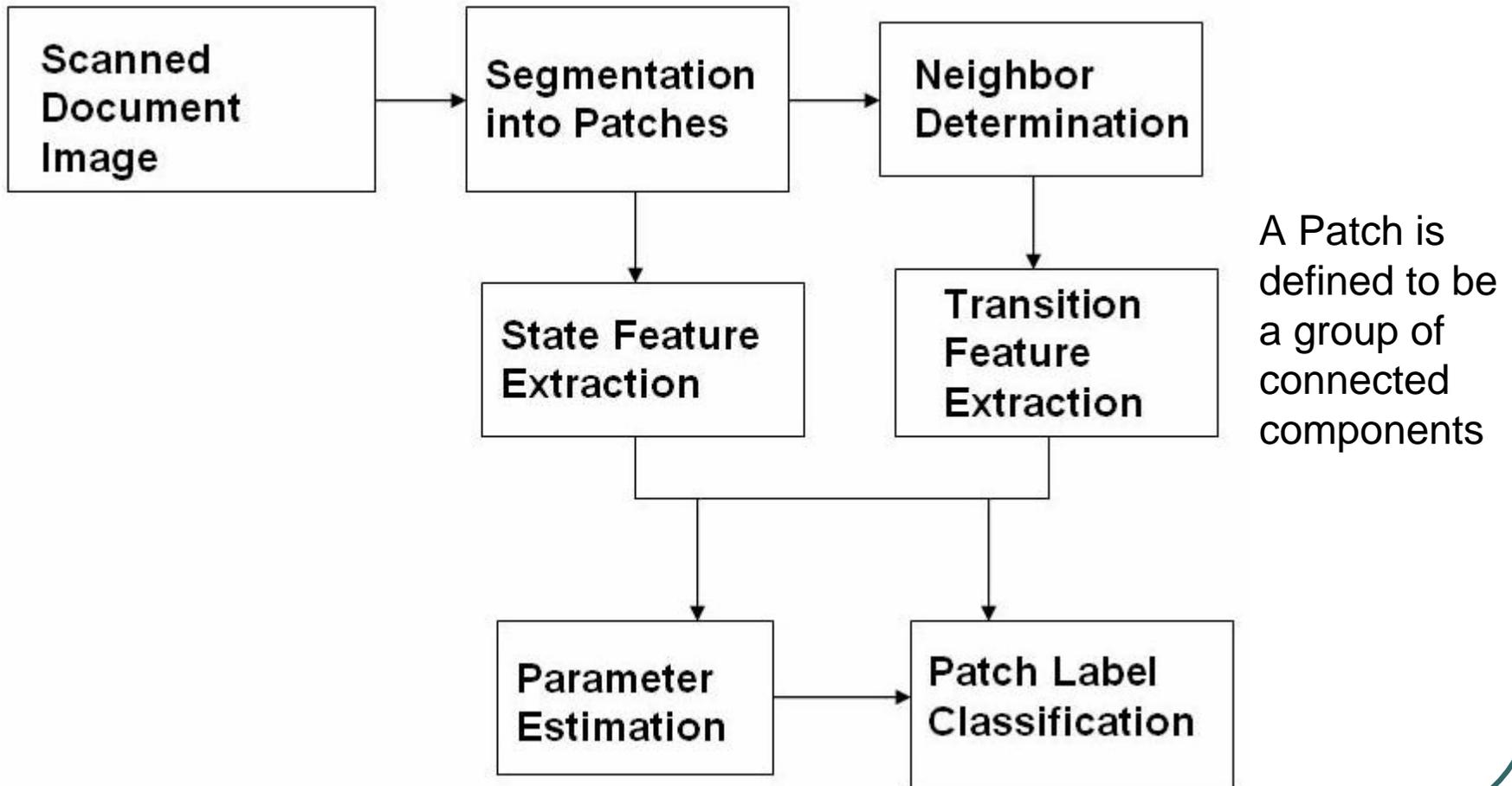


Noise



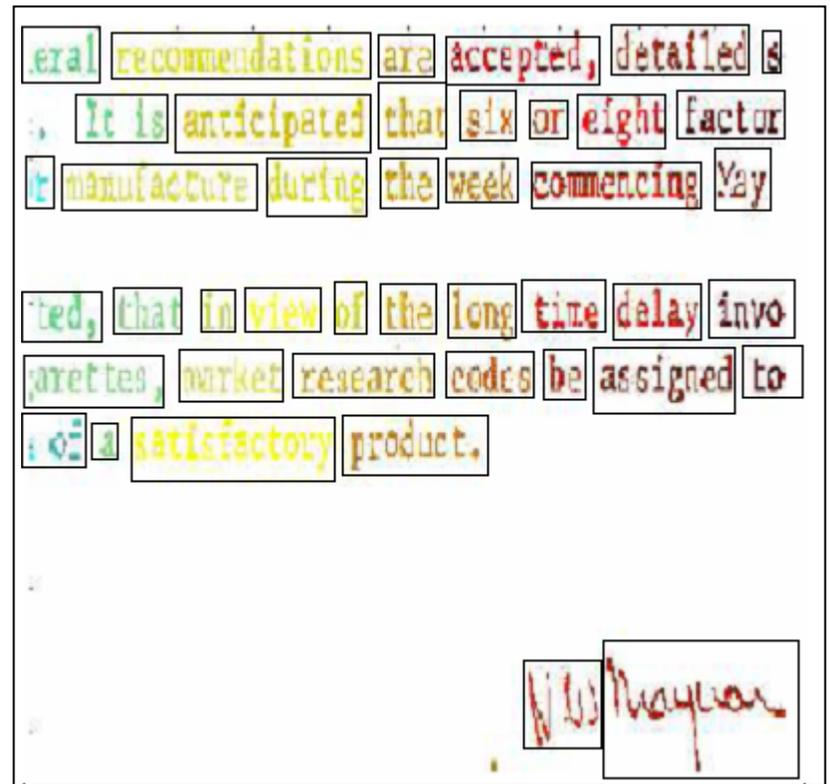
Handwriting

Block Diagram



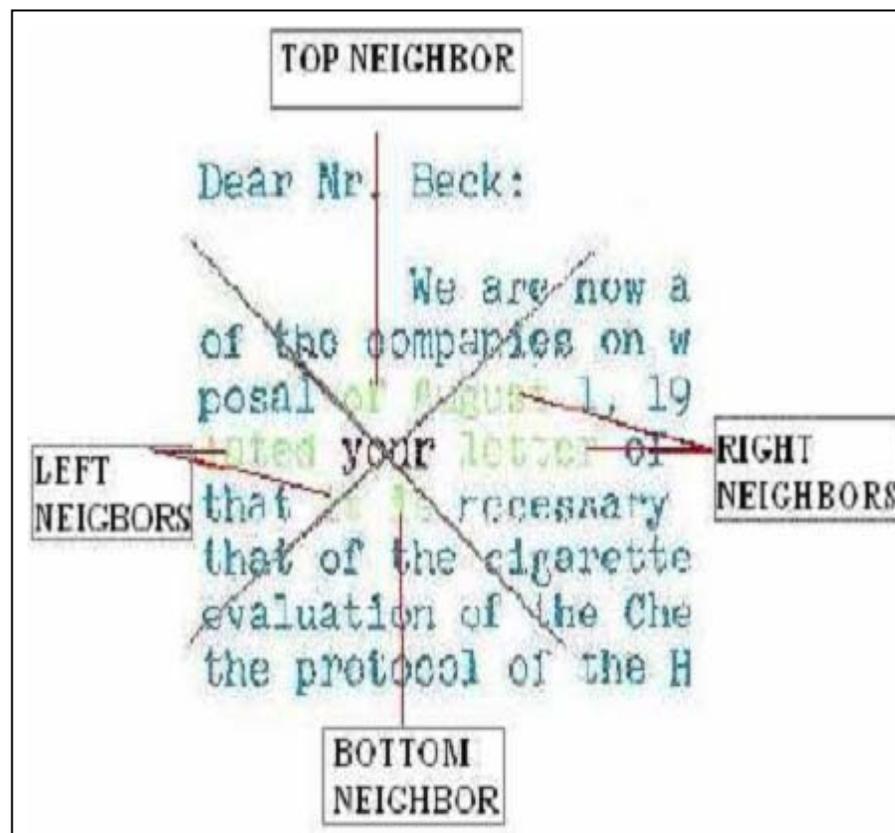
Segmentation

- The patches are generated using a region growing algorithm
- The size of each patch was optimized to represent approximately the size of a word
- The region within each rectangle represents a patch

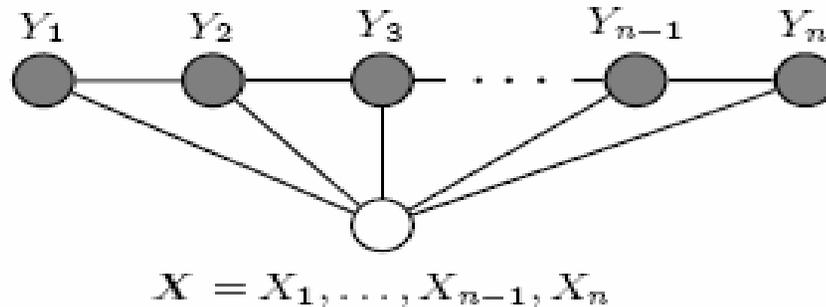


Neighbor Detection

- 6 neighbors are identified for each patch
- The closest(top/bottom) and the two closest(left/right) in terms of the convex-hull distance between the patches are identified as neighbors.



Advantage of using Conditional Random Fields



- CRF models the conditional probability $P(\text{label } Y \mid \text{observed } X)$ as compared to generative models like HMM which model the joint probability $P(Y, X)$
- It relaxes the assumption of conditional independence of the data given their labels as commonly used in Markov Random Fields

CRF Model

The probabilistic model of the Conditional Random Field used is

$$P(\mathbf{y}|\mathbf{x}, \theta) = \frac{e^{\psi(\mathbf{y}, \mathbf{x}; \theta)}}{\sum_{\mathbf{y}'} e^{\psi(\mathbf{y}', \mathbf{x}; \theta)}} \quad \begin{array}{l} \text{where } \mathbf{y}_i \in \{\text{Machine-print, Handwriting, Noise}\} \\ \theta : \text{CRF model parameters} \\ \mathbf{x} : \text{Observed document} \end{array}$$

$$\psi(\mathbf{y}, \mathbf{x}; \theta) = \sum_{j=1}^m \left(A(j, y_j, \mathbf{x}; \theta^s) + \sum_{(j,k) \in E} I(j, k, y_j, y_k, \mathbf{x}; \theta^t) \right)$$

where $A(j, y_j, \mathbf{x}; \theta^s)$: association potential

it associates the characteristics of that patch with its corresponding label

$I(j, k, y_j, y_k, \mathbf{x}; \theta^t)$: interaction potential

it captures the neighbor/contextual dependencies by associating pair wise interaction of the neighboring labels and the observed data

m : non-overlapping patches

CRF Parameter Estimation and Inference

- **Parameter estimation**

- Done by maximizing pseudo-likelihood parameters using conjugate gradient descent with line search optimization

$$\theta_{ML}^{\hat{}} \approx \arg \max_{\theta} \prod_{i=1}^M P(y_i | y_{N_i}, \mathbf{x}, \theta)$$

- **Inference**

- labels are assigned to each of the patches using Gibbs Sampling
- Initially random labels are assigned to each patch
- Over several iterations the label for each patch is generated from a distribution conditioned on the current values of its neighbors.

Features for HW/Print/Noise Classification

State Feature	Description
Height	Maximum height of the patch
Avg component width	The mean width of the connected components within a patch
Density	Density of foreground pixels within the patch
Aspect ratio	Width/Height of the patch
Gabor filter	8 features capturing the different stroke orientations
Variation of height	Variation in height within a patch
Width variation	Variation in width within a patch
Overlap	Sum of overlap in area between the connected components within a patch
Percentage of text above	Relative location of the patch with respect to the entire document
Number of components	Count of the connected components within a patch
Maximum component size	Maximum size of a component within a patch
Points in convex hull	Number of points in the convex hull of the patch
Maximum run length	The maximum horizontal run length within a patch
Avg run length	The average horizontal run length within a patch
Horizontal Transitions	A count of the number of times the pixel value transitions from white to black horizontally
Vertical Transitions	A count of the number of times the pixel value transitions from white to black vertically

Transition Feature	Description
Relative location	Assigned weights based on the relative location - top/bottom or right/left
Convex hull distance	The convex hull distance between the 2 patches
Ratio of aspect ratio	The ratio of the aspect ratio values of the 2 patches
Ratio of number of components	The ratio of the number of components present in the 2 patches

Output

Harley-Davidson Motor Co., Inc. 1300 W. Wisconsin St. Milwaukee, Wisconsin 53233-1600

May 2, 1986

PERSONAL EXHIBIT

Mr. Thomas Mau
Lorillard, Inc.
666 Fifth Avenue
New York, New York 10019

Dear Tom,

Attached is a copy of a proposed letter agreement with Brown & Williamson. This letter was worked out in a series of telephone conversations with B&W following their receipt of my letter of March 31, 1986.

In my opinion this is an acceptable compromise and I am prepared to sign it if it is acceptable to you. If it is, I would also modify the reference to the March 31 letter to our Agreement and substitute a signed copy of the May 1, 1986 letter.

Very truly yours,



Jeffrey L. Weinstein
Vice President - Parts and Accessories

/pat
Attachment

PRINTED TEXT

HANDWRITTEN
TEXT

NOISE

95507804

Data set

- On average each document consists of 140 patches.
- 26 documents consisting of 3700 patches were used for training
- 27 documents consisting of 3800 patches were used for testing

Results

Type of Label	Accuracy	Precision	
Machine Printed Text	98.36%	96.29%	
Handwritten Text	94.81%	90.18%	
Noise	89.80%	95.72%	
Overall	95.75%	-	

Comparison of error rates

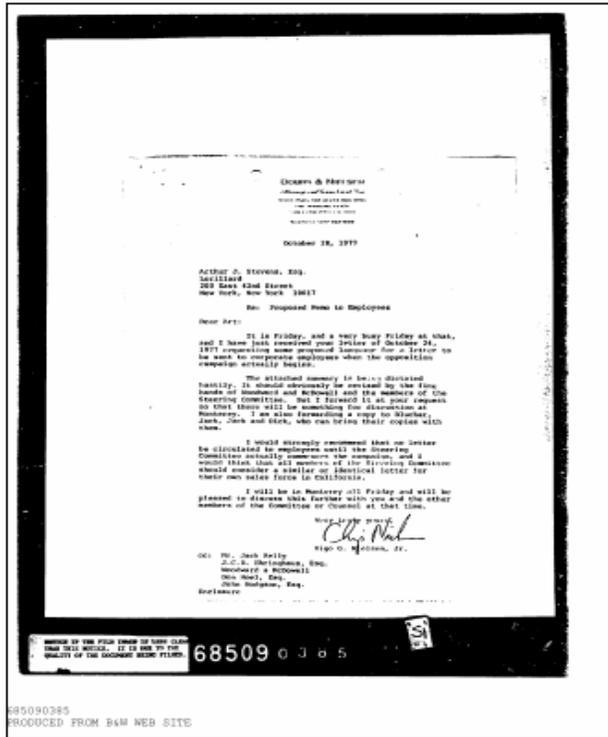
	CRF	Neural Network	Naive Bayes
Machine Printed Text	1.64%	2.35%	11.54%
Handwritten Text	5.19%	20.90%	25.04%
Noise	10.20%	15.00%	12.23%
Total	4.25%	7.04%	12.58%

Applications of labeled documents

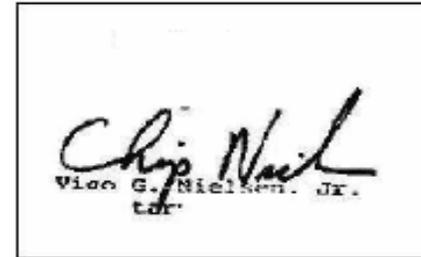
- Content based Image Retrieval
- Optical Handwriting Recognition
- Signature based biometric verification

Signature Retrieval

Step 1 - Signature Extraction

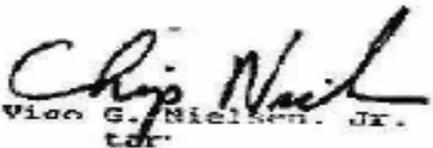
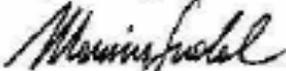
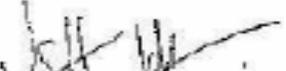
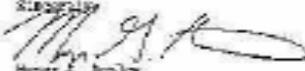
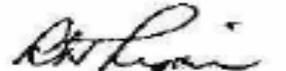


Original Document



Extracted Signature

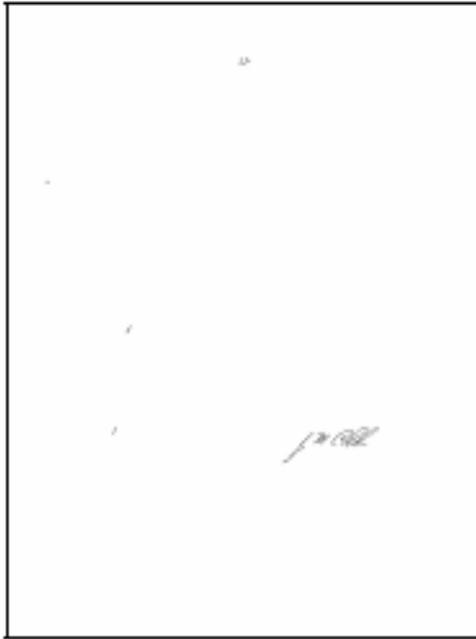
Step 2 – Signature Retrieval

	RETRIEVAL	Dissimilarity
<p>QUERY IMAGE</p> 	<p>Very truly yours,  Vigo G. Nielsen, Jr.</p>	0.26
	<p>Yours sincerely,  Martin Jodal, P.L.T.</p>	0.36
		0.39
	<p>Sincerely,  Chip Nielsen</p>	0.43
		0.49

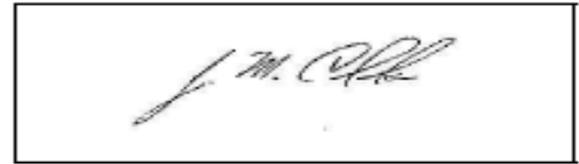
Indexing Steps

- **Signature Block Extraction**
 - Patch Classification using Conditional Random Fields
 - Post-processing
- **Noise Removal**
 - Support Vector Machines are used to remove noise or printed text overlapping the signature block
- **Signature Feature extraction**
 - Gradient, Structural and Concavity features

Post-processing

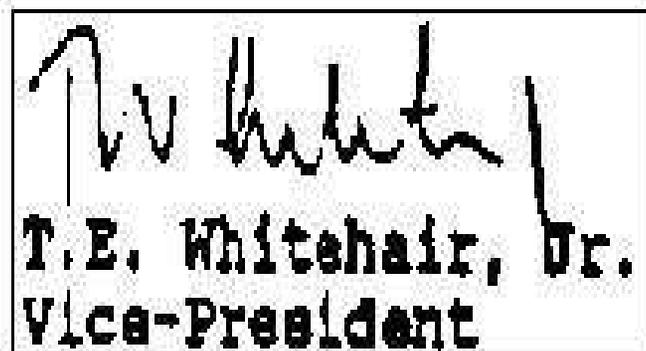


Handwritten components
pulled out from CRF labels

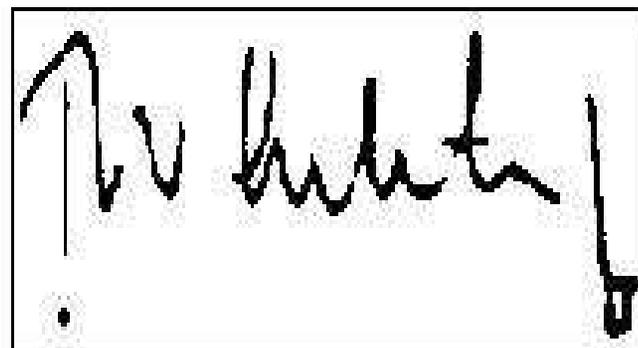


Extracted Signature

Noise Removal

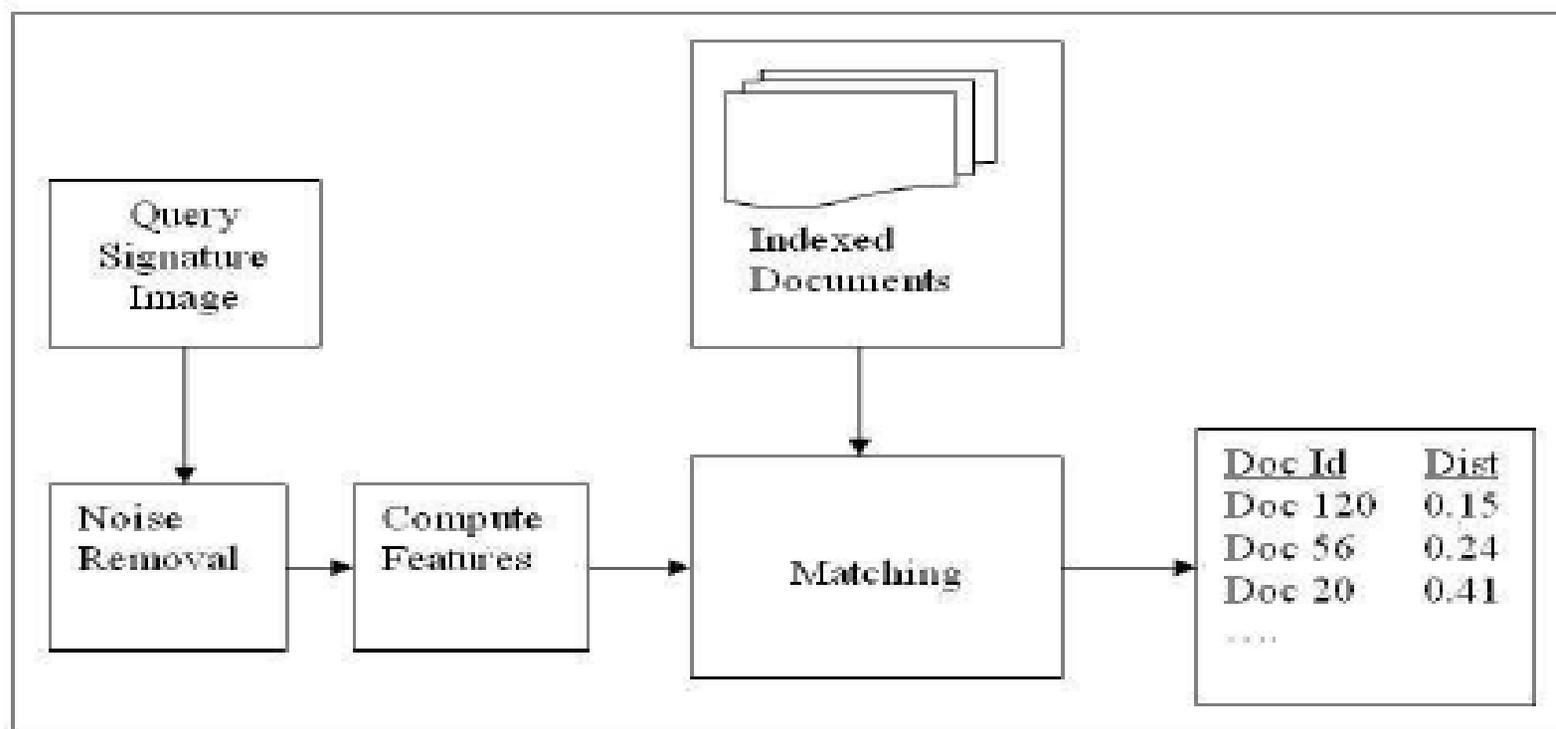


T.E. Whitehair, Jr.
Vice-President



T.E. Whitehair, Jr.

Document Retrieval



Matching Algorithm

Correlation Similarity Measure -

$$S(X,Y) = \frac{1}{2} + \frac{S_{11}S_{00} - S_{10}S_{01}}{2((S_{10} + S_{11})(S_{01} + S_{00})(S_{11} + S_{01})(S_{00} + S_{10}))^{1/2}}$$

where S_{ij} , $(i,j) = \{0,1\}$, is the number of occurrences where pattern i occurs in the first binary vector and pattern j occurs in the second vector in the same position.

Query Expansion using Automatic Relevance Feedback – The signature returned as the top choice result is used to expand the query

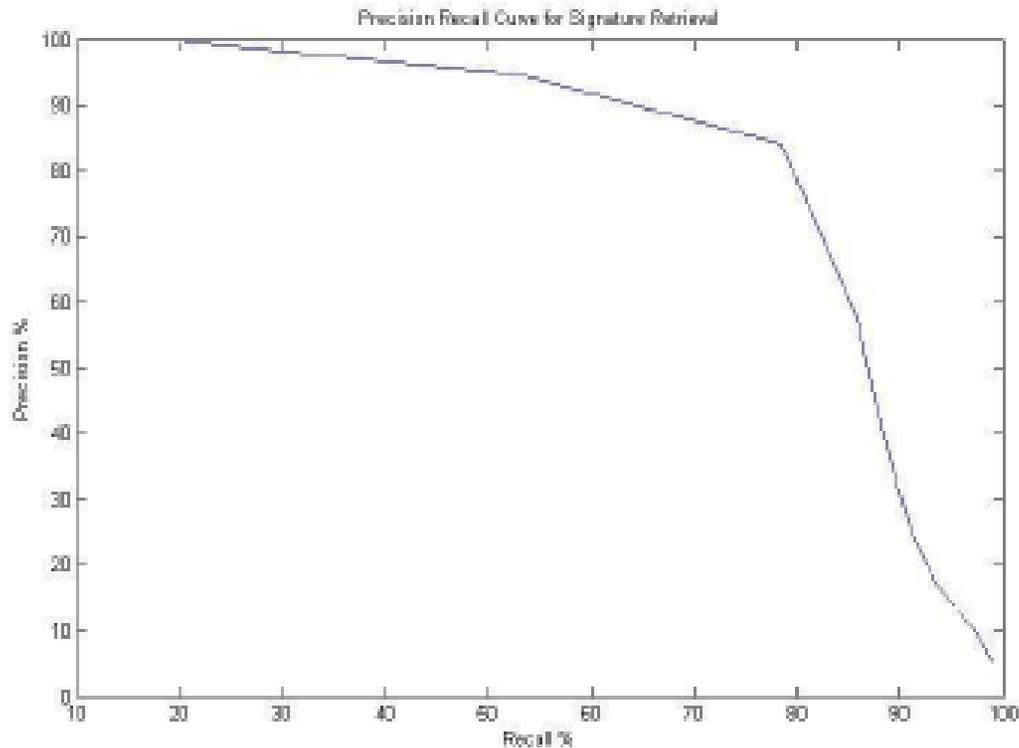
$$S_i = S(f(q), f(D_i))$$

D_i correspond to the document with the lowest S_i

q_{new} = signature image extracted from the document D_i

$$S_i(\{q, q_{new}\}, D_i) = \min\{S(q, D_i), S(q_{new}, D_i)\}$$

Precision-Recall Curves for signature retrieval results



Precision of 84.2% at Recall of 78.4% after Query Expansion

Conclusion

- The CRF model exploits spatial interdependencies in the data by using features from neighboring patches in addition to features specific to the current patch.
- 95.75% of the data were assigned the correct labels.
- The resulting labeled documents can be used in several applications.