

*Use of Conditional Random  
Fields for  
Signature-based Retrieval of  
Scanned Documents*

Harish Srinivasan, Shravya Shetty and  
Sargur Srihari

Center of Excellence for Document Analysis  
and Recognition (CEDAR)

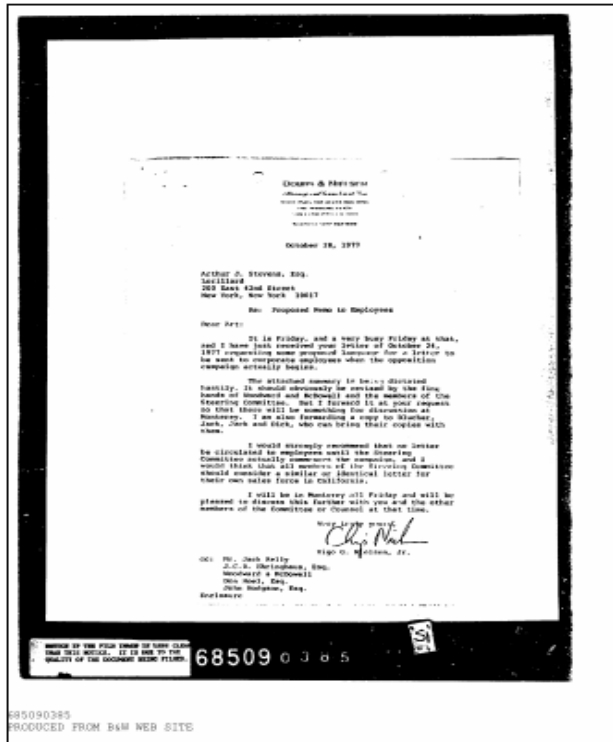
University at Buffalo  
State University of New York

# Overview

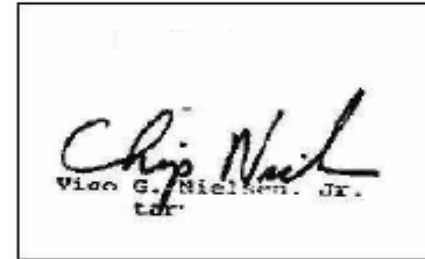
---

- Objective
- Signature Extraction
- Document Indexing
- Signature Retrieval
- Results
- Conclusion

# Objective 1 - Signature Extraction



Original Document



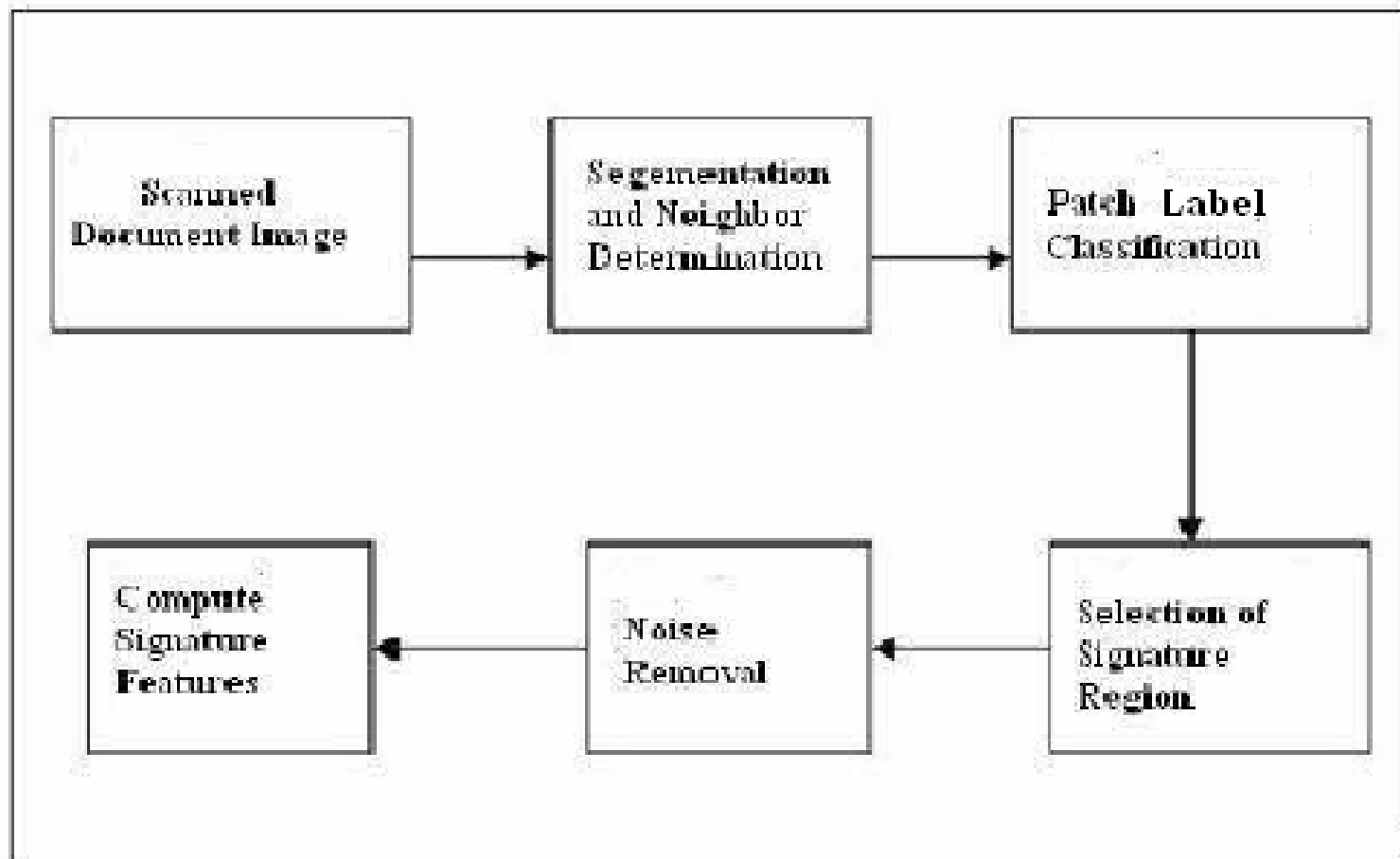
Extracted Signature

# Objective 2 - Signature Retrieval

---

	RETRIEVAL	Dissimilarity
<p>QUERY IMAGE</p> 		0.26
		0.36
		0.39
		0.43
		0.49

# Indexing



# Indexing Steps

---

- **Signature Block Location**
  - Segmentation and Neighbor determination
  - Features Extraction for Patch Label Classification
  - 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

# Segmentation

---

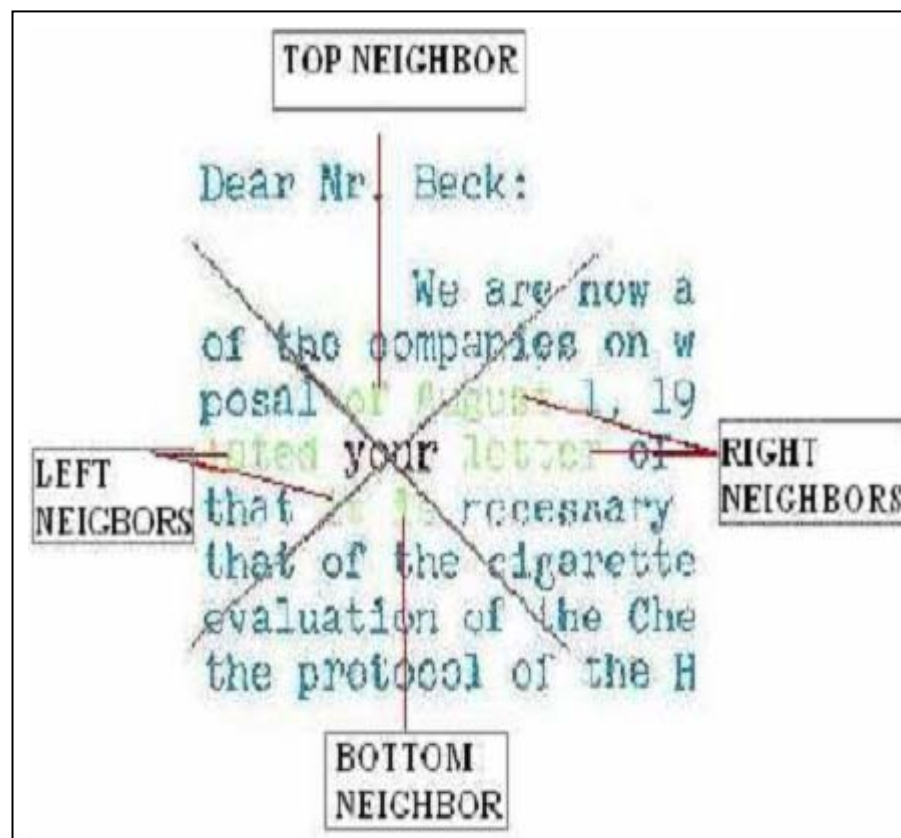
- The patches are generated using a region growing algorithm
- The size of each patch was optimized to represent approximately the size of a word

eral recommendations are accepted, detailed s  
It is anticipated that six or eight factor  
it manufacture during the week commencing Yay  
ted, that in view of the long time delay invo  
igarettes, market research codes be assigned to  
of a satisfactory product.

W. W. Mayhew

# 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.





# Conditional Random Field (CRF)

---

- **Model**

The probabilistic model of the CRF is given by

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

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

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

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

$m$  : 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

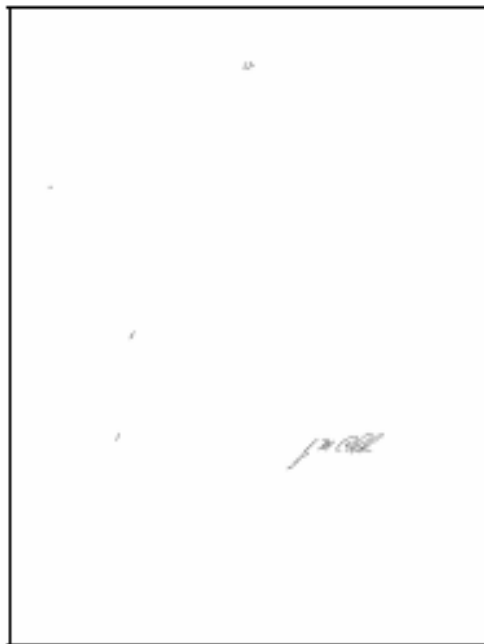
# Features for HW/Print/Noise Classification

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

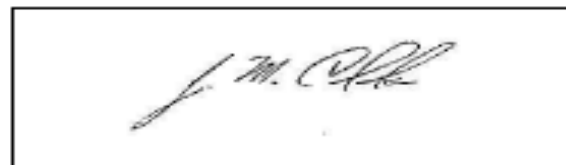
Transition Feature	Description
Relative location	Assigned weights based on the relative location - top/bottom or right/left
Convex hull distance	The convex hull distance between the 2 patches
Ratio of aspect ratio	The ratio of the aspect ratio values of the 2 patches
Ratio of number of components	The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

# Post-processing

---



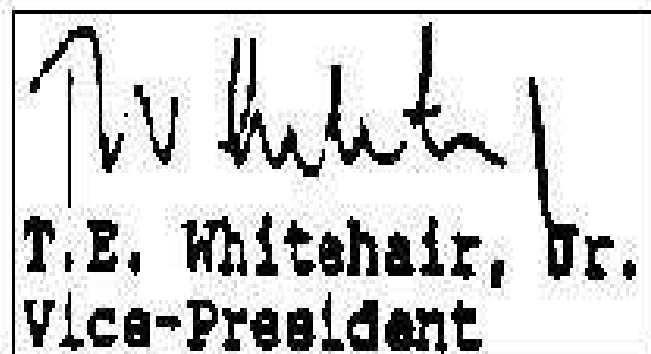
Handwritten components  
pulled out from CRF labels



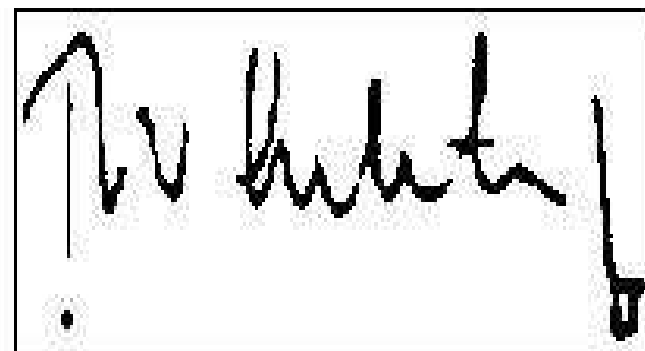
Extracted Signature

# Noise Removal

---



T.E. Whitehair, Dr.  
Vice-President

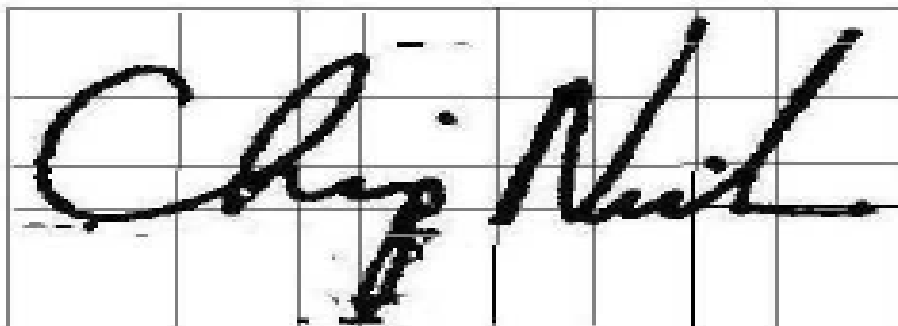


.

# Signature Features for Retrieval

8

4



Signature Image under a 4 x 8 Division

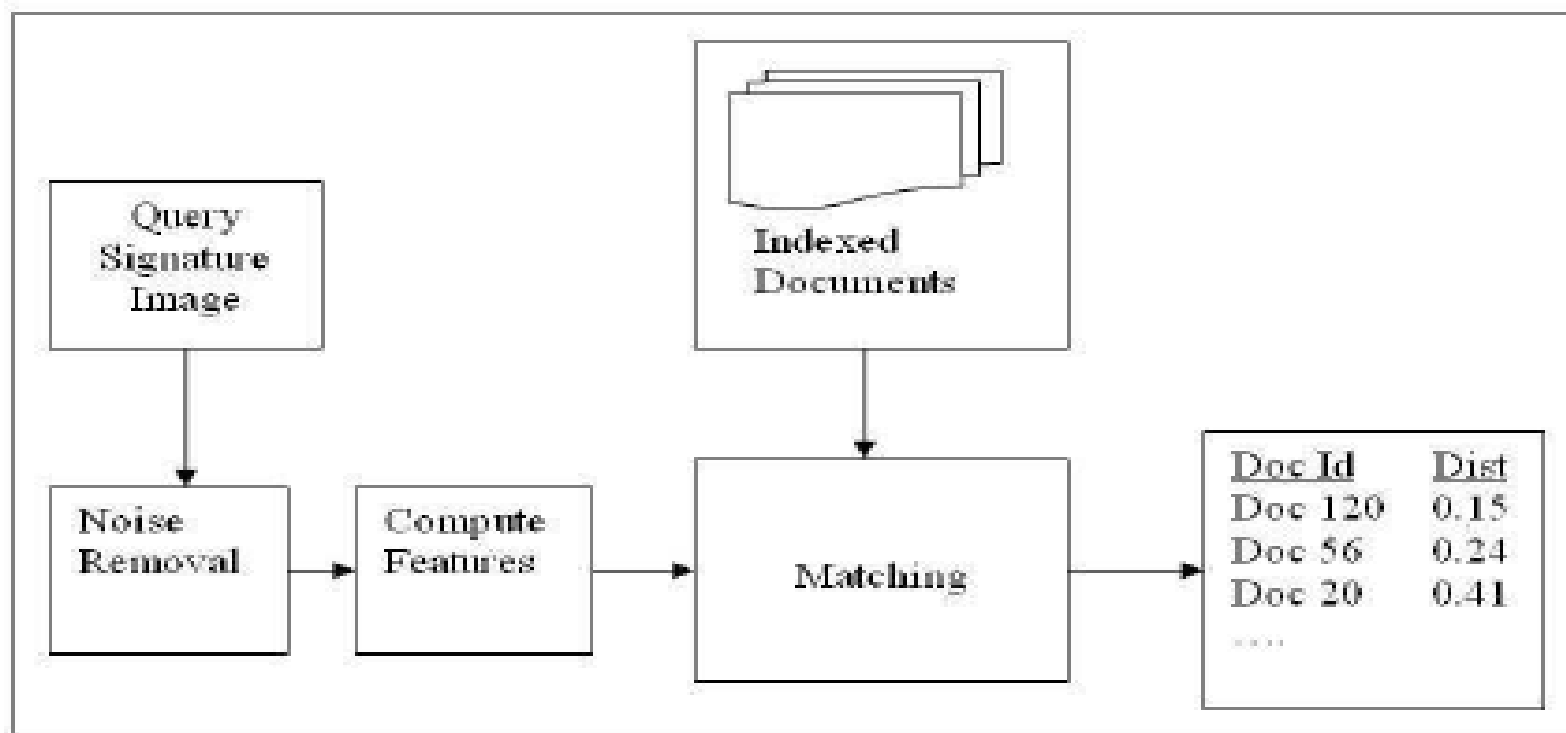
```
0011010011010110000110001110000110001111111110111111111
1000001110011110001110000000011100111100111101111001
111001111001011111011100001111000000111001111001111000
1110001111011110111000000000000101111001111111011111111
00100001111111111111111111000111000111011111000111111
11111100111100000000000000100111111111111110111111111
111111111111111111111101110111111001000001110000000
00111100001000000000110000000001111100111111000000
00000000000010111111000100101000001101011000000011
1100001101111110011011100000000000000000000111100
01100110000100001100000100101000000001100110000011
01000000111110001011111000110001110000011110100
00011101000001100000000000000000100000110001011
1111001111100110000011000000000000000000000000000000
0000000000001000101011101111110001000000000000010
1000000011000010000000000010010110001100010000000
0000000000000011001010010000000000010000011000000
11001000000000110011011100100110000000000000000000000
0000100101000001000000000000000001000000100000
```

1024 bit Binary Feature Vector  
Extracted

Gradient (12 bits): Structural (12 bits): Concavity (8 bits)

# 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 -

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

$D_i$  correspond to the document with the lowest  $S_i$

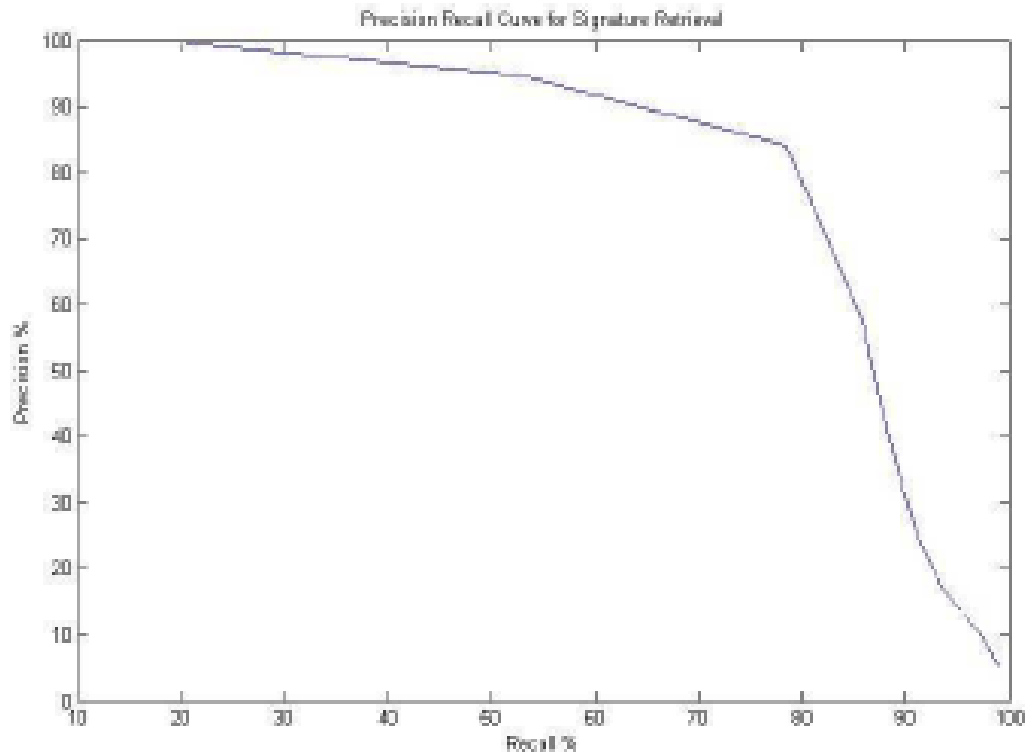
$q_{new} =$  signature image extracted from the document  $D_i$

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



# Precision-Recall Curves for signature retrieval results

---



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

# Conclusion

---

- **Signature Extraction is Robust**
  - due to modeling neighboring spatial dependencies in the labels as well as the observed data using CRF
- **Retrieval has satisfactory Performance**
  - Precision and Recall of 84% and 78% respectively.