



cse@buffalo

Indexing and Retrieval of Handwritten Document Images

Huaigu Cao

hcao3@cubs.buffalo.edu



cse@buffalo

Outline

- Problem Statement
- Proposed Approach
- Results
- Future Works

Problem Statement

*Information
Retrieval
Techniques*

*Document
Analysis and
Recognition
Techniques*

machine-print (99%)
handwritten (good
quality: 50-70%; poor
quality: 20-40%)

Document image
retrieval:
✓ machine-print
✗ handwritten

Text Retrieval:
webpage, library
catalogs, etc

- Search engine for handwritten document images performs poorly because of the low recognition rate
- Objective: Improve IR performance on handwritten data, especially that of poor quality



Document Retrieval – An Overview

- The similarity between a document and a query can be roughly measured by the number of occurrences of the query terms in the document, for example:

$$\{d_j\} = \{d_1 = \text{"pt has a trauma"}, d_2 = \text{"pt has breath difficulty"}\}$$

$$\{t_i\} = \{t_1 = \text{"pt"}, t_2 = \text{"has"}, t_3 = \text{"a"}, t_4 = \text{"trauma"}, t_5 = \text{"breath"}, t_6 = \text{"difficulty"}\}$$

$$freq_{i,j} = \begin{bmatrix} 1, 1, 1, 1, 0, 0 \\ 1, 1, 0, 0, 1, 1 \end{bmatrix} \quad q = \text{"breath difficulty"}, qtf_{i,q} = [0, 0, 0, 0, 1, 1]$$

- The frequency $freq_{i,j}$ is the most crucial for a document retrieval system

Estimating $freq_{i,j}$ from Noisy Handwriting Recognition Results

- Index built on top- n Word Recognition Results [Lee, ICDAR05]
- Assigning different weights to candidates

$$freq_{i,j} = \sum_k \Pr(t_i | w_k) \quad [\text{Rath04}]$$

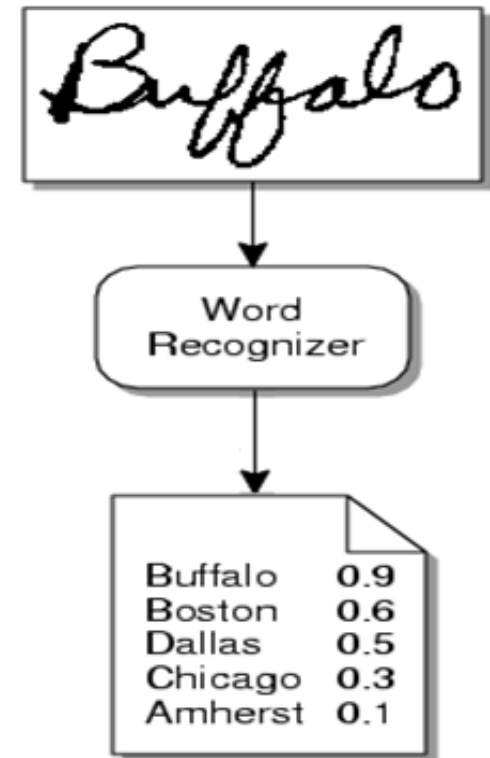
or

$$freq_{i,j} = \sum_k \frac{Const}{\text{RANK}(t_i, w_k)} \quad [\text{Howe05}]$$

w_1, w_2, \dots : word images in document d_j

$\Pr(t_i | w_k)$: word recognition posterior probability

- Perfect word segmentation?





.cse@buffalo

Outline

- Problem Statement
- Proposed Approach
- Results
- Future Works



Improved Estimation of $freq_{i,j}$ Incorporating Word Segmentation and Language Models [Our Approach]

$$\hat{freq}_{i,j} = \sum_{\vec{w}} \Pr(\vec{w} | \vec{o}) \cdot \sum_{\vec{\tau}} \Pr(\vec{\tau} | \vec{w}) \cdot freq_{i,\vec{\tau}}$$

$\vec{o} = [o_1 o_2 \dots o_T]$: observation series

$\vec{w} = [w_1 w_2 \dots w_L]$: word images

$\vec{\tau} = [\tau_1 \tau_2 \dots \tau_L]$: terms (OCR results)

$\Pr(\vec{w} | \vec{o})$: word sequence segmentation probability

$\Pr(\vec{\tau} | \vec{w})$: word sequence recognition probability

$freq_{i,\vec{\tau}}$: number of t_i in $\vec{\tau}$

Solved by *dynamic programming*

UB Word Gap Labeling

.cse@buffalo

Features for each gap between two CC's

- fv_1 : Euclidian distance between bounding boxes
- fv_2 : Shortest white run length between two CC's
- fv_3 : Distance between the convex hulls



By Bayes' rule

$$\Pr(\text{Valid} | \vec{fv}) = \frac{\Pr(\text{Valid}) p(\vec{fv} | \text{Valid})}{\Pr(\text{Valid}) p(\vec{fv} | \text{Valid}) + \Pr(\text{Non-valid}) p(\vec{fv} | \text{Non-valid})}$$

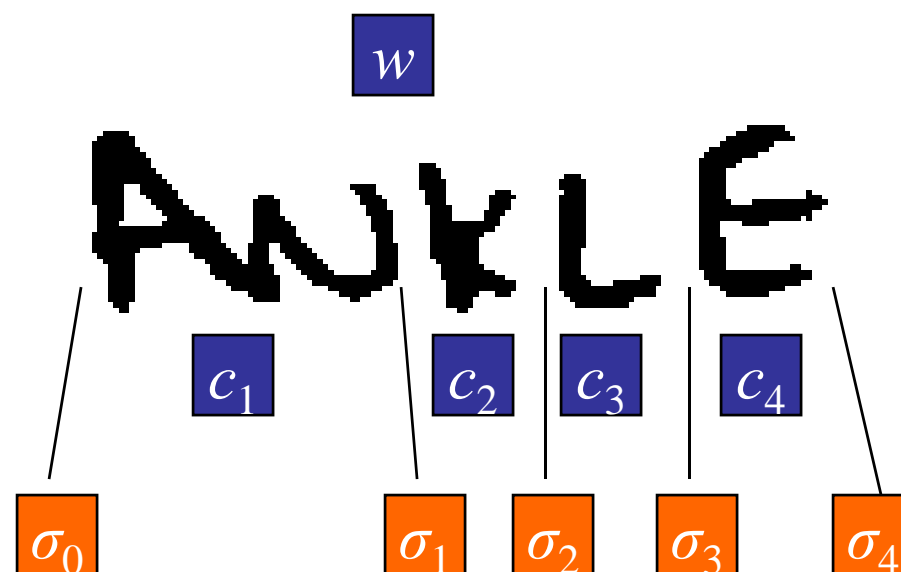
The likelihoods are estimated non-parametrically using *Parzen window*

Word Spotting Using Word Gap Labeling

Word image (w) is represented by connected components :

$$w \rightarrow c_i, c_{i+1}, \dots, c_j$$

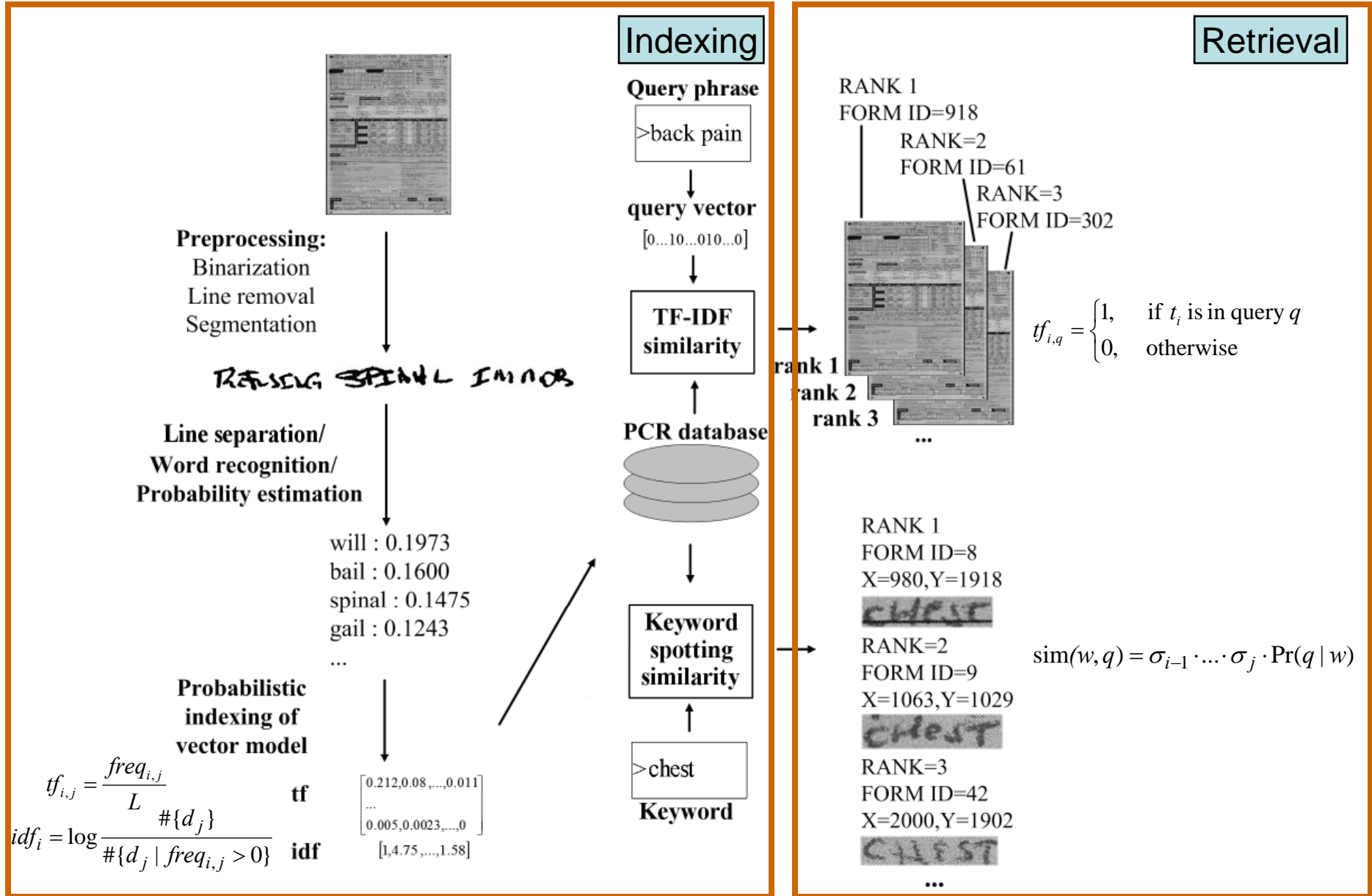
Word gap probability (σ_k):
The probability of the gap between c_k and c_{k+1} being a word gap



Word - query similarity :

$$\text{sim}(w, q) = \sigma_{i-1} \cdot \dots \cdot \sigma_j \cdot \Pr(q | w)$$

Search Engine Diagram





cse@buffalo

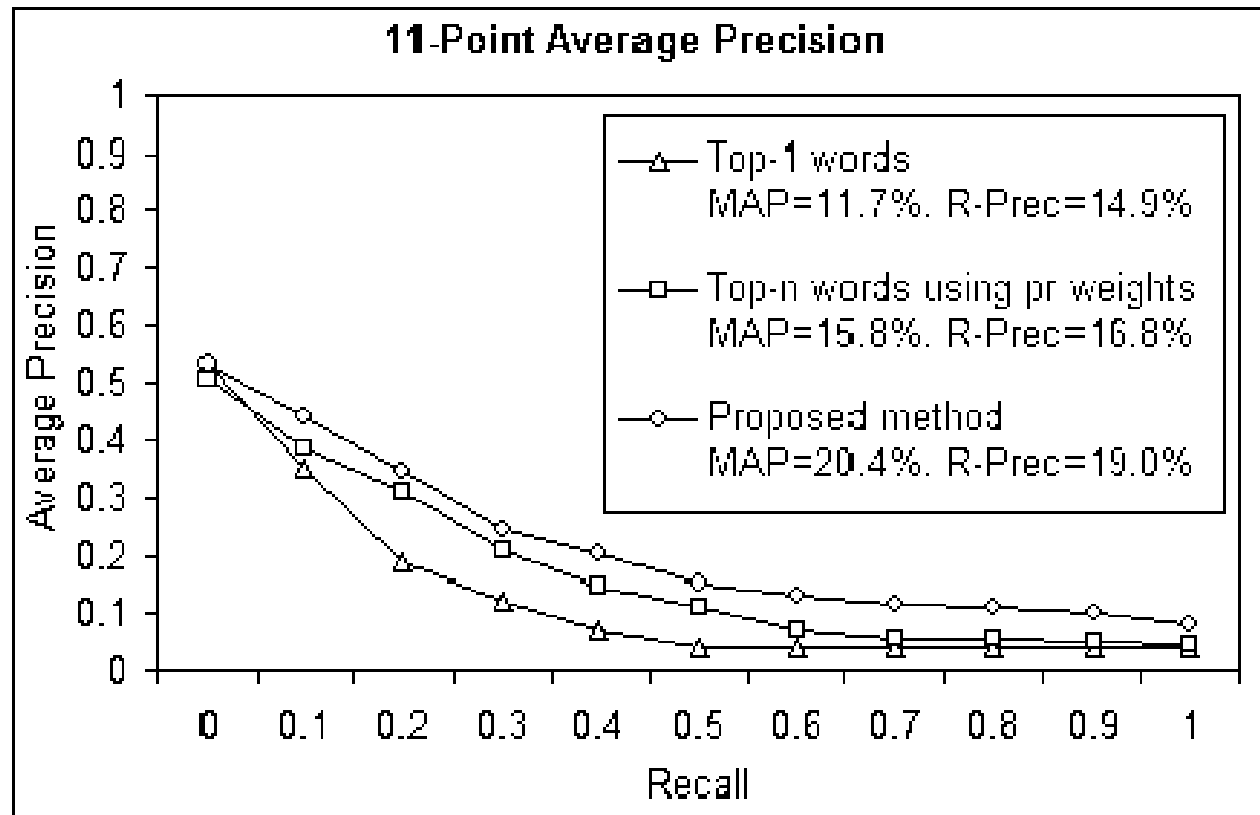
Outline

- Problem Statement
- Proposed Approach
- Results
- Future Works



Document Retrieval Performance: Recall – Precision

342 Documents; 28 Queries

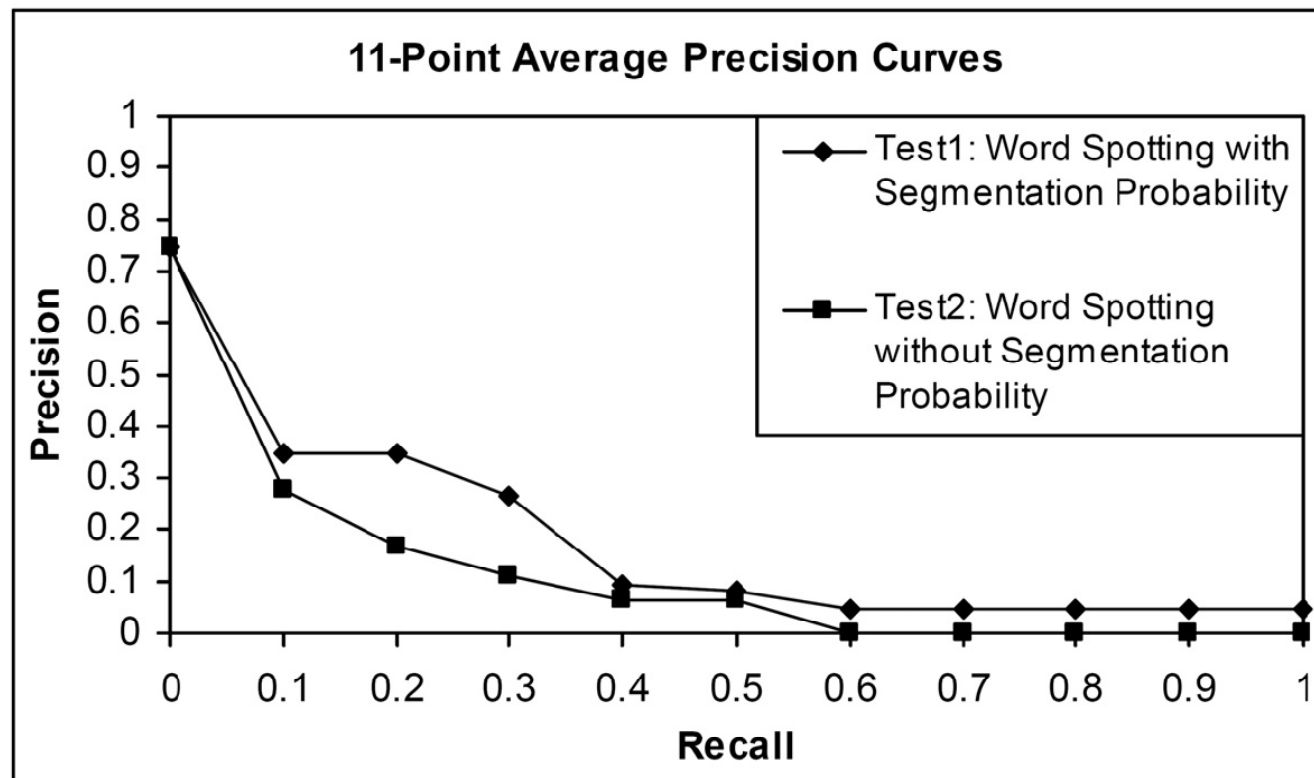




cse@buffalo

Word Spotting Performance: Recall – Precision

342 Documents; 33 Queries





.cse@buffalo

Outline

- Problem Statement
- Proposed Approach
- Results
- Future Works



Future Work

- Investigation of other high-level applications using the ranked OCR output
 - Machine translation for the handwritten documents
 - Handwritten document categorization