

# Digital Access of Handwritten Documents

**Venu Govindaraju**

Anurag Bhardwaj

Huaigu Cao

[Venu@cubs.buffalo.edu](mailto:Venu@cubs.buffalo.edu)

# Outline

- **Recognition**
  - Postal Application
  - Paradigms
  
- **Search**
  - OCR accuracy
    - Fusion
    - Lexicon Reduction
    - Statistical Topic Models
  - Document Search
  - Word Spotting

## Challenge of Handwriting

<sup>the</sup>  
S. Winchester: October 28<sup>th</sup> 1755.  
Parole Hampton.

The officers who came down from Fort Cumberland with Colonel Washington, are immediately to go Recruiting; and they are allowed until the 5<sup>th</sup> of December; at which time if they do not punctually appear at the place of Rendezvous assigned them, they will be tried by a Court Martial, for disobedience of Orders.

Ich hätte Ihre Bitte gewiss früher erfüllt, wenn Ihre Briefe nicht in der Hülle ihrer Geschwister aus meinem Blickfeld verschwunden wären. Denn Graphologie hat mich stets interessiert, wenn ich mich auch nie systematisch dafür interessiert habe, ich mache es mit Schriften wie mit Gesichtern, ich urteile instinktiv.

Hochachtungsvoll

A. Einstein.

A FOUND IN WHEELCHAIR, AT  
THE LAST MOMENT OF DEATH, HE  
WANTED HER TO BE EQUAL

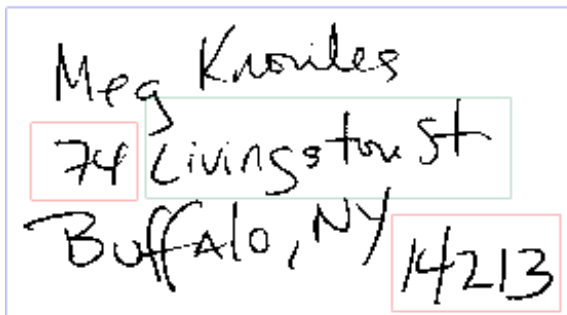
## Motivation

- **Vast, irreplaceable, culturally vital legacy collections of historical documents are competing ineffectively for attention with billions of digital documents**
- **Thus historical archives are threatened with neglect, perceived irrelevance, .... & eventually, oblivion?**

**Threat: *'If it's not in Google, it doesn't exist!'***

# Postal Context (138 mil records)

Typical Mail-piece Address Block :



ZIP and Primary Number Recognition

ZIP Code: 14213 Database Query  
 Primary Number: 74

Lexicon Entry (Street Name)	Addon
BRADLEY ST	1024
COLONIAL CIR	1467
DEWITT ST	1534
GRANT ST	1905
HAMPSHIRE ST	2012
HAWLEY ST	1053
LAFAYETTE AVE	1347
LIVINGSTON ST	1653
MANCHESTER PL	1251
PO BOX	0074
POOLEY PL	1128
POTOMAC AVE	1146
PUTNAM ST	1650

Delivery Point Code: 14213165374 ← Word Recognition (Street Name) And Delivery Point Code Derivation

Lex	Top 1	Top 2
10	96.5	98.7
100	89.2	94.1
1000	75.3	86.3

**LDR**

## ZIP Code

- 30% of ZIP Codes contain a single street name
- 5% of ZIP Codes contain a single primary number
- 2% of ZIP Codes contain a single add-on

### <ZIP Code, primary number>

- Maximum number of records returned is 3,071

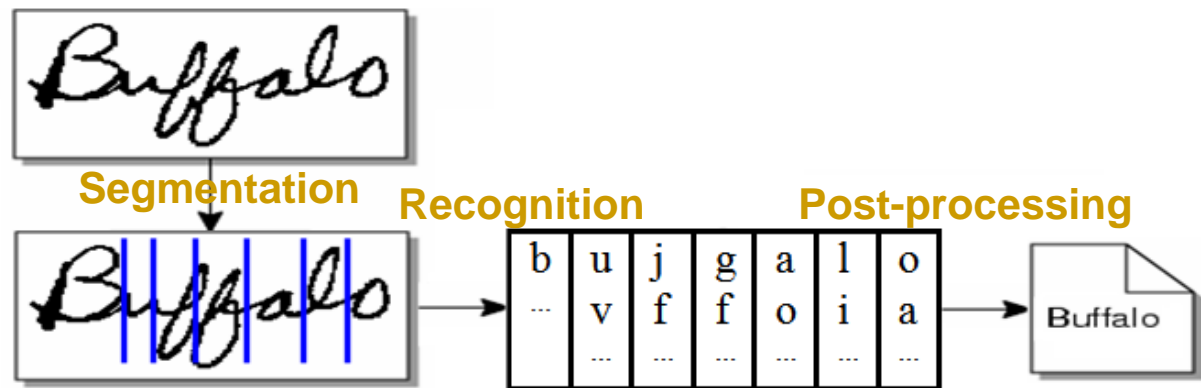
### <ZIP Code, add-on>

- Maximum number of records returned is 3,070

# Paradigms

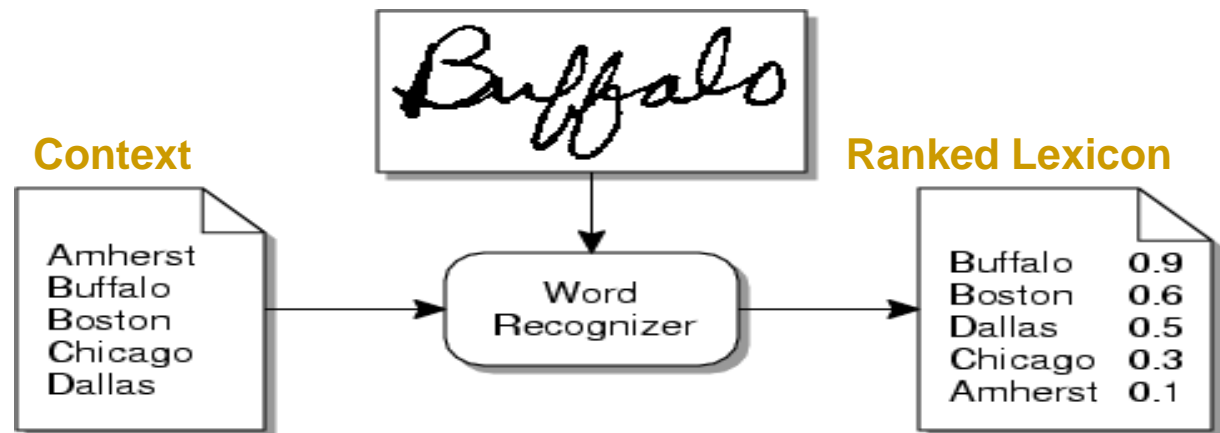
Lexicon Free OCR

LFR

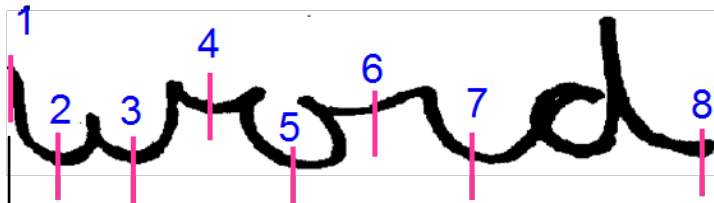


Lexicon Driven OCR

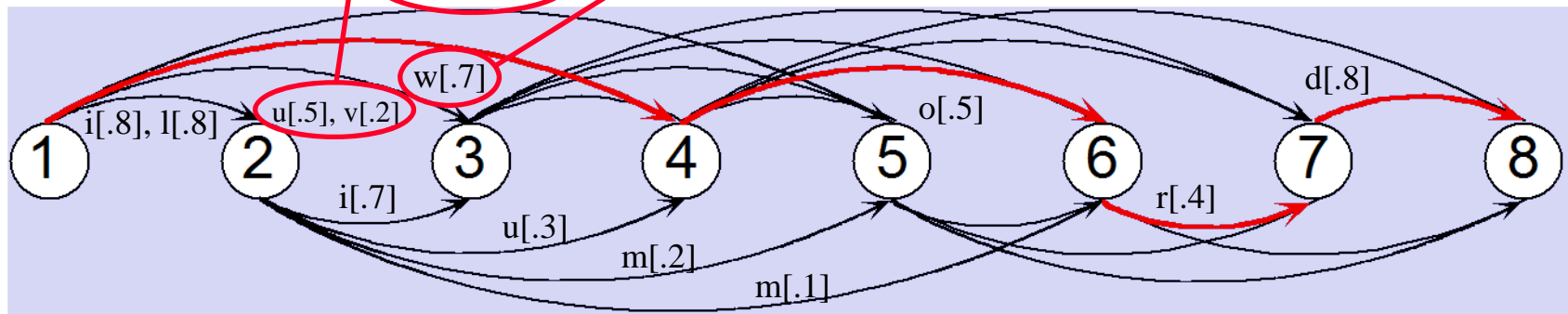
LDR



## Lexicon Free (LFR)

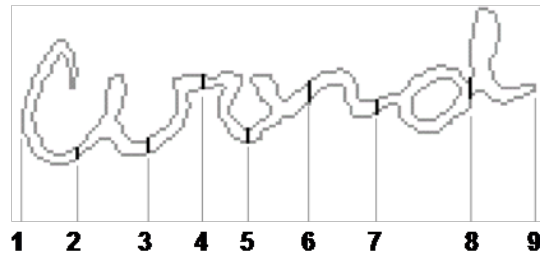


- Image from 1 to 3 is a **i** with 0.5 confidence
- Image from segment 1 to 4 is a **w** with 0.7 confidence
- Image from segment 1 to 5 is a **w** with 0.6 confidence and an **m** with 0.3 confidence



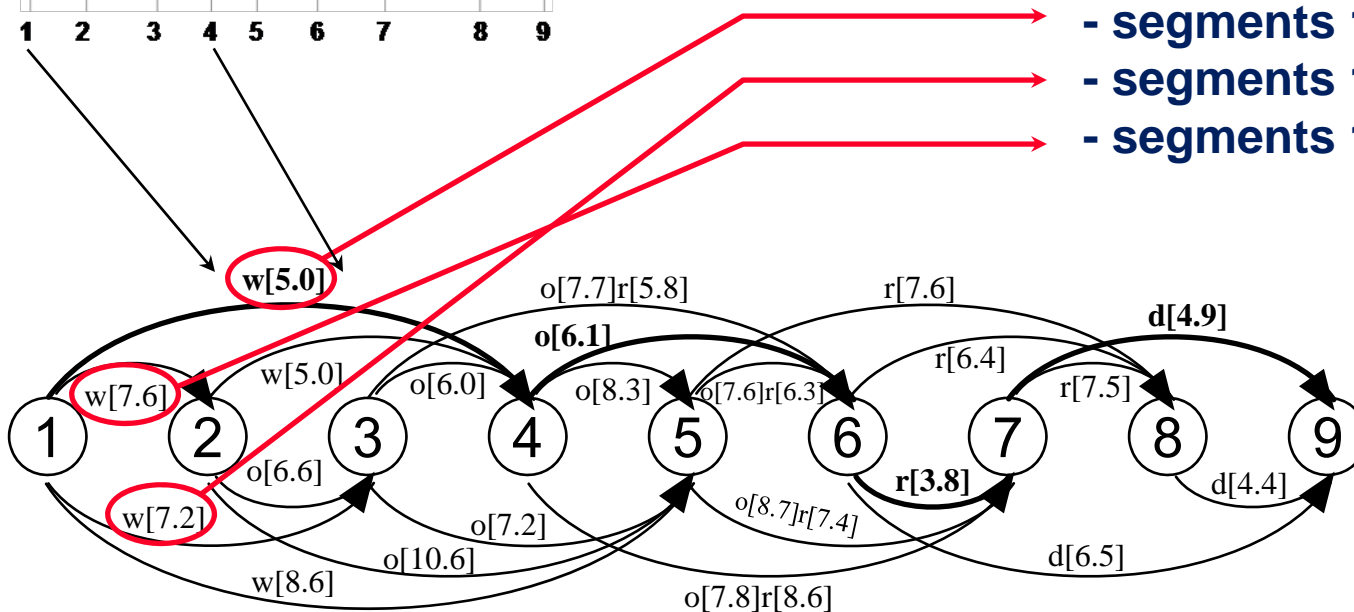
**Find the best path in graph from segment 1 to 8**

## Lexicon Driven (LDR)



Distance between lexicon entry 'word' first character 'w' and the image between:

- segments 1 and 4 is 5.0
- segments 1 and 3 is 7.2
- segments 1 and 2 is 7.6



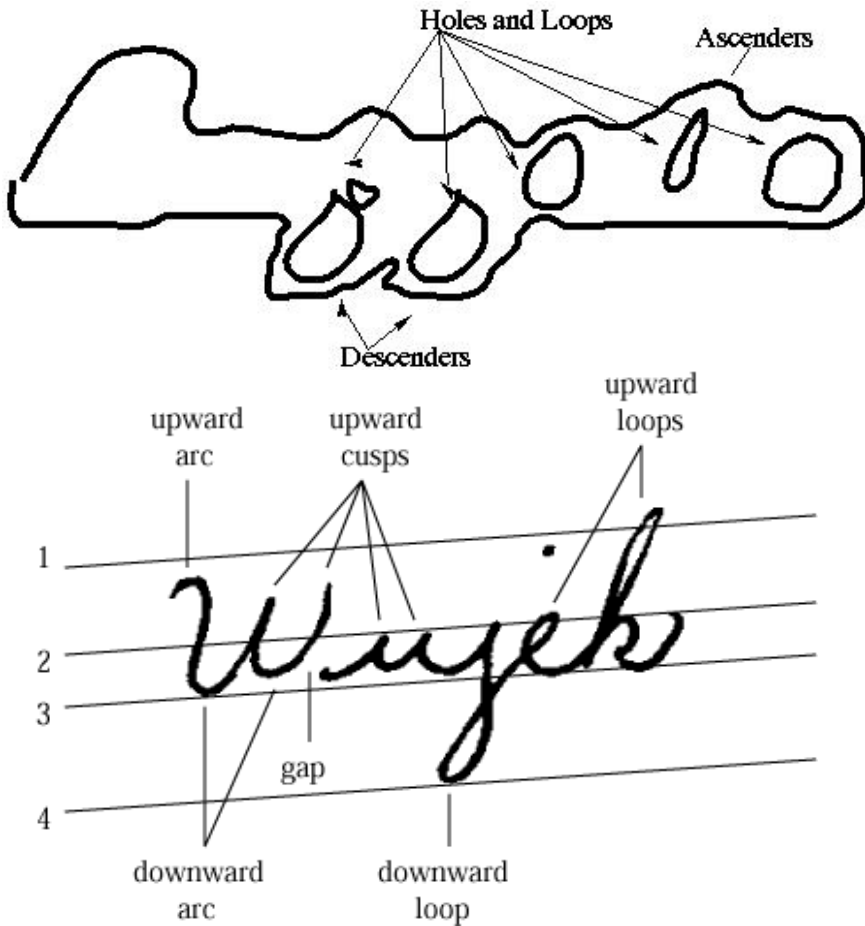
Find the best way of accounting for characters 'w', 'o', 'r', 'd' by consuming all segments 1 to 8

## Interactive Models (LDR)

### 2-way interaction

- a) Amherst
- b) Buffalo
- c) Boston

- a) San Jose
- b) Buffalo
- c) Washington
- d) None of the above



## Search for Handwritten Documents

- Lexicons are typically large: >5K
- Need around 70% accuracy

Lexicon	Good Quality		Historical		Medical
	10K	1K	10K	1K	4K
Top 1 (%)	57	67	12	28	20
Top 3 (%)	69	72	22	44	27
Top 10 (%)	74	75	32	72	42

### Strategy

- Reduce lexicon size using **topic categorization** (DAS 06;08)
- Use **Top-N** choices returned by OCR (ICDAR 07)

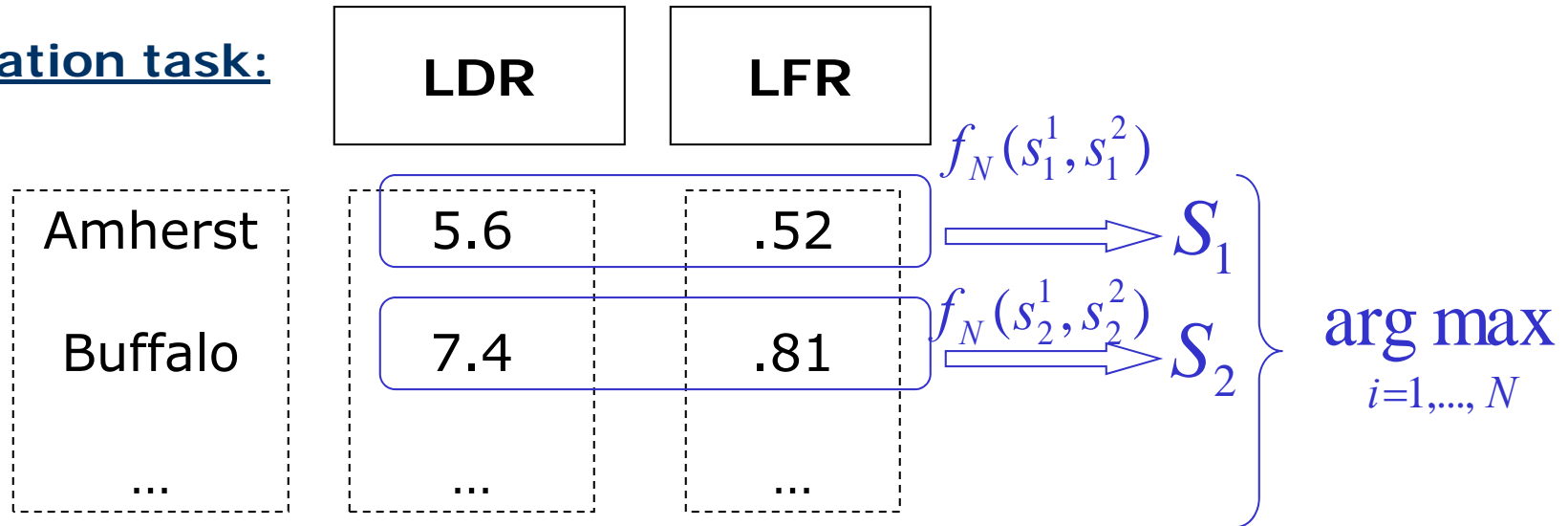
# Outline

- **Recognition**
  - Postal Application
  - Paradigms
- **Search**
  - OCR accuracy
    - Fusion
    - Lexicon Reduction
    - Statistical Topic Models
  - Document Search
  - Word Spotting

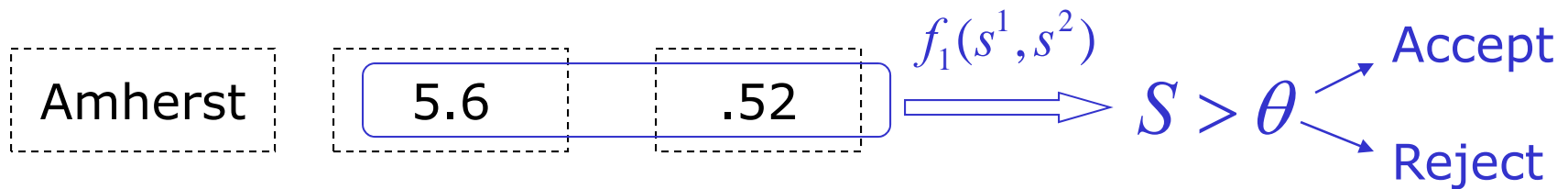
# Fusion of Recognizers

## Type III

### Identification task:



### Verification task:



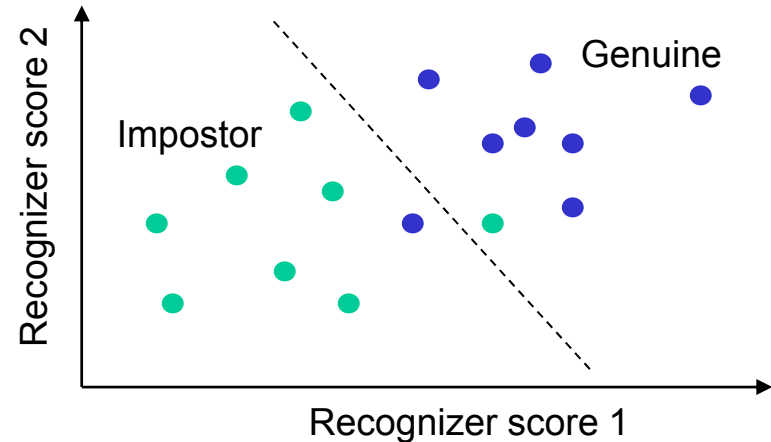
$$f_N \equiv f_1 ?$$

## Traditional Fusion Rules

- Sum rule  $f_1(s^1, s^2) = s^1 + s^2$
- Weighted sum rule  $f_1(s^1, s^2) = w_1 s^1 + w_2 s^2$
- Product rule  $f_1(s^1, s^2) = s^1 \times s^2$
- Max rule  $f_1(s^1, s^2) = \max(s^1, s^2)$
  
- Rank-based methods  $s_i^1 \rightarrow r_i^1 = \text{rank}(s_i^1, \{s_1^1, \dots, s_N^1\})$   
 $f_1(s_i^1, s_i^2) = r_i^1 + r_i^2$   
 $f_1(s_i^1, s_i^2) = P(r_i^1, r_i^2 \mid \text{gen})$

# Likelihood Ratio Verification Tasks $f_V$

- 2 classes: imposter and genuine
- Pattern classification task



**Minimum risk criteria: optimal decision boundaries coincide with the contours of likelihood ratio function:**

$$f_{lr}(s^1, s^2) = \frac{p_{gen}(s^1, s^2)}{p_{imp}(s^1, s^2)} \quad f_V = f_{lr}$$

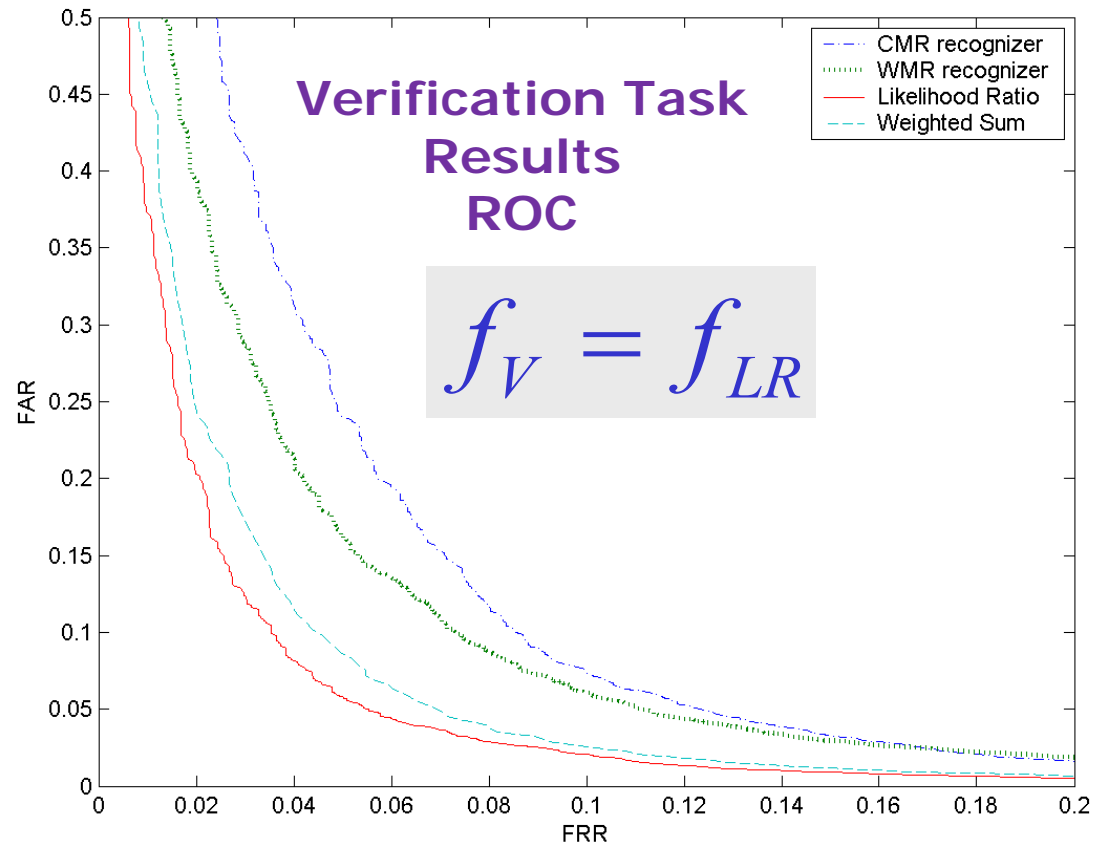
**Metaclassification with NN, SVM, etc. also possible**

# Optimal Combination functions

## Identification Task Results

Top choice correct rate

LFR is correct	54.8%
LDR is correct	77.2%
Both are correct	48.9%
Either is correct	83.0%
<b>Likelihood Ratio</b>	<b>69.8%</b>
<b>Weighted Sum</b>	<b>81.6%</b>

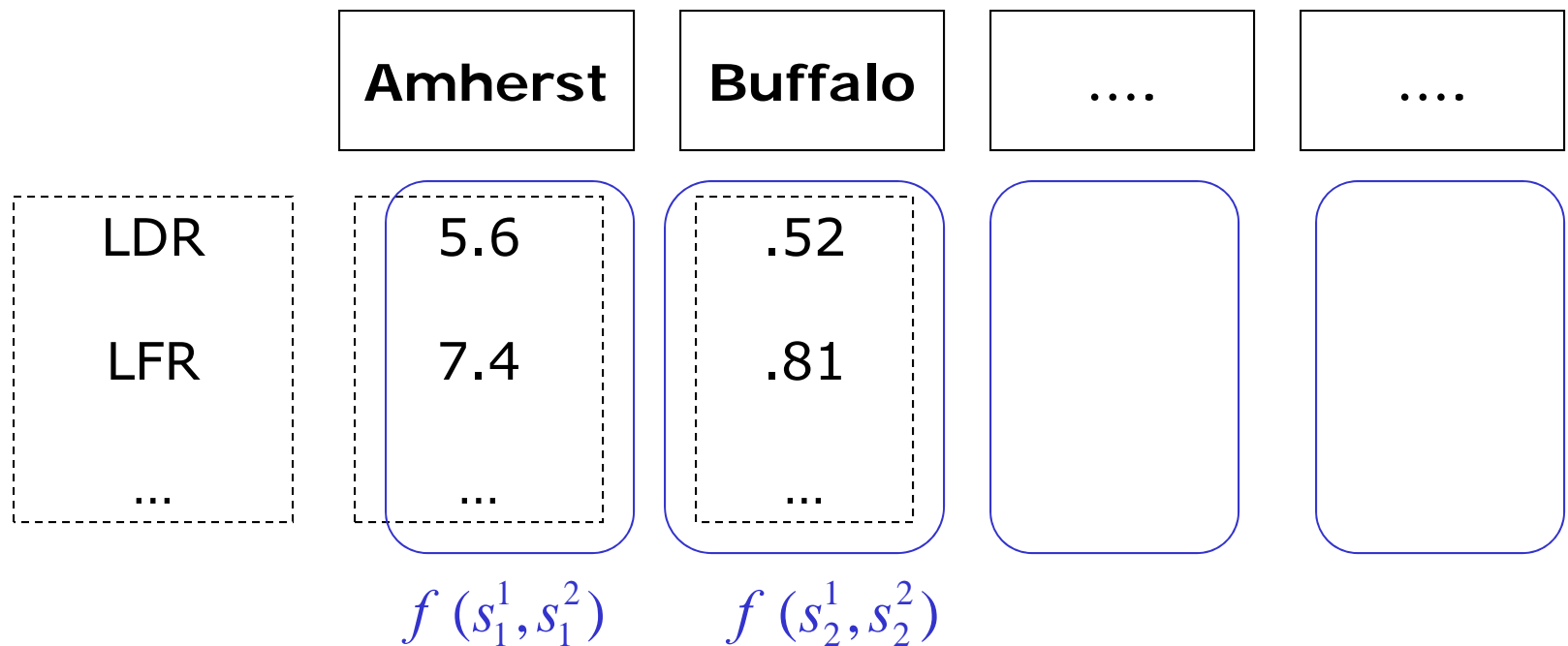


- LR combination  $f_V$  is worse than single matcher

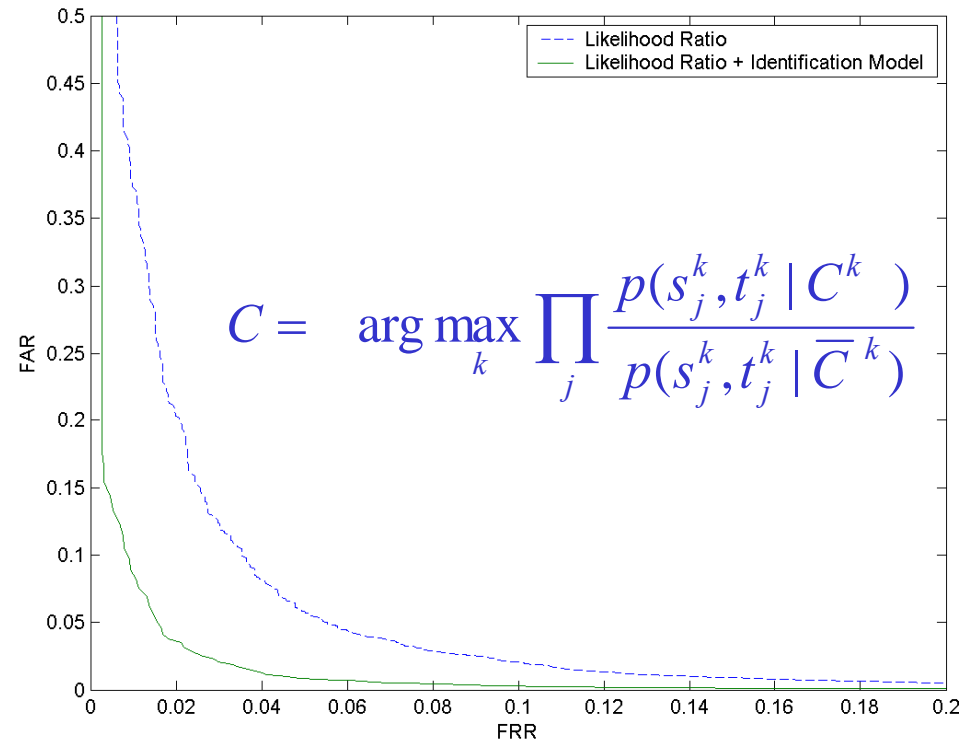
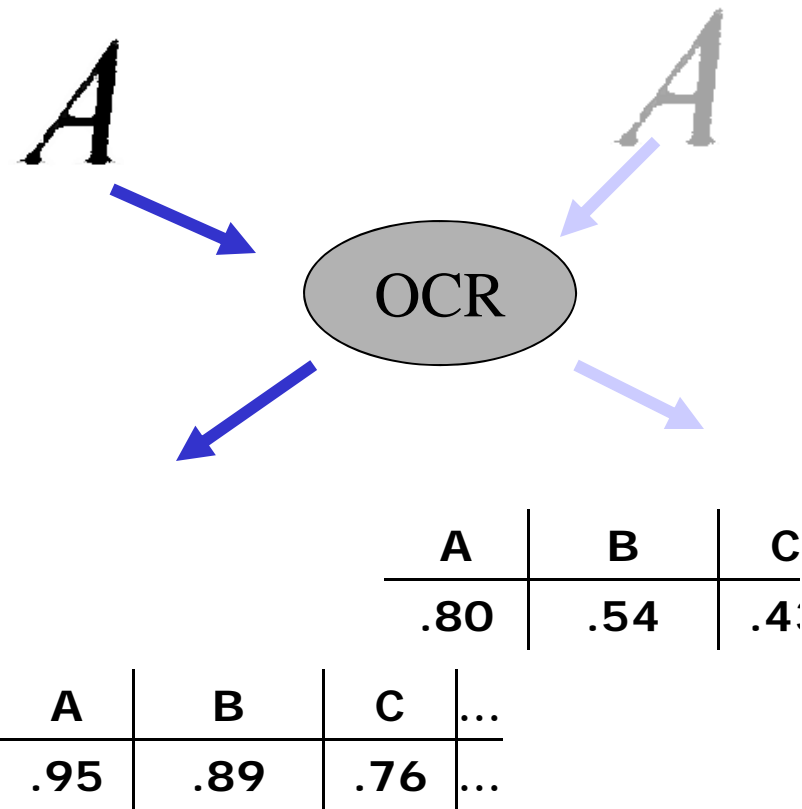
# Independence of Scores

In a single trial

$$S_i = f(s_i^1, s_i^2, \dots, s_i^M, \{s_k^1, s_k^2, \dots, s_k^M\}_{k \neq i})$$



# Dependencies



## Iterative Methods

- Initialize a combination function  $f(s^1, s^2, \dots, s^M)$
- Get scores from the same identification trial (for all trials)
- Update function so Genuine score better than any impostor score

**Best Impostor Function**

$$f() = \frac{p_{gen}(s_i^1, s_i^2, \dots, s_i^M)}{p_{imp}(s_i^1, s_i^2, \dots, s_i^M)}$$

**Sum of Logistic Functions**

$$f() = \frac{1}{1 + e^{-(\alpha_1 s^1 + \alpha_2 s^2 + \dots + \alpha_M s^M + \alpha_{M+1})}}, \alpha_j \geq 0$$

	Likelihood Ratio	Weighted sum	Best Impostor Likelihood Ratio	Logistic Sum	Neural Network
LFR & LDR	69.84	81.58	80.07	81.43	81.67

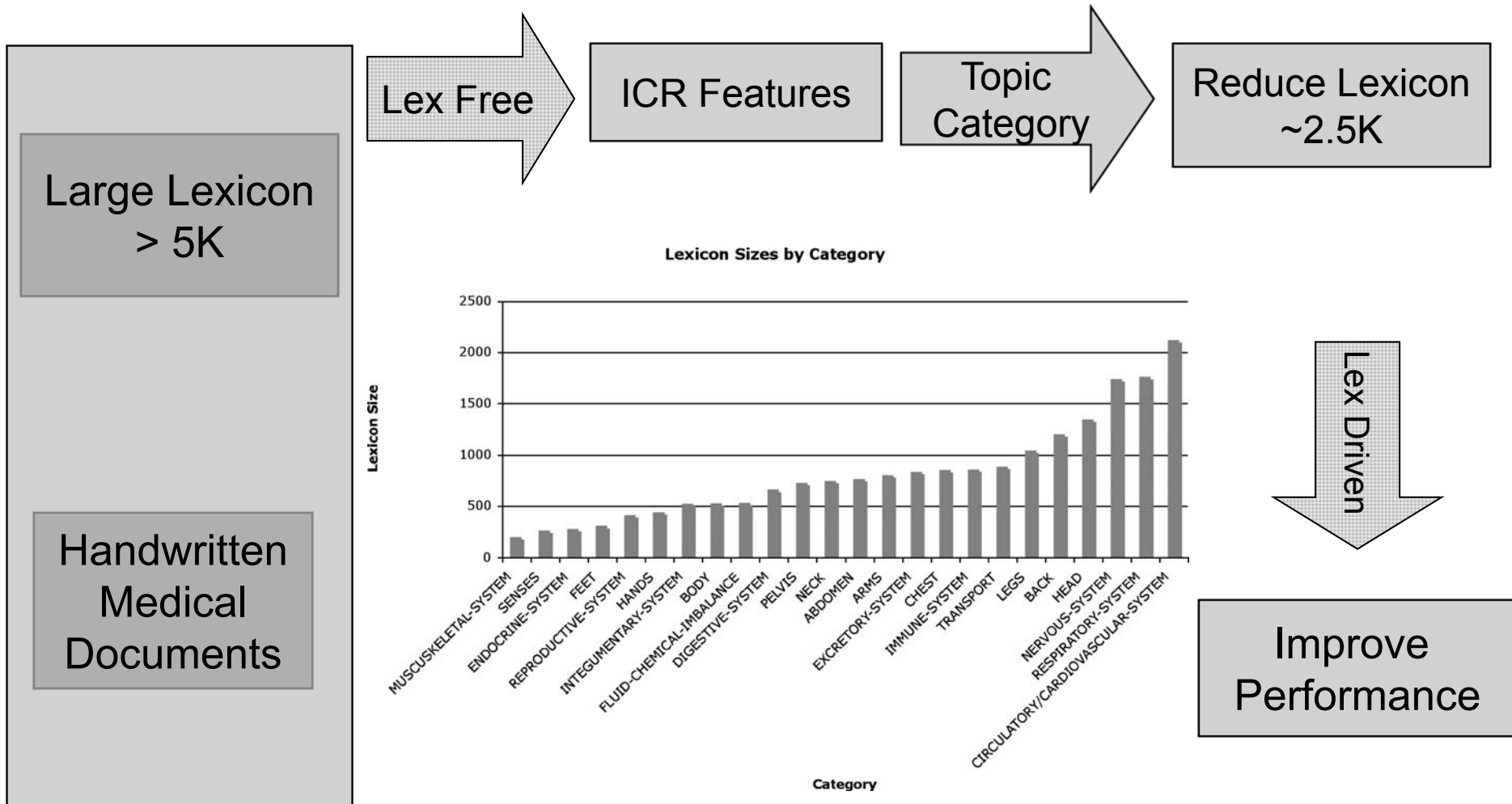


## Search Engine for Medical Forms

TIME	RESP	PULSE	B.P.	LEVEL OF CONSCIOUSNESS	GCS	R	PUPILS	L	SKIN	STATUS
1320	Rate: 20 <input checked="" type="checkbox"/> Regular <input type="checkbox"/> Shallow <input type="checkbox"/> Labored	Rate: 74 <input checked="" type="checkbox"/> Regular <input type="checkbox"/> Irregular	130 P	<input checked="" type="checkbox"/> Alert <input checked="" type="checkbox"/> Voice <input type="checkbox"/> Pain <input type="checkbox"/> Unresp.	15	<input checked="" type="checkbox"/>	<input type="checkbox"/> Normal <input type="checkbox"/> Dilated <input type="checkbox"/> Constricted <input type="checkbox"/> Sluggish <input type="checkbox"/> No-Reaction	<input type="checkbox"/>	<input type="checkbox"/> Unremarkable <input type="checkbox"/> Cool <input checked="" type="checkbox"/> Warm <input type="checkbox"/> Moist <input checked="" type="checkbox"/> Dry <input type="checkbox"/> Pale <input type="checkbox"/> Cyanotic <input type="checkbox"/> Flushed <input type="checkbox"/> Jaundiced	<input type="checkbox"/> C <input type="checkbox"/> U <input checked="" type="checkbox"/> P <input type="checkbox"/> S
<b>OBJECTIVE PHYSICAL ASSESSMENT</b>										
84 y/o ♂ COA/B HEENT NO trauma/ABNER/PEARL										
Pt. is hard of hearing, Chest: ⊕ Bilat. Exp. ABD: soft x4										
Ext: no trauma ⊕ pulse/motion/neuro. Some sweating noted										
<b>COMMENTS</b>										
D stick 74 Pt. Tx supine w/ ALS during Tx.										
p. TX ER RM <sup>3</sup> w/ care → ER Staff										

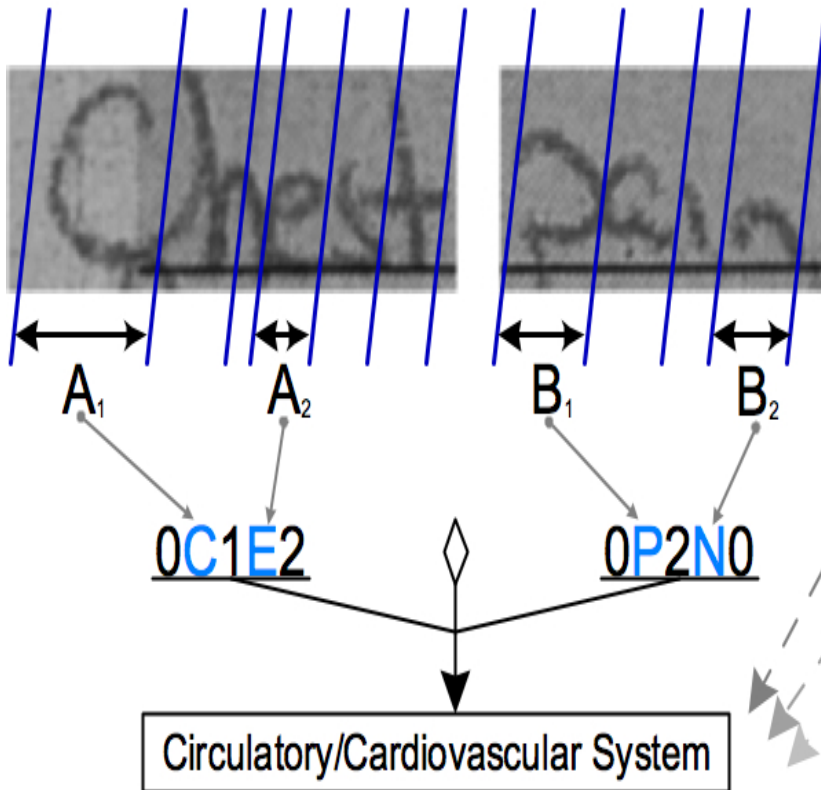
- Find all people who reported asthma problems in NY
- How many people with high blood pressure are on medication X?
- Is there an epidemic breaking?

# Lexicon Reduction



# ICR Features Index

Cohesive Phrase: CHEST PAIN

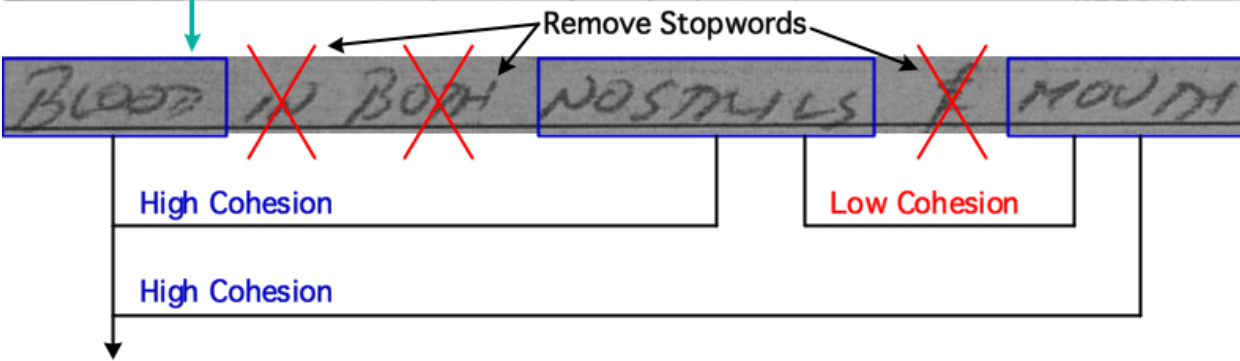


(ID: 342-2)		*C*S*	*A*N*
(ID: 407-1)		*C*M*	*P*A*
(ID: 407-2)		*H*O*	*A*W*
(ID: 473)		*C*F*	*L*A*
(ID: 643)		*C*H*	*A*W*
(ID: 695-1)		*L*S*	*L*R*
(ID: 98)		*C*S*	*S*S*
(ID: 606)		*C*W*	*D*M*
(ID: 695-2)		*C*F*	*W*L*

Many terms mapped to category

# Topic Features

OBJECTIVE PHYSICAL ASSESSMENT FOUND SUPINE ON GROUND - PE → CAO x 3 - PEARL - SMALL LAC ON NOSE  
 BLOOD IN BOTH NOSTRILS & MOUTH - LUNGS CM - BOLLY - FIRM PULSATING MASS  
 MOVING ALL 4 EXT - SPINAL PRECAUTIONS, → AND 4 ROUTINE BLS/ALS  
 COMMENTS Tx 2 IV CANTOURE. BLOOD GLUCOSE @ 95MG/DL, Tx 3 DON 12000  
 (A) SYNCOPE FACIAL INJURIES, POSS. AAA.  
 MEDICAL REASON FOR AMBULANCE TRANSPORT: O2/SPINAL/LXG/IV/EVAL/Tx → REPORT GIVEN TO KAREN -RN.



Generate All Possible Character Combinations (Remove Duplicates)

Blood bl, bo, ~~bo~~, bd, lo, ~~lo~~, ld, oo, od, ~~od~~  
 Nostrils no, ns, nt, nr, ni, nl, ~~ns~~, os, ot, or, oi, ol, ~~os~~, st, sr, si, sl, ss ...  
 Mouth mo, mu, mt, mh, ou, ot, oh, ut, uh, th

Pairing All High Cohesion Phrases and Construct nm-gram Terms

(Blood, Nostrils) (bl, no), (bl, ns), (bl, nt), (bl, nr) ...  
 (Blood, Mouth) (bl, mo), (bl, mu), (bl, mt), (bl, mh) ...

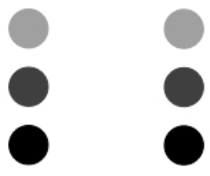
DIGESTIVE-SYSTEM

FQ	CHSN	PHRASE
30	0.72	PAIN INCIDENT
5	0.31	PAIN TRANSPORTED
42	0.54	PAIN CHEST
52	0.81	STOMACH PAIN
9	0.25	HOME PAIN
6	0.43	VOMITING ILLNESS

$$cohesion(w_a, w_b) = z \cdot \frac{f(w_a, w_b)}{\sqrt{f(w_a) * f(w_b)}}$$

# Topic Categorization

	Digestive System	Circulatory System	Immune System	Head	Neck	Back	Pelvis	Legs	Arms	Pseudo-Category Vector
(bl, no)	0	0	0	.95	0	0	0	0	0	0
(bl, ns)	.07	0	0	.34	0	.34	0	0	0	2
(bl, nt)	0	0	0	.95	0	0	0	0	0	1
(bl, nr)	0	.06	0	.44	0	0	0	.20	0	1
(bl, mo)	0	0	0	.33	0	.56	0	0	0	1
(bl, mu)	0	0	0	.41	0	0	0	.51	0	0
(bl, mt)	0	.05	0	.44	0	.19	0	0	0	2
(bl, mh)	0	0	0	.95	0	0	0	0	0	0
(bo, no)	.18	0	0	.57	0	.25	0	0	0	3
(bo, ns)	0	0	0	.42	0	.49	0	0	0	1


 Vector Q is now added to the original normalized/weighted matrix X

$$B_{t,c} = \frac{A_{t,c}}{\sqrt{\sum_{e=1}^n A_{t,e}^2}}$$

$$IDF(t) = \log_2 \frac{n}{c(t)}$$

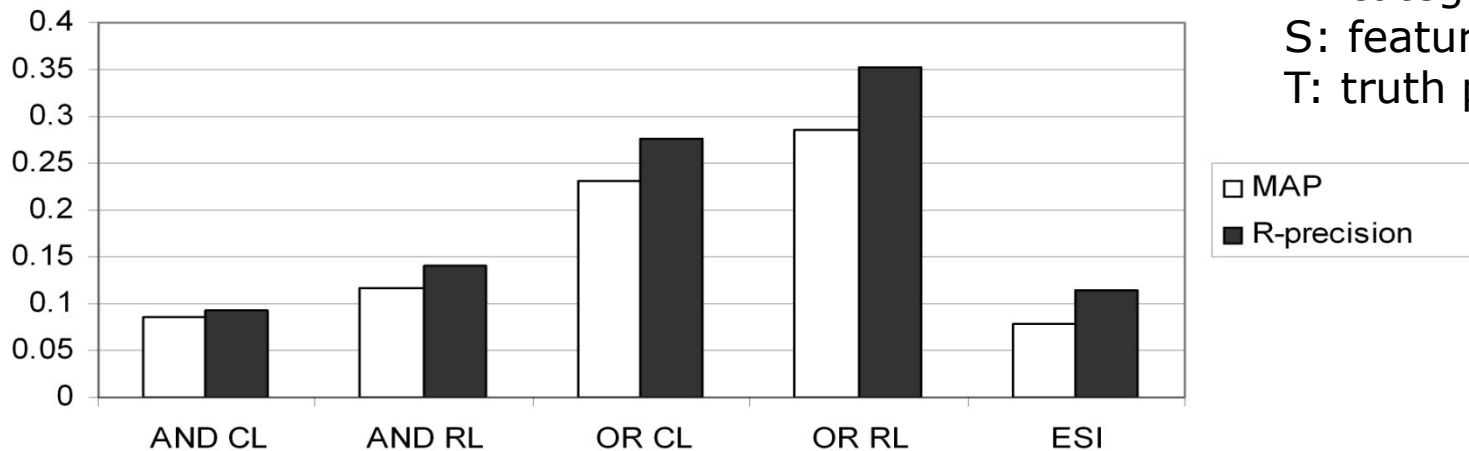
$$X_{t,c} = IDF(t) \cdot B_{t,c}$$

Cosine similarity between trained topic vectors and test document

## Results

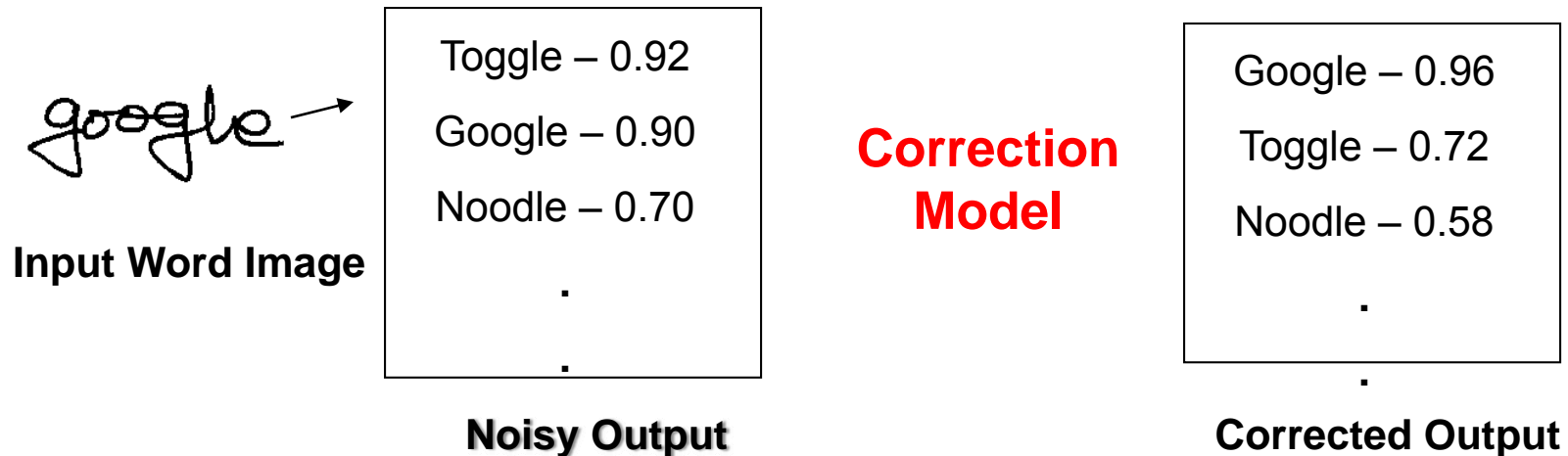
	CLT to RLT	CL to RL	CLT to ALT	CLT to SLT
<b>HR</b>	↑7.48%	↑7.42%	↑17.58%	↑7.42%
<b>Error Rate</b>	↓10.78%	↓10.88%	↓24.53%	↓10.21%

MAP and R-prec comparison chart



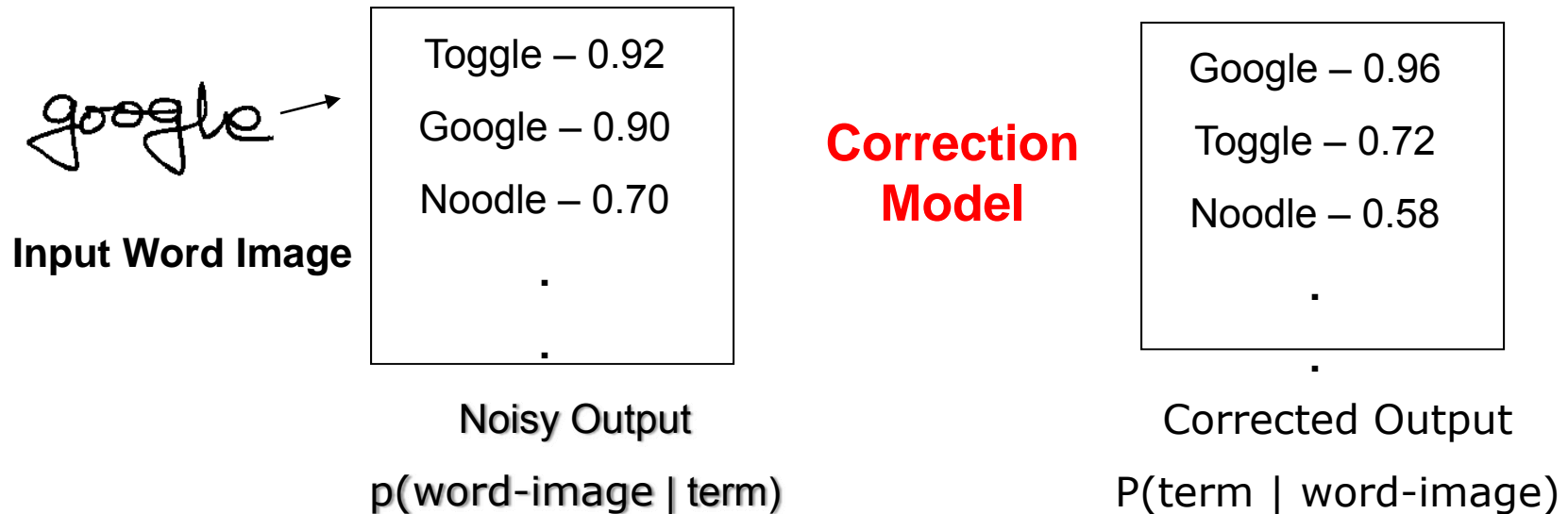
C: complete lexicon  
 R: reduced lexicon  
 A: category given  
 S: features synthetic  
 T: truth present

## Statistical Topic Modeling



- Train topic categorization maximum entropy model
- Generate topic distribution of test document
- Use topic distribution to score each topic as new prior
- Compute posterior probability of word recognition
- Improves from 32% to 40% on IAM dataset

## Statistical Topic Modeling

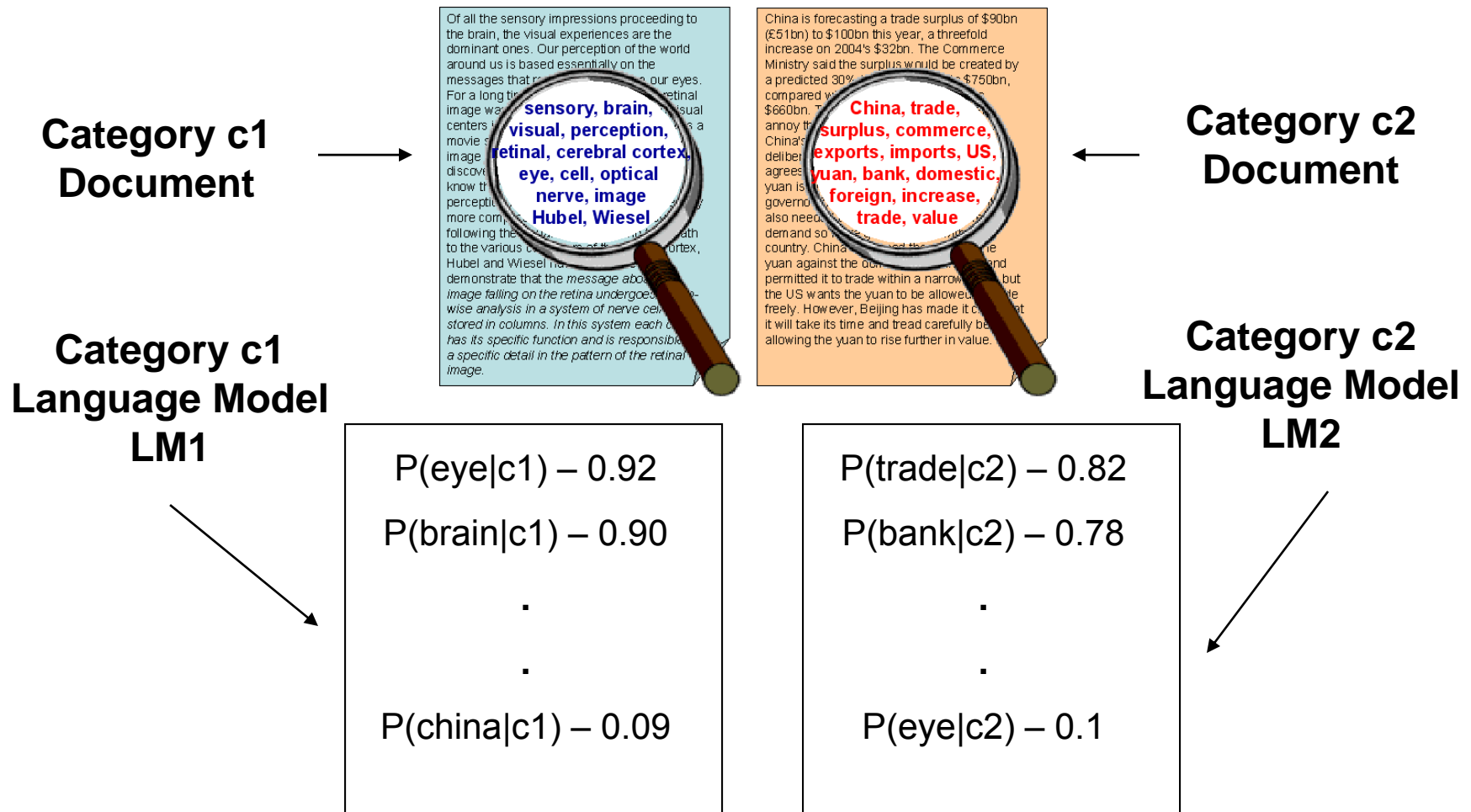


$$P(\text{term} \mid \text{word-image})$$

$$= P(\text{word-image} \mid \text{term}) \times P(\text{term})$$

$$= P(\text{word-image} \mid \text{term}) \times \sum \{ P(\text{term} \mid \text{LM}_i) \times P(\text{LM}_i) \}$$

# Language Model = $P(t | LM_i)$

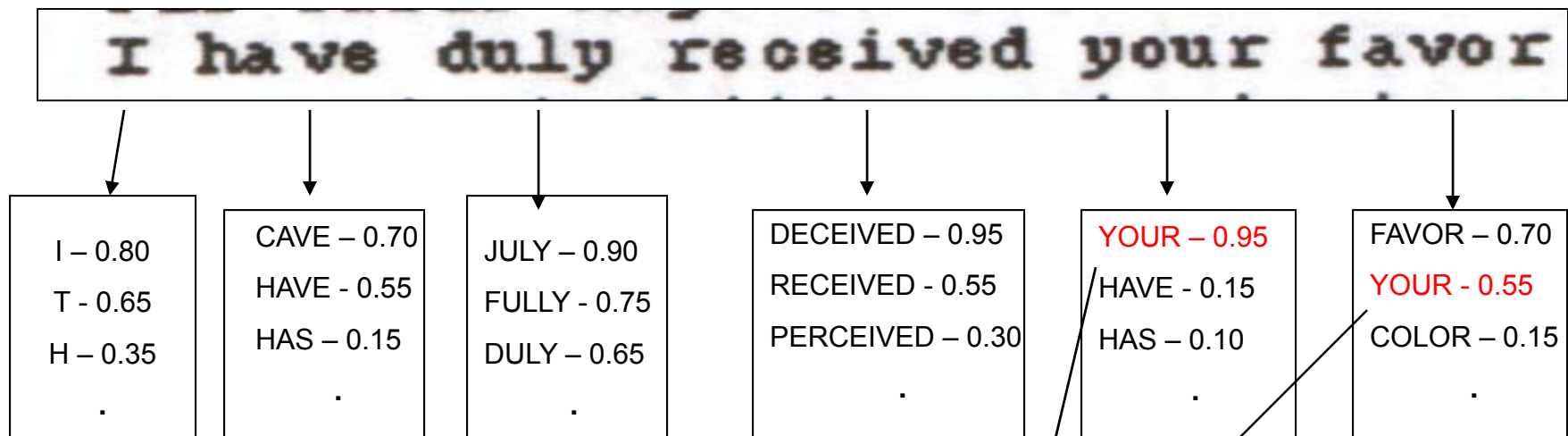


## Topic Distribution = $P(LM_i)$

- Train the Max-Entropy model - fix  $\lambda_i$

$$P(c | d) = \frac{e^{\sum \lambda_i f_i(d,c)}}{\sum_c e^{\sum \lambda_i f_i(d,c)}}$$

$f_i$  is feature (e.g., normalized word counts)

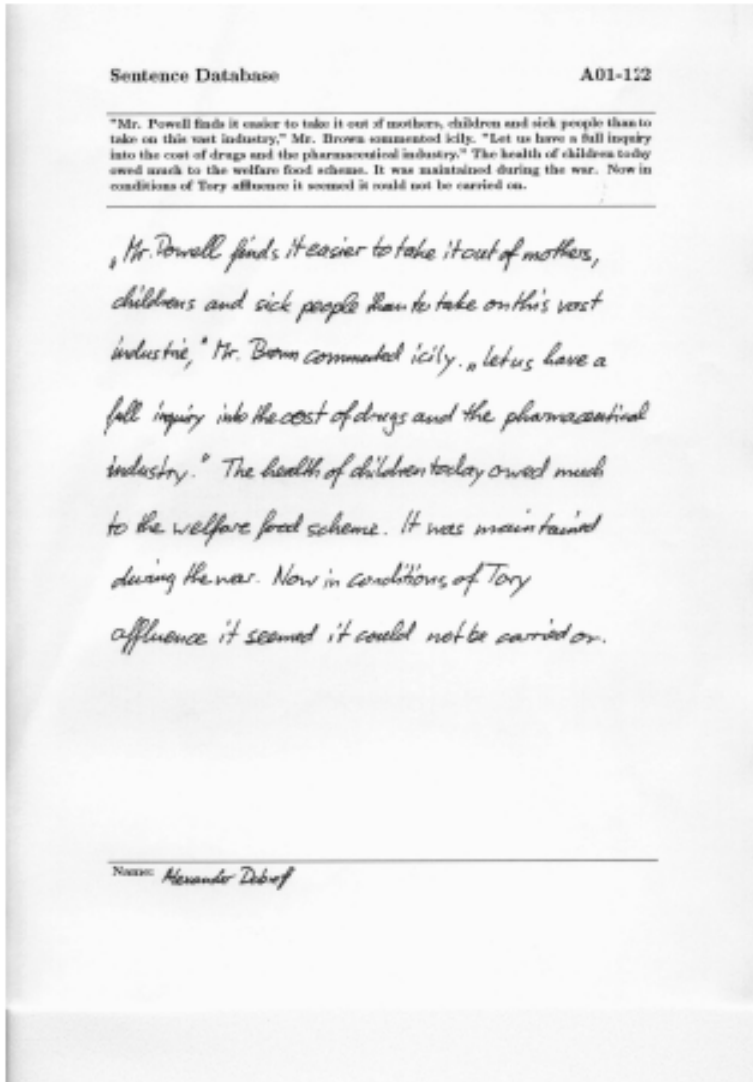


$$\text{count}(\text{YOUR}) = 0.95 + 0.55 = 1.50$$

## Experiments

Method	Word Recognition
Raw	32.33 %
Raw + 3-LM	35.95 %
Raw + Topic-LM	40.63 %

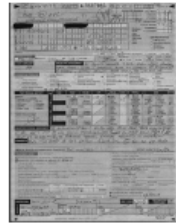
- Corpus : **IAM Database**
- Word Recognizer : **WMR**
- Topic Categorization : **Mallet**
- LM : **CMU-Cambridge LM toolkit**
- # of training docs : **380**
- # of train categories : **13**
- # of test docs : **70**
- # of test word images : **4033**



# Outline

- **Recognition**
  - Postal Application
  - Paradigms
- **Search**
  - OCR accuracy
    - Fusion
    - Lexicon Reduction
    - Statistical Topic Models
  - Document Search
  - Word Spotting

**Indexing**



**Preprocessing:**  
 Binarization  
 Line removal  
 Segmentation

*REUSEG SPINAL IMAGOR*

**Line separation/  
 Word recognition/  
 Probability estimation**

*Handwriting  
 Recognition*

will : 0.1973  
 bail : 0.1600  
 spinal : 0.1475  
 gail : 0.1243  
 ...

**Probabilistic  
 indexing of  
 vector model**

**tf**  $\begin{bmatrix} 0.212, 0.08, \dots, 0.011 \\ \dots \\ 0.005, 0.0023, \dots, 0 \end{bmatrix}$   
**idf**  $[1, 4.75, \dots, 1.58]$

**Retrieval**

**Query phrase**

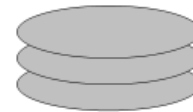
>back pain

**query vector**

$[0 \dots 10 \dots 010 \dots 0]$

**TF-IDF  
 similarity**

**PCR database**



**Keyword  
 spotting  
 similarity**

>chest

**Keyword**

**RANK 1**

FORM ID=918

**RANK=2**

FORM ID=61

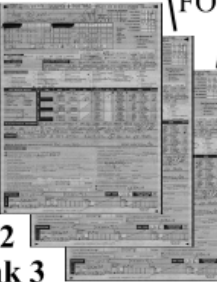
**RANK=3**

FORM ID=302

rank 1

rank 2

rank 3



**RANK 1**

FORM ID=8

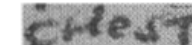
X=980, Y=1918



**RANK=2**

FORM ID=9

X=1063, Y=1029



**RANK=3**

FORM ID=42

X=2000, Y=1902



...

## Vector IR Model (TF-IDF)

Set of terms  $\{t_i\}$ ;

Set of documents  $\{d_j\}$  of length  $\{L_j\}$

- Term Frequency (TF)  $tf_{i,j} = \frac{freq_{i,j}}{L_j}$
- Inverted Document Frequency-IDF

$$idf_i = \log \frac{\#\{d_j\}}{\#\{d_j \mid freq_{i,j} > 0\}}$$

- Query TF  $tf_{i,q} = \begin{cases} 1, & \text{if } t_i \text{ is in query } q \\ 0, & \text{otherwise} \end{cases}$

- Similarity  $\text{sim}(d_j, q) = \sum_i tf_{i,j} \cdot idf_i \cdot tf_{i,q}$

$q = \{ \text{"back"}, \text{"pain"} \}$

terms	$tf_{i,j}$	$idf_i$	$tf_{i,q}$
	.	.	0
	.	.	...
	.	.	0
back	0.024	4.1	1
	.	.	0
	.	.	...
	.	.	0
pain	0.008	2.4	1
	.	.	0
	.	.	$\text{sim}(d_j, q)$
	.	.	0

$\times$   $\times$   $\rightarrow \Sigma$


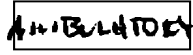
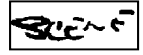
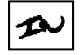
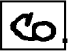
## OCR- Based IR [Prior Work]

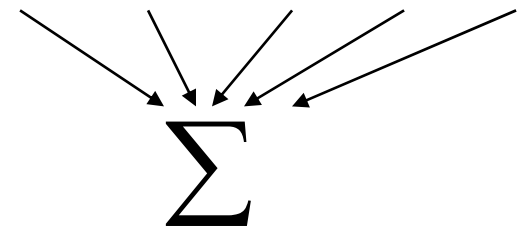
- [Mittendorf, SDAIR96; Ohta, ICDAR97; Jing, CL02]
  - Model common OCR errors
  - Machine print
- [Rath, SIGIR04]
  - Learn word pdf (over image features)
  - Requires large annotated training set
  - **ASSUME**
    - Perfect segmentation
    - Single writer
- [Howe, SIGIR2005]
  - **ASSUME**
  - Ranks obey Zipfian distribution
  - Segmentation assumed perfect

# Estimation

$$E\{freq_{i,j}\} = \sum_{k=1}^L \Pr(t_i | w_k)$$

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

Doc  $d_j$          ..  
 Pr("pain" |  $w_k$ )    0.02    0.01    0.2    0.01    0.01



$$E\{freq_{\text{"pain"},j}\}$$

## Required (Assumed) Inputs

Word segmentation result

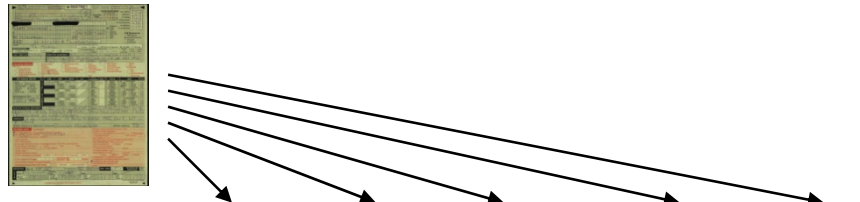
Word recognition likelihoods

## Estimating Term Frequency

$$E\{freq_{i,j}\} = \sum_{I_w} \Pr(I_w) \Pr(t_i | I_w) + \varepsilon$$

$\{t_i\} : \{t_0 = \text{"head"}, t_1 = \text{"arm"}, t_2 = \text{"pelvis"}, \dots\}$

$$\begin{aligned}
 &E\{freq_{1,j}\} \\
 &= \sum_{I_w} \Pr(I_w) \Pr(\text{"arm"} | I_w) \\
 &= 1 \times 0.05 + 1 \times 0.7 + \\
 &0.5 \times 0.01 + 1 \times 0.07 + \dots
 \end{aligned}$$

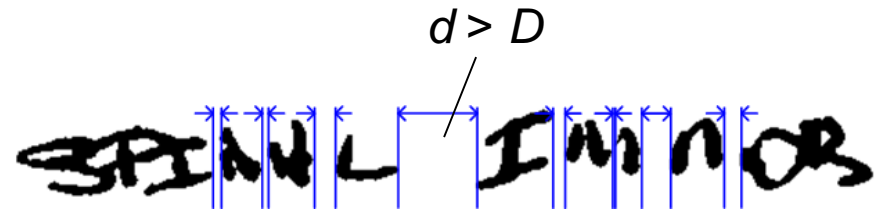


$I_w$	Clear	Arm	pelvis	Stable	...
$\Pr(I_w)$	1	1	0.5	1	...
$\Pr(\text{"head"}   I_w)$	0.2	0.07	0.002	0.01	...
$\Pr(\text{"arm"}   I_w)$	0.05	0.7	0.01	0.07	...
$\Pr(\text{"pelvis"}   I_w)$	0.01	0.01	0.8	0.03	...
...	...	...	...	...	...

# Estimating Segmentation

$$\Pr(I_w)$$

- Word Segmentation
- Gap between adjacent connected components above a threshold  $D$



- Generate multiple hypotheses with multiple  $D$



- If hypothesis  $I_w$  overlaps
- $m$  other hypotheses, then

$$\Pr(I_w) = \frac{1}{m+1}$$

3 hypotheses

$m$	1	2	1
$\Pr(I_w)$	$\frac{1}{2}$	$\frac{1}{3}$	$\frac{1}{2}$

# Word Recognition

## Prior Work

$$\Pr(t_i | I_w)$$

- **Top-Rank (Top-S candidates involved)**

$$\Pr(t_i | I_w) = \begin{cases} \frac{1}{S}, & \text{if } 1 \leq \text{rank}(t_i) \leq S \\ 0, & \text{otherwise} \end{cases}$$

- **Weighted Top-Rank (R = rank( $t_i$ ))**

$$\Pr(t_i | I_w) = \text{top - R OCR rate} - \text{top - (R - 1) OCR rate}$$

- **Empirical**

$$\Pr(t_i | I_w) = \frac{\Pr(t_i) \cdot e^{-\frac{d_i^2}{2\sigma^2}}}{\sum_i \Pr(t_i) \cdot e^{-\frac{d_i^2}{2\sigma^2}}}$$

## Results

- 16303 word images (342 forms)
- Automatic segmentation: 63% correct; 32% under; 5% over
- Lexicon 4405 words
- Query: 22 (1-3 words)

$n =$	1	2	5	10
Top $n$ word recognition rate	20%	27%	35%	42%

Estimate of $\Pr(t_i I_w)$	OCR'ed text (no Estimate)	Top-Rank (S=5)	Top-Rank (S=7)	Weighted Top-Rank	Distance Modelling	Ground Truth (no Estimate)
MAP/R-Precision(%) (human word segmentation)	15.9/17.4	16.8/22.2	19.0/23.4	19.4/25.3	22.1/25.6	70.6/70.7
MAP/R-Precision(%) (not handling seg error)	<b>6.7/7.4</b>	12.8/16.9	14.0/18.8	15.9/20.3	17.6/20.8	N/A
MAP(%) (handling seg error)	7.8/7.7	13.9/18.6	15.5/21.2	16.3/20.1	<b>18.7/22.6</b>	N/A

## MMSE Estimation

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

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

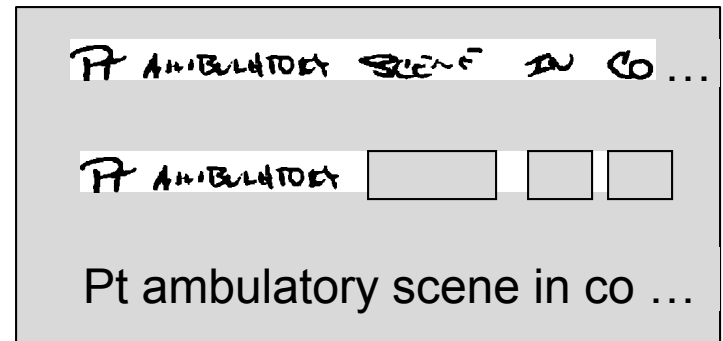
$\vec{w} = [w_1 w_2 \dots w_L]$  : word segmentation hypothesis

$\vec{\tau} = [\tau_1 \tau_2 \dots \tau_L]$  : a decoded term sequence

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

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

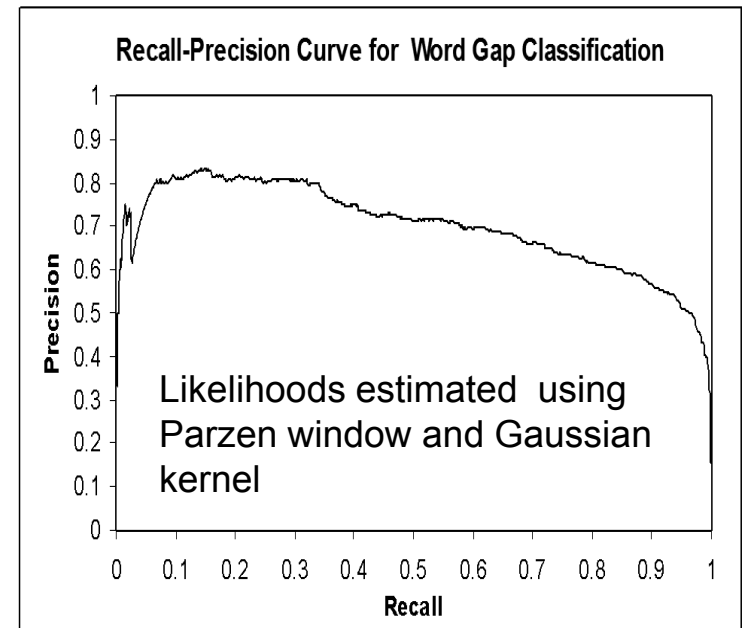
$\#_{t_i}(\vec{\tau})$  : number of  $t_i$  in  $\vec{\tau}$   $\#_{\text{"scene"}}(\text{"Pt ambulatory scene in co..."} ) = 1$



## Word Gap

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

- $fv_1$ : Euclidian dist between bounding boxes
- $fv_2$ : Shortest white run between two CC's
- $fv_3$ : Distance between convex hulls

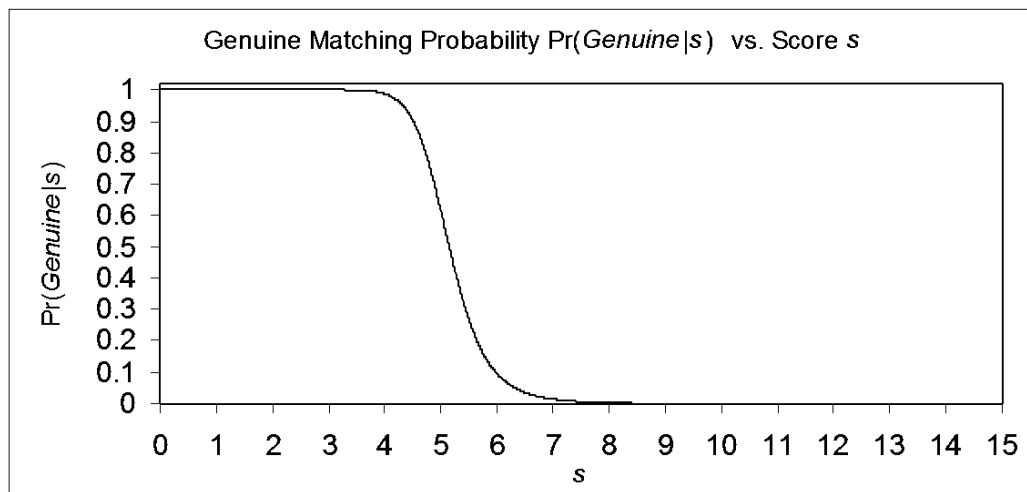


## Word Recognition Likelihood

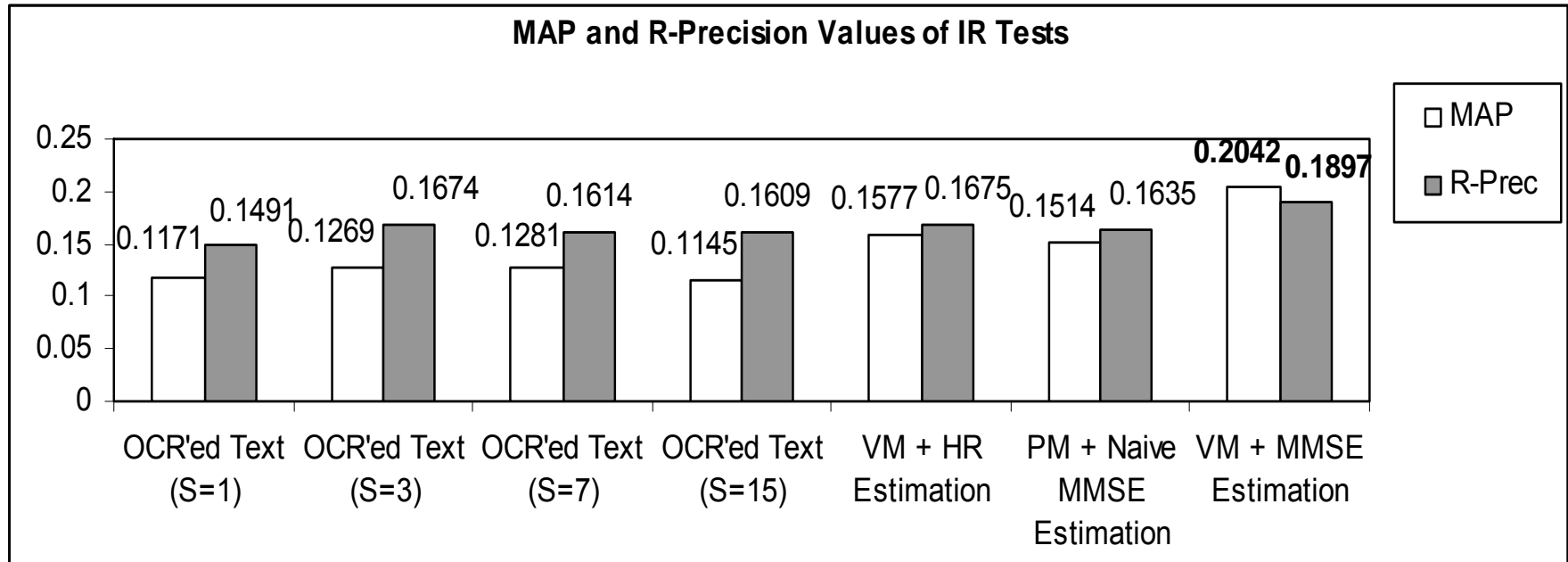
$s(w, t_i)$

$UBM(s) = \Pr(\text{Genuine} | s) =$

$\frac{\Pr(\text{Genuine}) \Pr(s | \text{Genuine})}{\Pr(\text{Genuine}) \Pr(s | \text{Genuine}) + \Pr(\text{Imposter}) \Pr(s | \text{Imposter})}$



$$p(w | t_i) \propto UBM(s(w, t_i))$$



S: number of top-candidates retained in the OCR'ed text

VM: Vector Model

PM: a Probabilistic IR Model

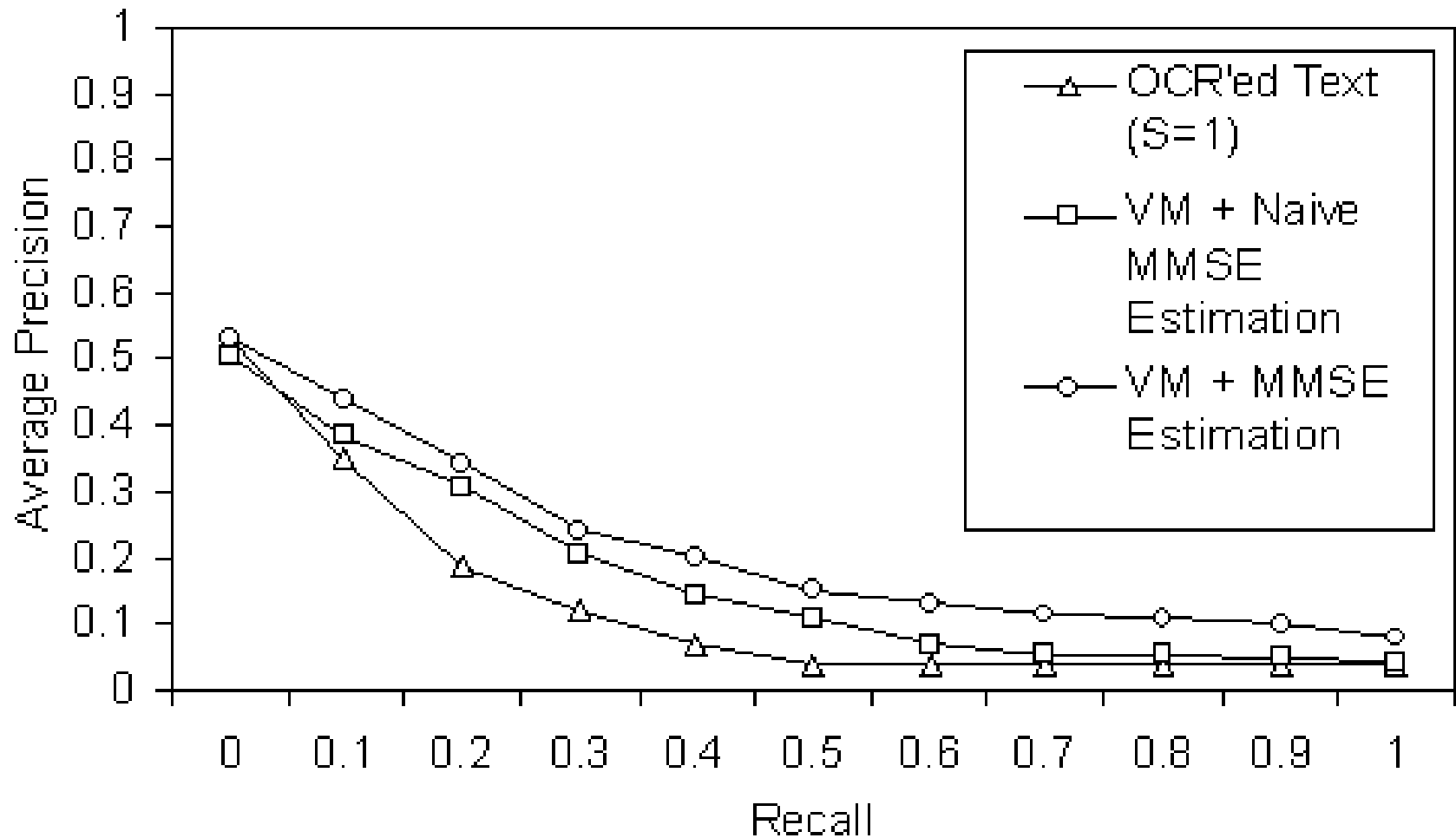
Naive MMSE Estimation:

MMSE Estimation:

$$E\{freq_{i,j}\} = \sum_{k=1}^L \Pr(\tau_k | w_k)$$

$$E\{freq_{i,j}\} = \sum_{\vec{w}} \Pr(\vec{w} | \vec{\sigma}) \cdot \sum_{\vec{\tau}} \Pr(\vec{\tau} | \vec{w}) \cdot \#_{t_i}(\vec{\tau})$$

### 11-Point Average Precision



# Word Spotting

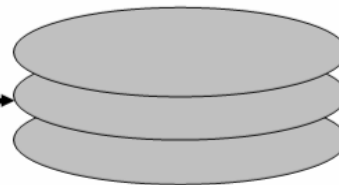
## Document Images



**Preprocessing:**  
Noise Removal  
Page Segmentation  
Word Segmentation

**Feature Extraction  
from Word Images**

**Storing into  
Database**



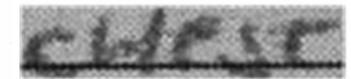
## User query

>chest

**Similarity**

**Output**

RANK 1  
FORM ID=8  
X=980,Y=1918



RANK=2  
FORM ID=9  
X=1063,Y=1029

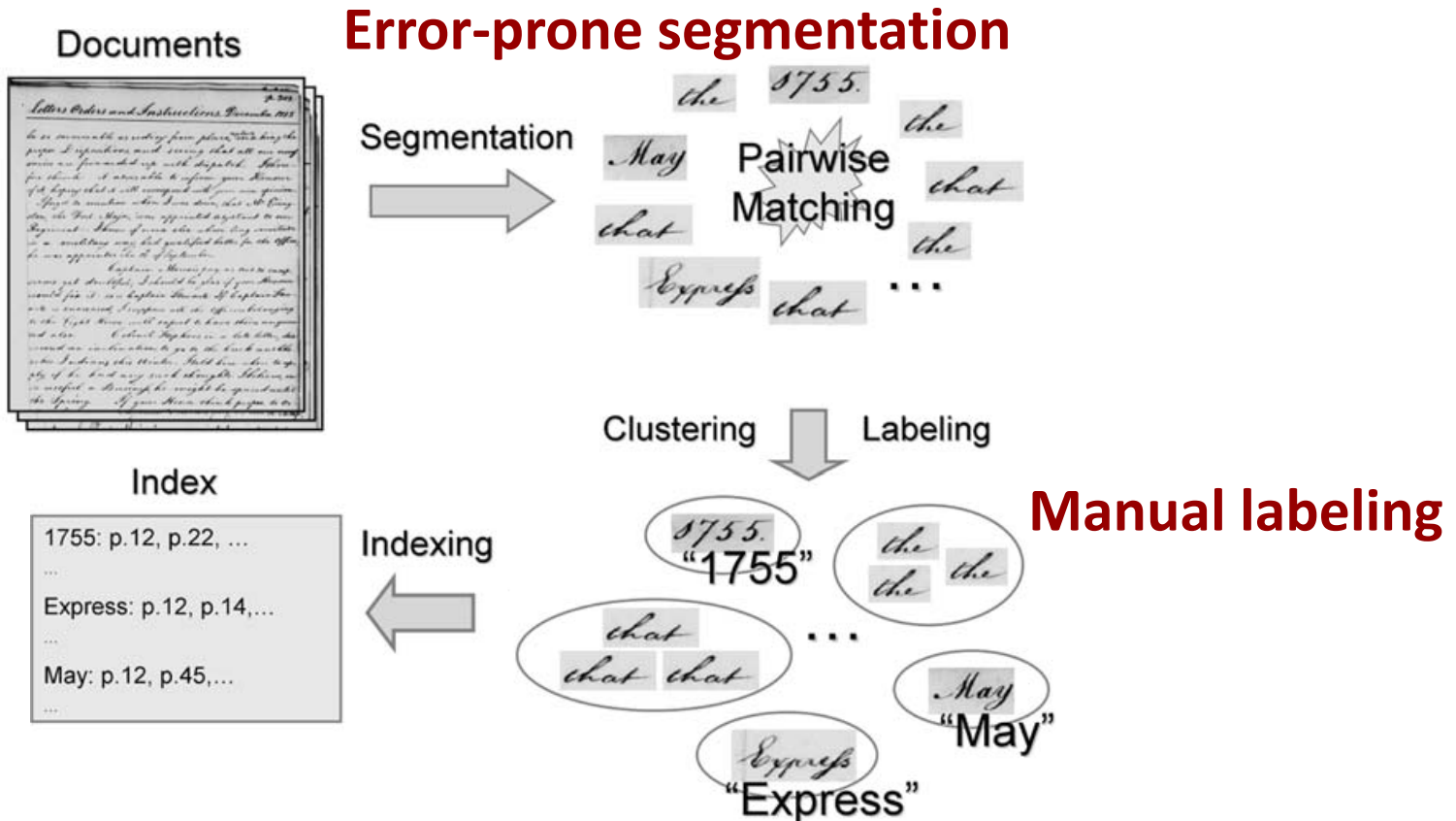


RANK=3  
FORM ID=42  
X=2000,Y=1902



⋮

# Image Based Methods



**Poor performance in multiple writer scenarios**

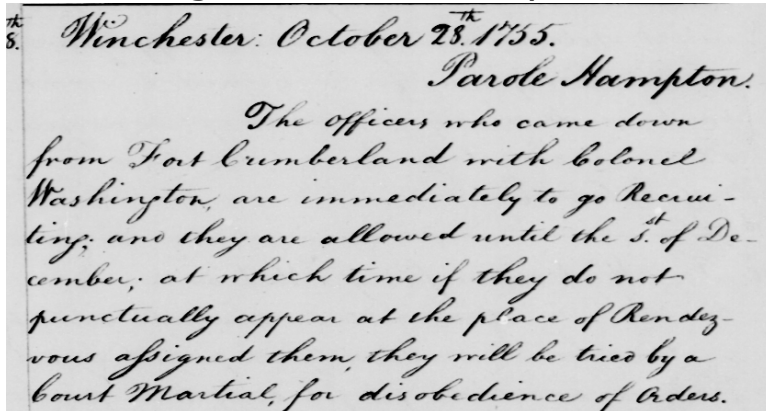


# Keyword Spotting

## Previous Work

- **Corpus**

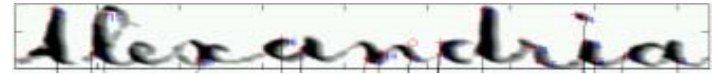
Washington's manuscripts



- **MAP Performance**

40.98% (2372 good quality images)  
16.5% (3262 poor quality images)

**Query: Both Image and Text**



Script specific Upper/ lower profile structural features



Observation density

$$\Pr(wrd, fv)$$

Posterior word recognition probability

$$\Pr(wrd | fv) = \frac{\Pr(wrd, fv)}{\sum_{wrd} \Pr(wrd, fv)}$$

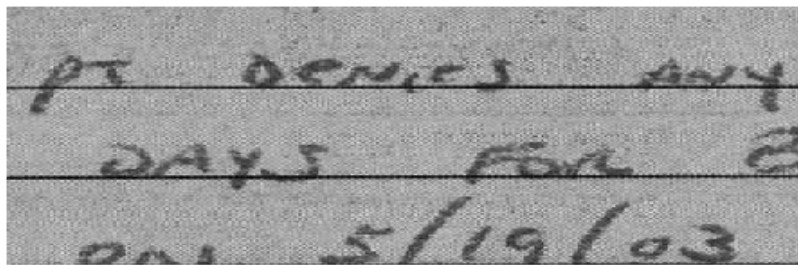


# Keyword Spotting

## Previous Work

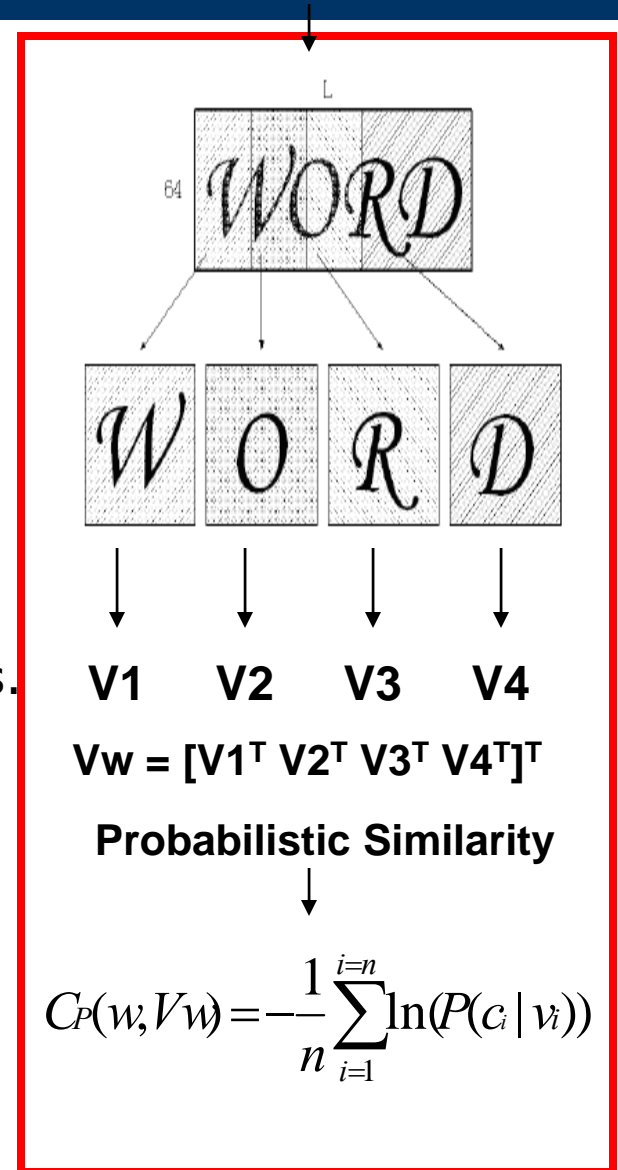
### Template Free Word Spotting

- Matching Gabor features of two word images.
- Posterior probability estimate from SVM OCR
- **Corpus:**  
12 medical forms with 5295 character images.  
101 samples of 6 keywords

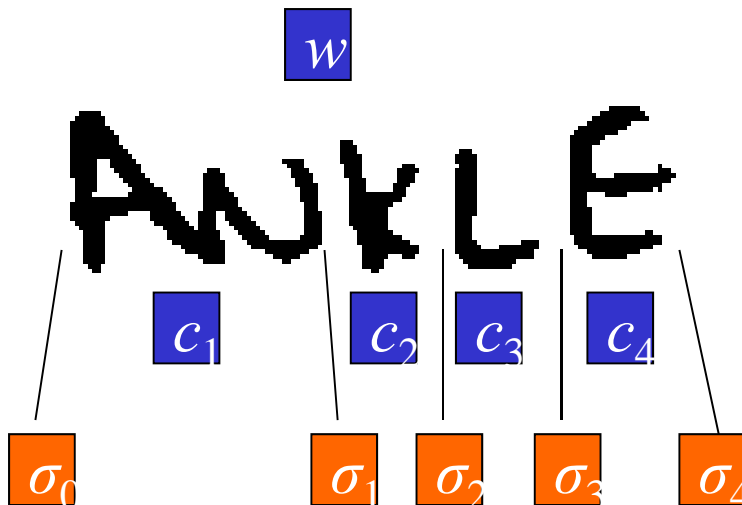


### MAP Performance

67.1% compared to 12.6% by DTW



# Probabilistic Indexing

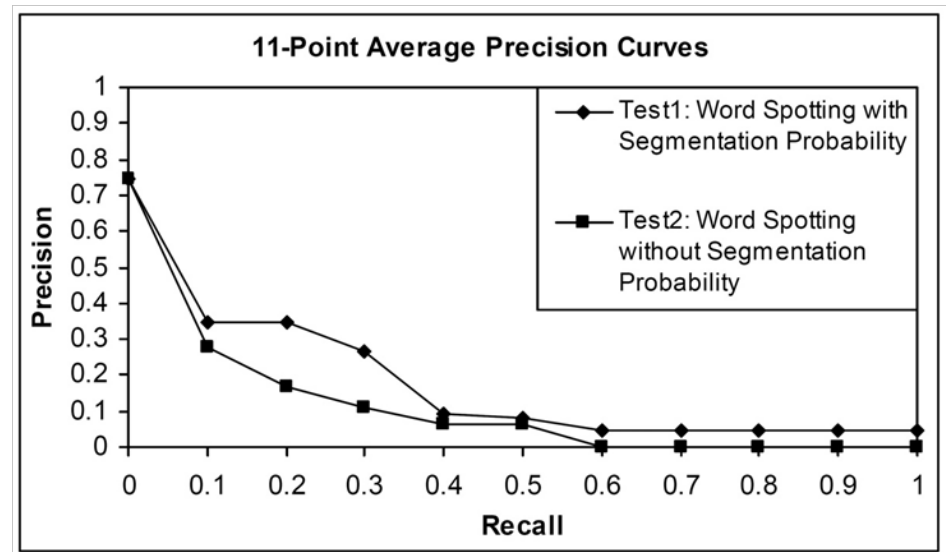


$$w \rightarrow c_i, c_{i+1}, \dots, c_j$$

Word gap probability ( $\sigma_k$ ) between  $c_k$  and  $c_{k+1}$  components

Word - query similarity :  $\text{sim}(w, q) =$

$$\sigma_{i-1} \times (1 - \sigma_i) \times \dots \times (1 - \sigma_{j-1}) \cdot \sigma_j \times \Pr(q | w)$$



## Summary

- Handwriting Recognition remains a challenging task despite success in postal applications
- Need for improved search technologies to access handwritten documents on the web
- Statistical topic models can help in document categorization and lexicon reduction
- Document indexing can be performed by MMSE modeling that integrates segmentation, language models, and recognition
- Word spotting can be performed by indexing image level features and on OCR results

Thank You