




AUTOMATIC RECOGNITION OF HANDWRITTEN MEDICAL FORMS FOR SEARCH ENGINES

Robert Milewski
Venu Govindaraju

TOC


- 
- ◆ Introduction
 - ◆ Prior Work
 - ◆ Binarization
 - ◆ Topic Categorization: Training
 - ◆ Topic Categorization: Testing
 - ◆ Results
 - ◆ Applications



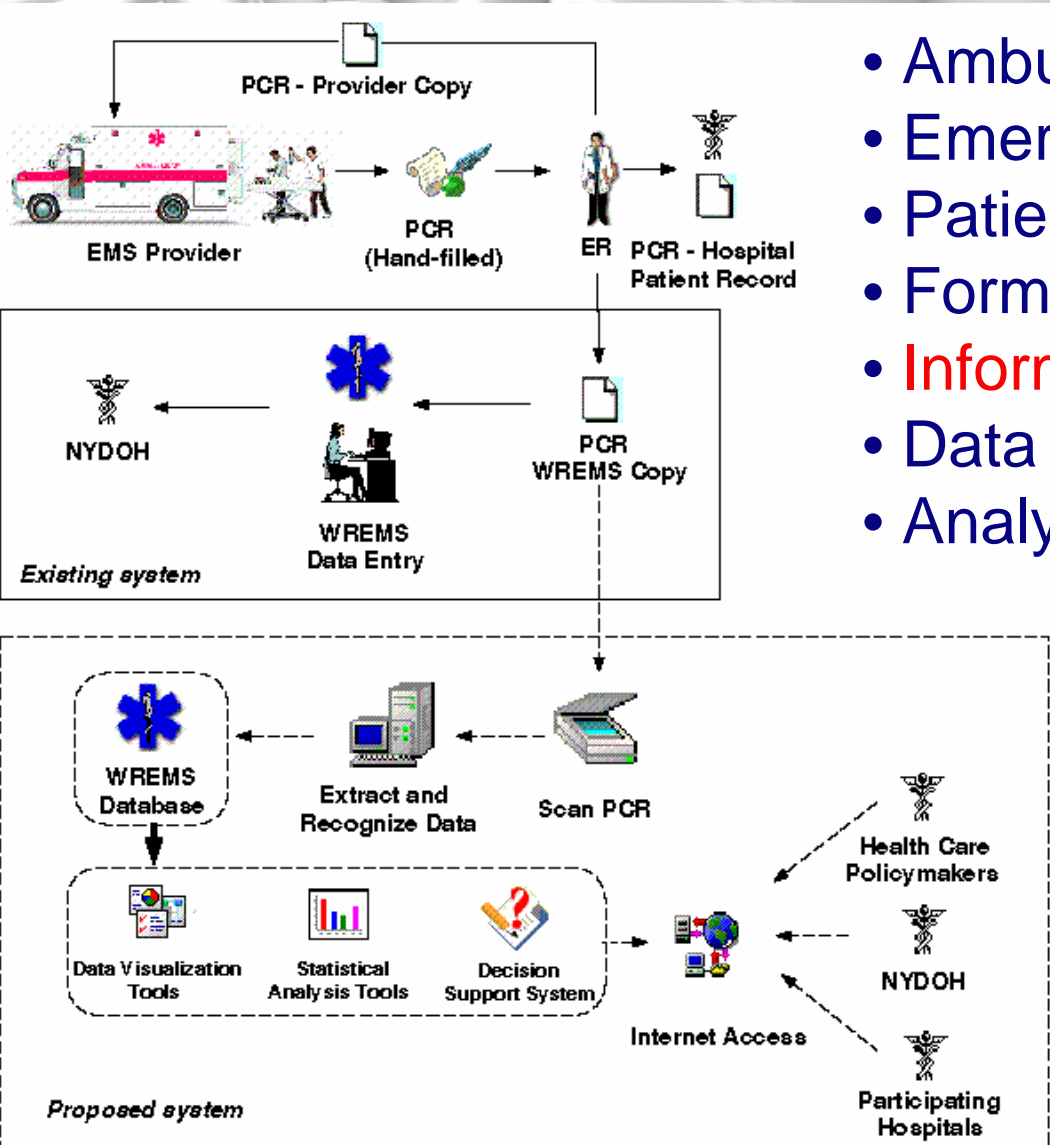
Problem Statement

Address unconstrained recognition and retrieval of handwritten forms by focusing on the NYS PCR medical form.

Objective


- 
- ◆ We will explore:
 - ◆ linguistic models
 - ◆ topic categorization
 - ◆ lexicon reduction

Envisioned Journey of the PCR



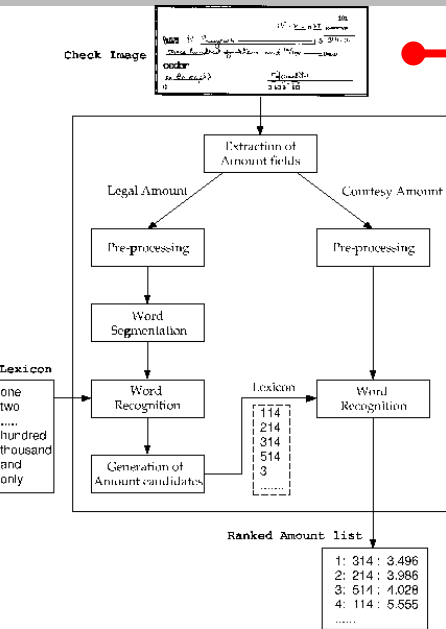
- Ambulance rescues patient
- Emergency environment complex
- Patient information documented
- Form passed to emergency room
- **Information extraction**
- Data disseminated
- Analysis performed

TOC

- 
- ◆ Introduction
 - ◆ **Prior Work**
 - ◆ Binarization
 - ◆ Topic Categorization: Training
 - ◆ Topic Categorization: Testing
 - ◆ Results
 - ◆ Applications

(Kim and Govindaraju, 1996)
 (Madhavanath et al 1996)
 (Kim and Govindaraju, 1997)
 (Tomai, Zhang and Govindaraju, 2002)

Evolution of Form Analysis



Bank Check

Describe the activity at location where employed. /

Scheels

(For example: hospital, newspaper publishing, mail order house, auto engine manufacturing, retail bakery)

c. Is this mainly — Fill ONE circle

Manufacturing Other (agriculture, construction, service, government, etc.)

Wholesale trade

Retail trade

19. Occupation

a. What kind of work was this person doing? →

Custodian

(For example: registered nurse, personnel manager, supervisor of order department, gasoline engine assembler, cake loer)

b. What were this person's most important activities or duties? →

Custodial

(For example: patient care, directing hiring policies, supervising order clerks, assembling engines, icing cakes)

Census Form

UNITED STATES POSTAL SERVICE

ENGINEERING
 8403 Lee Highway
 Manassas VA 20108-8101

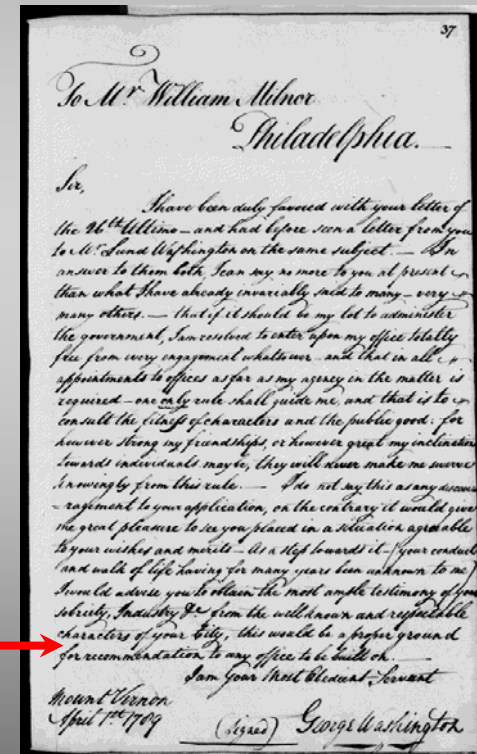
Official Business

First-Class Mail
 Postage & Fees Paid
 USPS
 Permit No. G-10

Mr. Edward J. Kuebert
 7438 Cedar Run Drive
 Warrenton, VA 20187


20187+2246

Postal Mail



Historic Documents

Survey of Lexicon Reduction

- 
- ◆ Bank Applications
 - ◆ Dollar amount constrains lexicon
 - ◆ Postal Applications
 - ◆ City/State/Street databases constrains lexicon
 - ◆ Word Length / Word Shape
 - ◆ Affected by rescue environment
 - ◆ Performance only Reduction
 - ◆ Some models significantly improve speed but drop recognition by as much as 30%

Holistic Lexicon Reduction

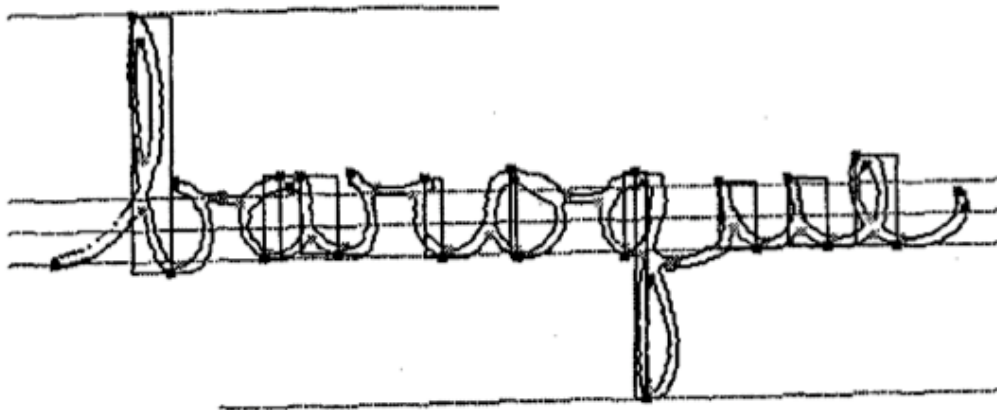


Figure 1: Extracting downstrokes from an offline image: The figure shows an offline image with local maxima and minima highlighted on the upper and lower contours of the image respectively. Extracted downstrokes are indicated by bounding rectangles.

- ◆ Used for lexicon reduction
- ◆ Features
 - ◆ ascenders
 - ◆ descenders
 - ◆ word shape
- ◆ Medical environment complicates these features



Call-Routing Problem

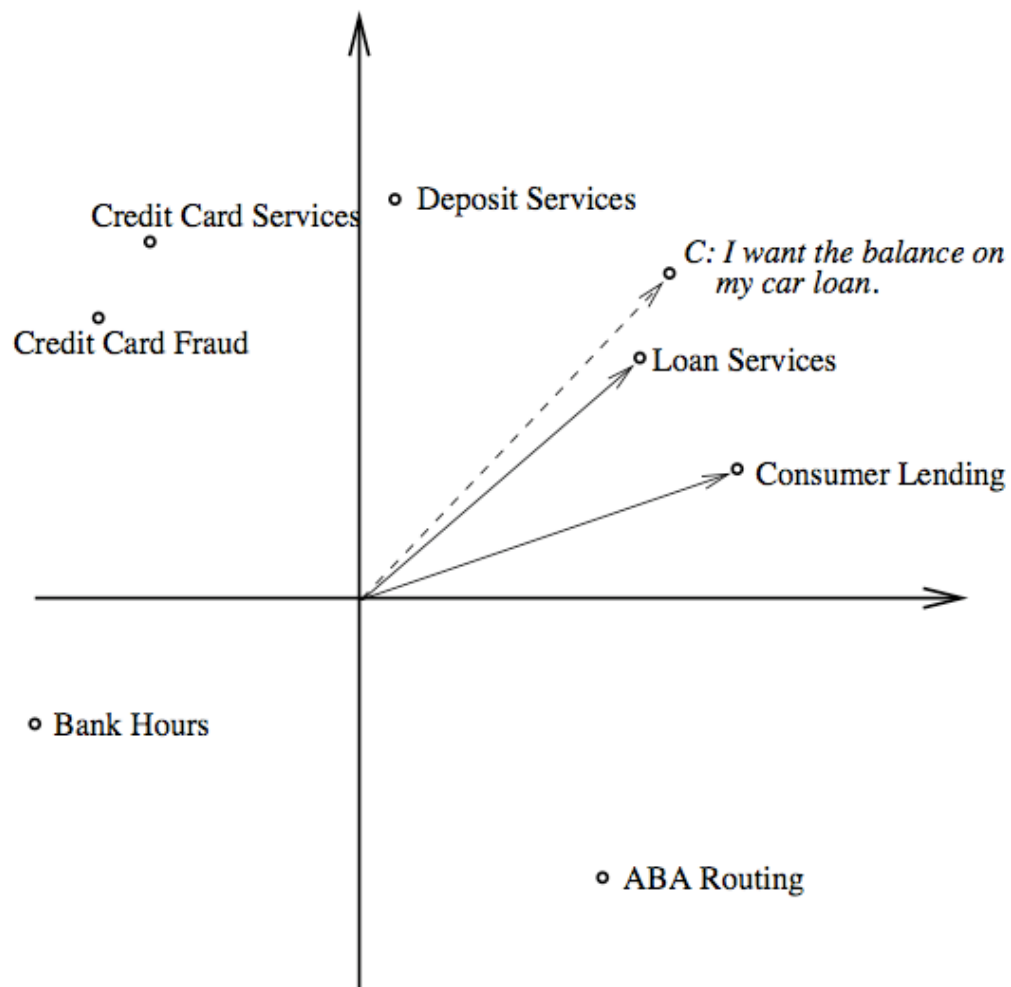


Figure 2
2-Dimensional Vector Representation for the Routing Module

Used R-SVD approach.

Call-Routing:

Input: Voice Recognition

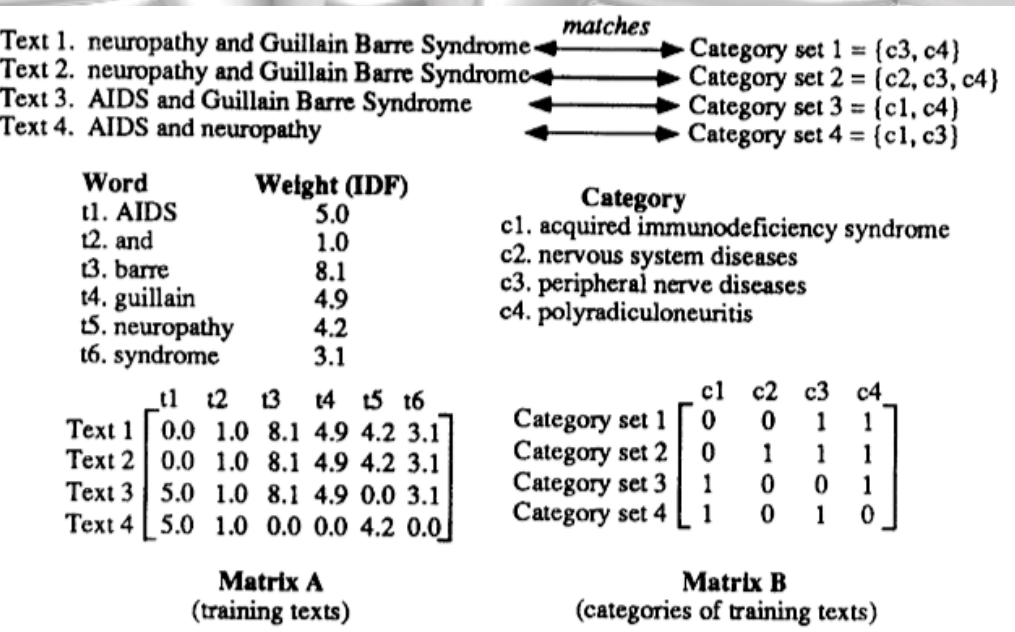
Output: Destination

Medical Forms:

Input: Handwriting Recognition

Output: Topic Category

Medical Text Categorization



- ◆ Used SVD approach
- ◆ Applied on known text
- ◆ Learned word-category associations on medical text
- ◆ Used MEDLINE and MAYO

Fig. 1. The matrix representation of matched text/category-set pairs.

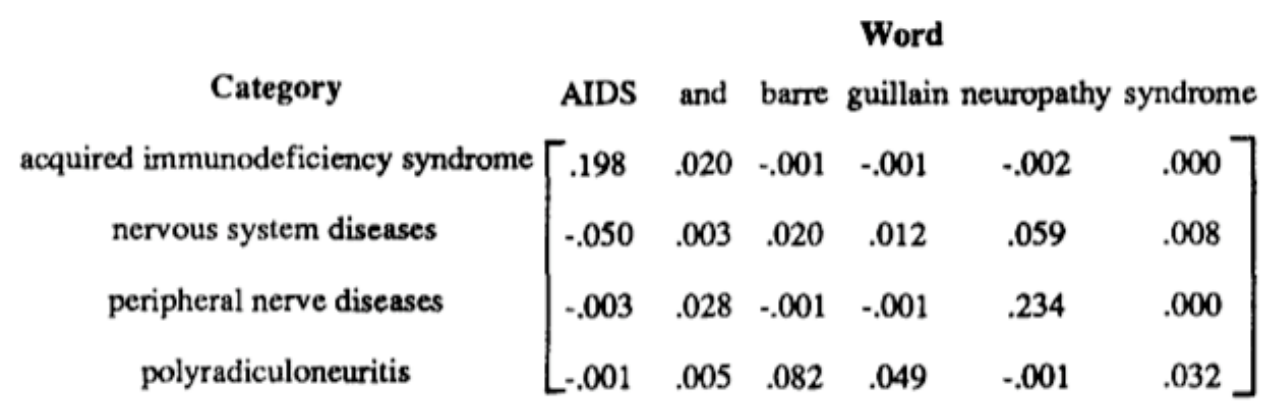



Fig. 2. A word-category association matrix obtained by solving the LLSF of the training set from Figure 1.

TOC

- 
- ◆ Introduction
 - ◆ Prior Work
 - ◆ **Binarization**
 - ◆ Topic Categorization: Training
 - ◆ Topic Categorization: Testing
 - ◆ Results
 - ◆ Applications

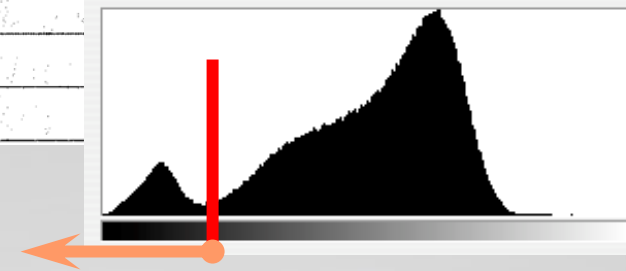
Histogram Classification

OBJECTIVE PHYSICAL ASSESSMENT 75 yo male PAO x 3, HEENT pt wears glasses neck. ⊖TD, ⊖JVD
Chest = expansion, clear sounds abd gut, stable pelvis, PMS x 4 ↑↓ extremities
② arm head trap, prior IV site, 1° vitals 1° assessment enroute, tx to Buffalo Mercy
COMMENTS without incident. PT walked from table to stretcher -

Original Image

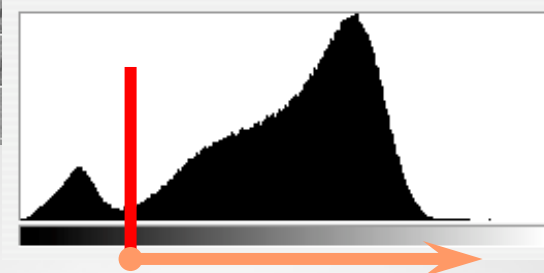
OBJECTIVE PHYSICAL ASSESSMENT 75 yo male PAO x 3, HEENT pt wears glasses neck. ⊖TD, ⊖JVD
Chest = expansion, clear sounds abd gut, stable pelvis, PMS x 4 ↑↓ extremities
② arm head trap, prior IV site, 1° vitals 1° assessment enroute, tx to Buffalo Mercy
COMMENTS without incident. PT walked from table to stretcher -

Thresholded to keep only form



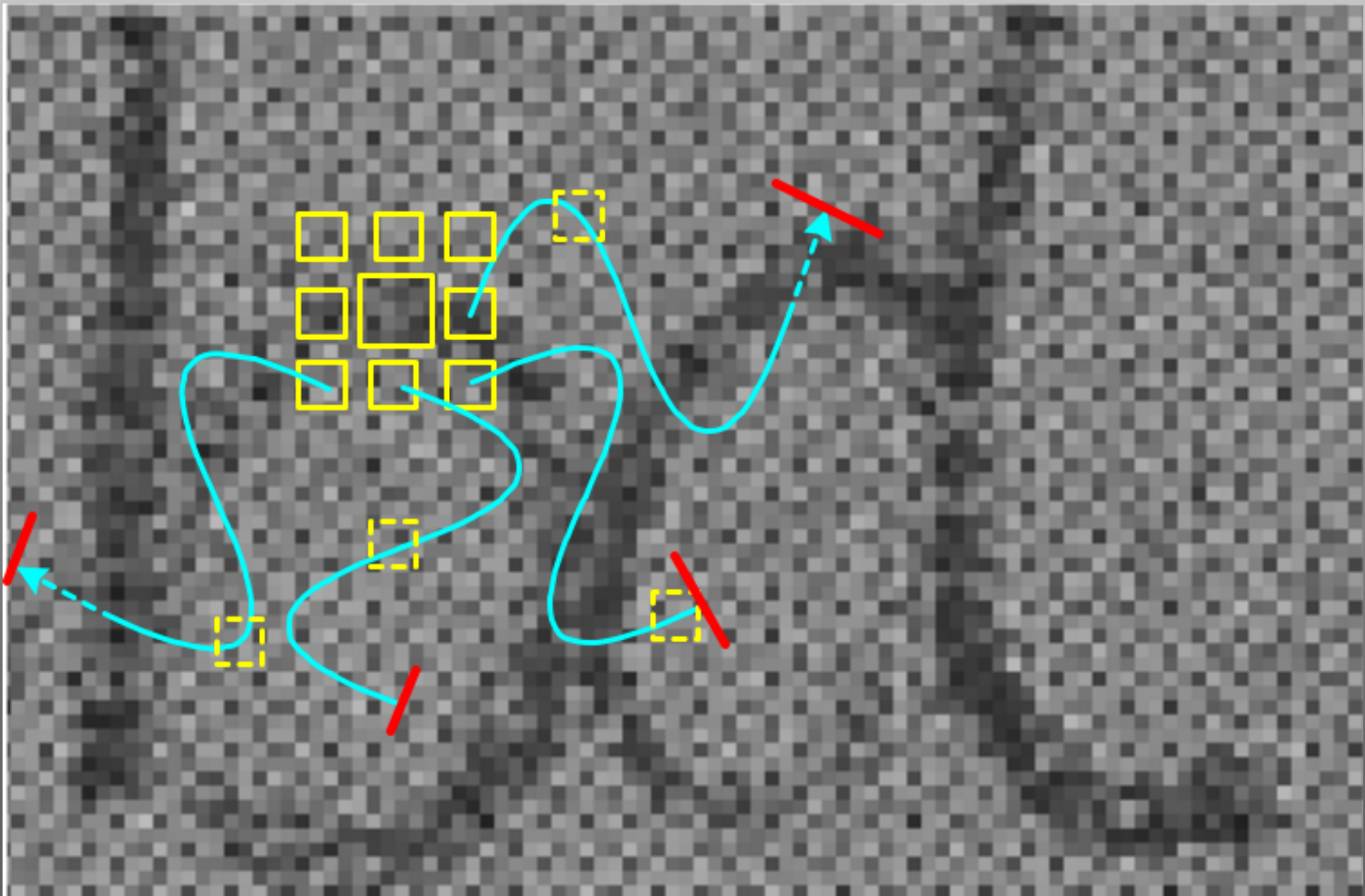
OBJECTIVE PHYSICAL ASSESSMENT 75 yo male PAO x 3, HEENT pt wears glasses neck. ⊖TD, ⊖JVD
Chest = expansion, clear sounds abd gut, stable pelvis, PMS x 4 ↑↓ extremities
② arm head trap, prior IV site, 1° vitals 1° assessment enroute, tx to Buffalo Mercy
COMMENTS without incident. PT walked from table to stretcher -

Thresholded to keep only carbon

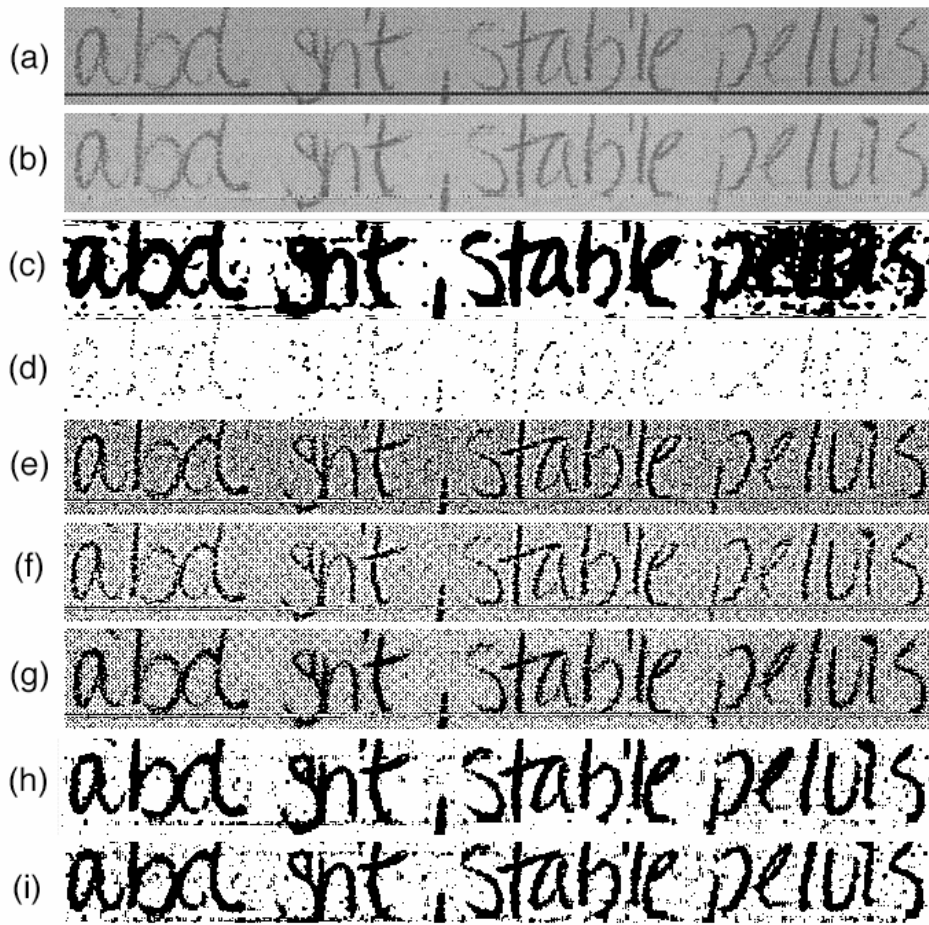


Sinusoidal Coverage

- ◆ Only 4 of the 8 paths are shown.
- ◆ The algorithm is computed on a larger surface; the surface of an entire form block.
- ◆ Desire to stay close to the origin since paper fluctuations occur occur further out.

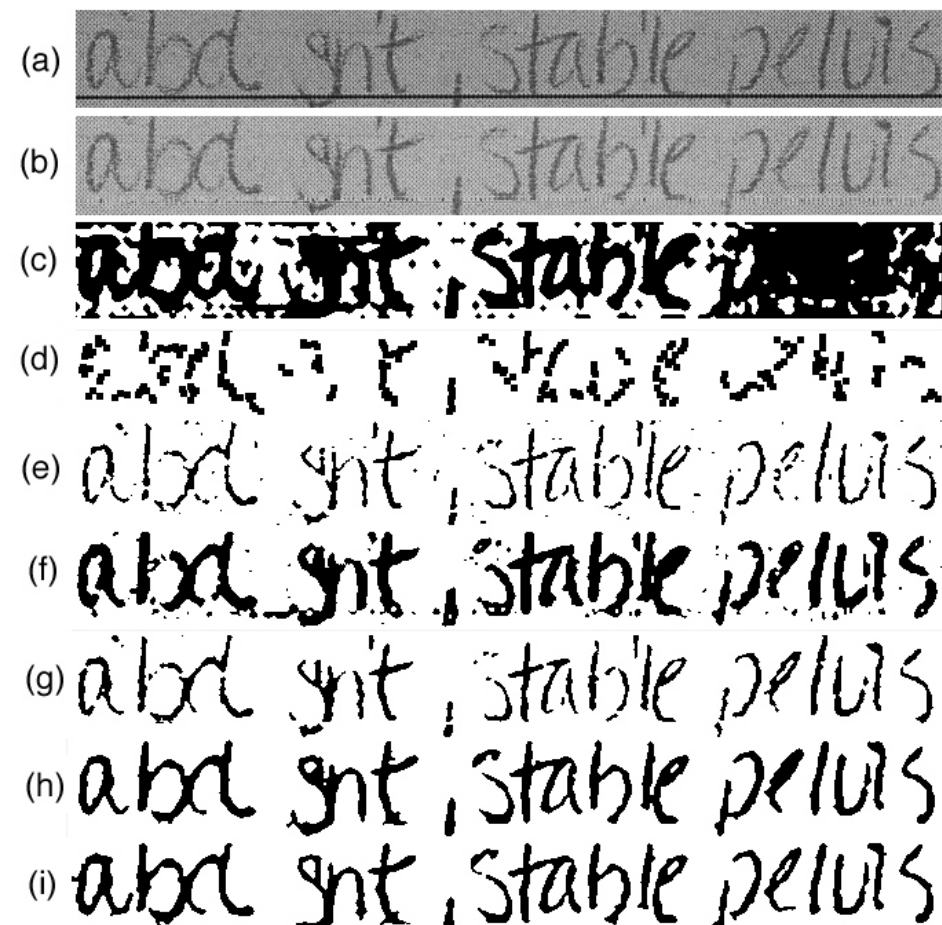


Visual Comparison



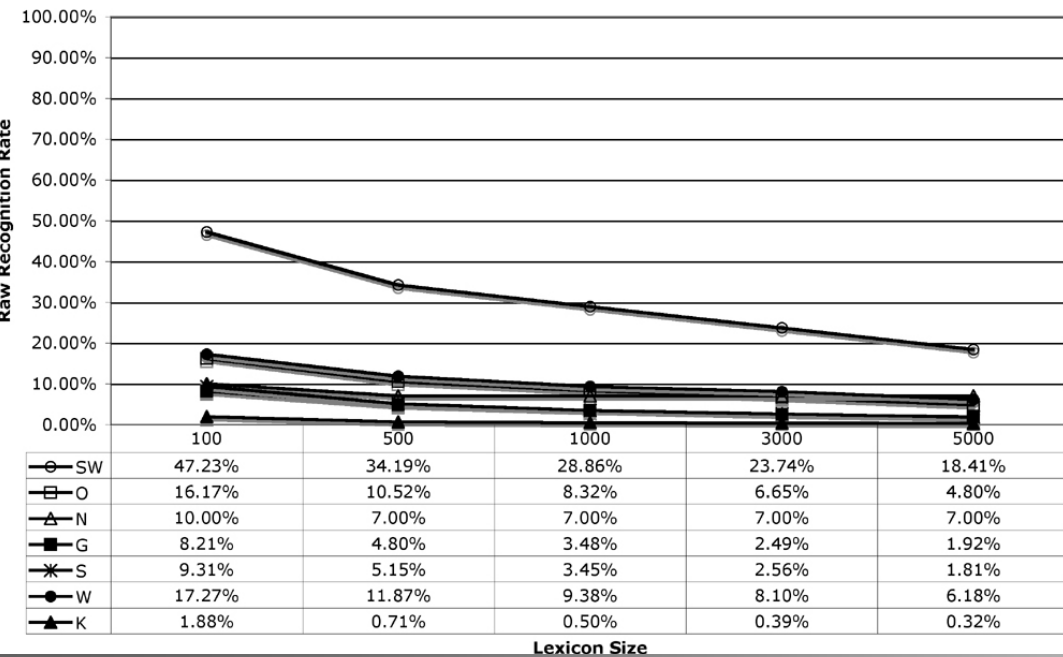
Binarization Only

- (a) Original image (b) Original image with form drop out (c) Wu/Manmatha Binarization (d) Kamel/Zhao Binarization (e) Niblack Binarization (f) Sauvola Binarization (g) Otsu Binarization (h) Gatos (i) Sine Wave Binarization

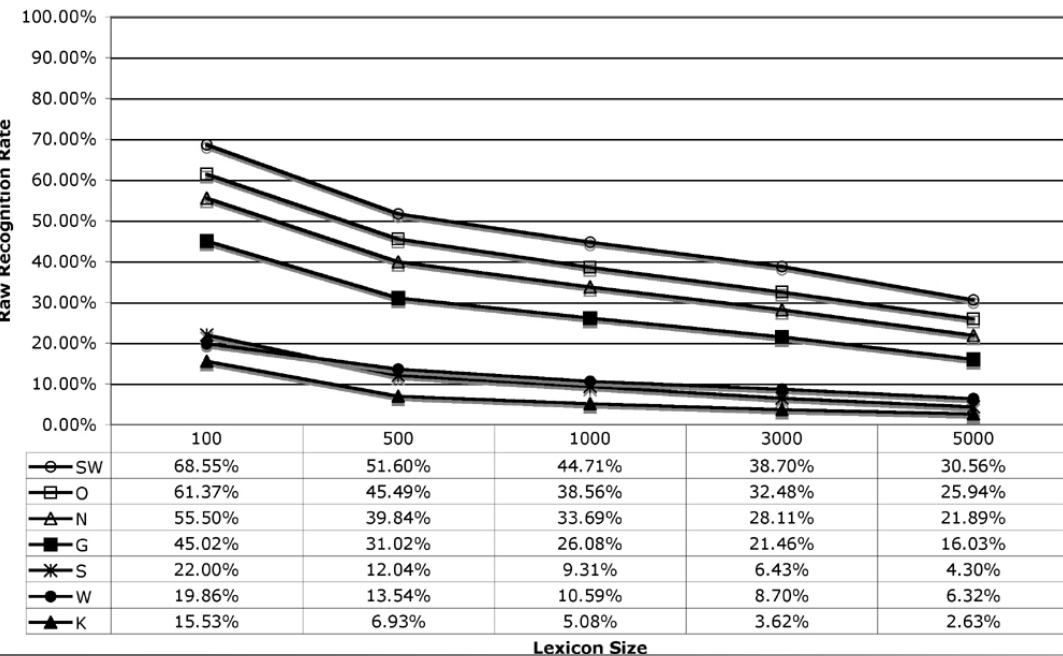


Binarization + Processing

Binarization Only



Binarization and Post-Processing



Performance

Legend

- (SW) Sine Wave Binarization
- (O) Otsu Binarization
- (N) Niblack Binarization
- (G) Gatos Binarization
- (S) Sauvola Binarization
- (W) Wu/Manmatha Binarization
- (K) Kamel/Zhao Binarization


Improvement

Binarization: + 11-31%
 Binarization/PP: + 4.5-7.25%

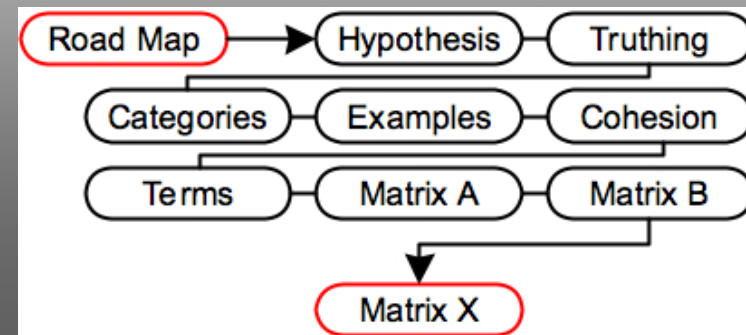
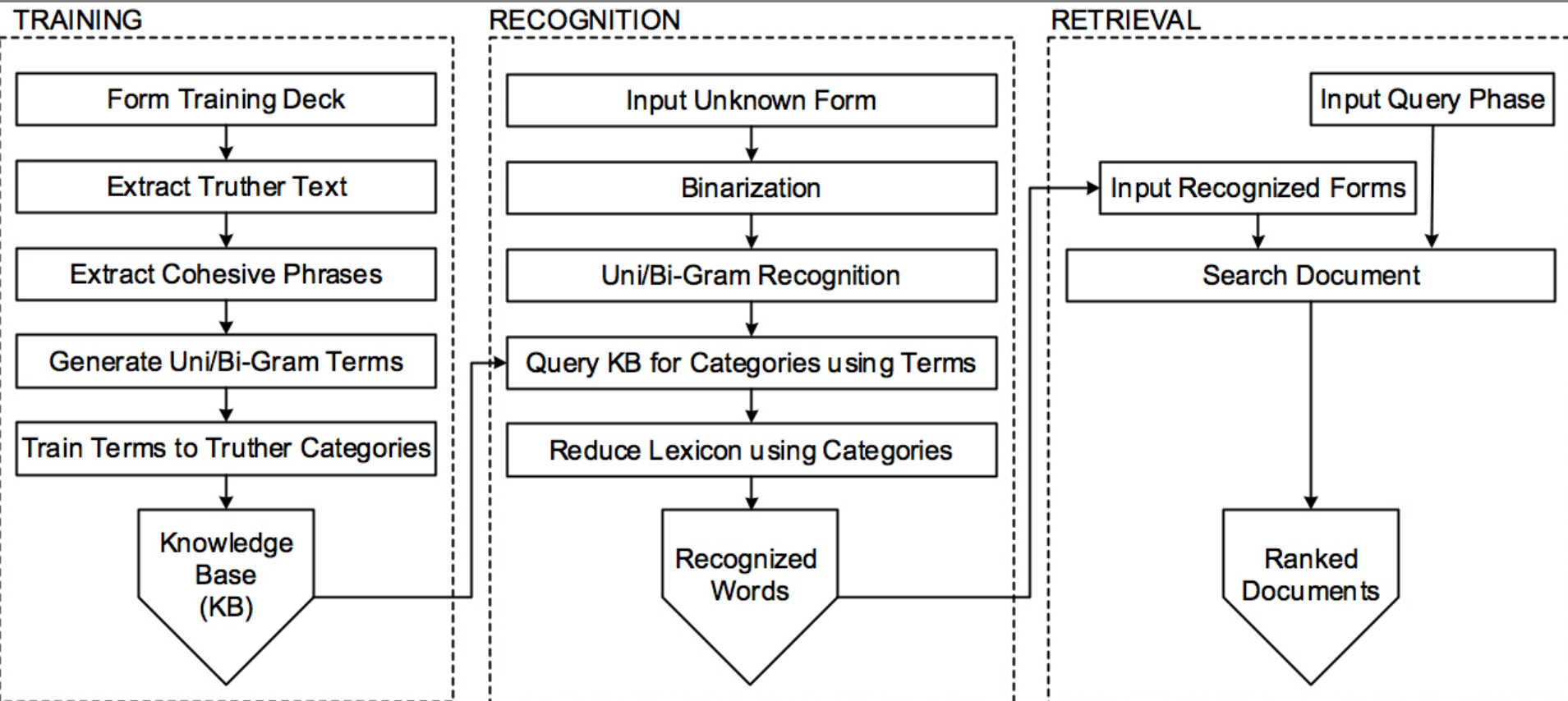
Environment

PCR's 62
 Words 3,000

TOC

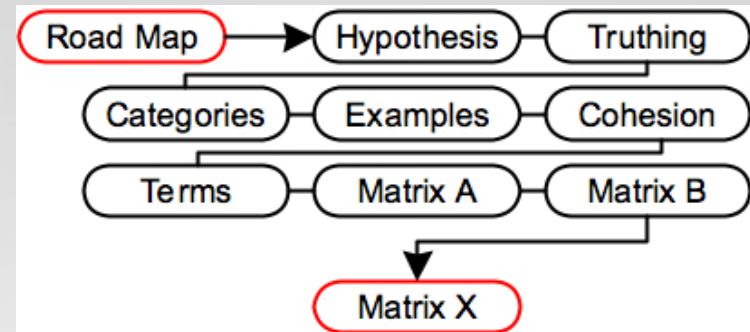
- 
- ◆ Introduction
 - ◆ Prior Work
 - ◆ Binarization
 - ◆ **Topic Categorization: Training**
 - ◆ Topic Categorization: Testing
 - ◆ Results
 - ◆ Applications

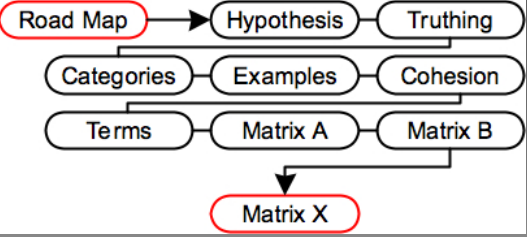
Road Map



Hypothesis

A sequence of confidently recognized characters, extracted from an image of handwritten medical text, can be used to represent a topic category.





PCR Truthing

10 Body Systems

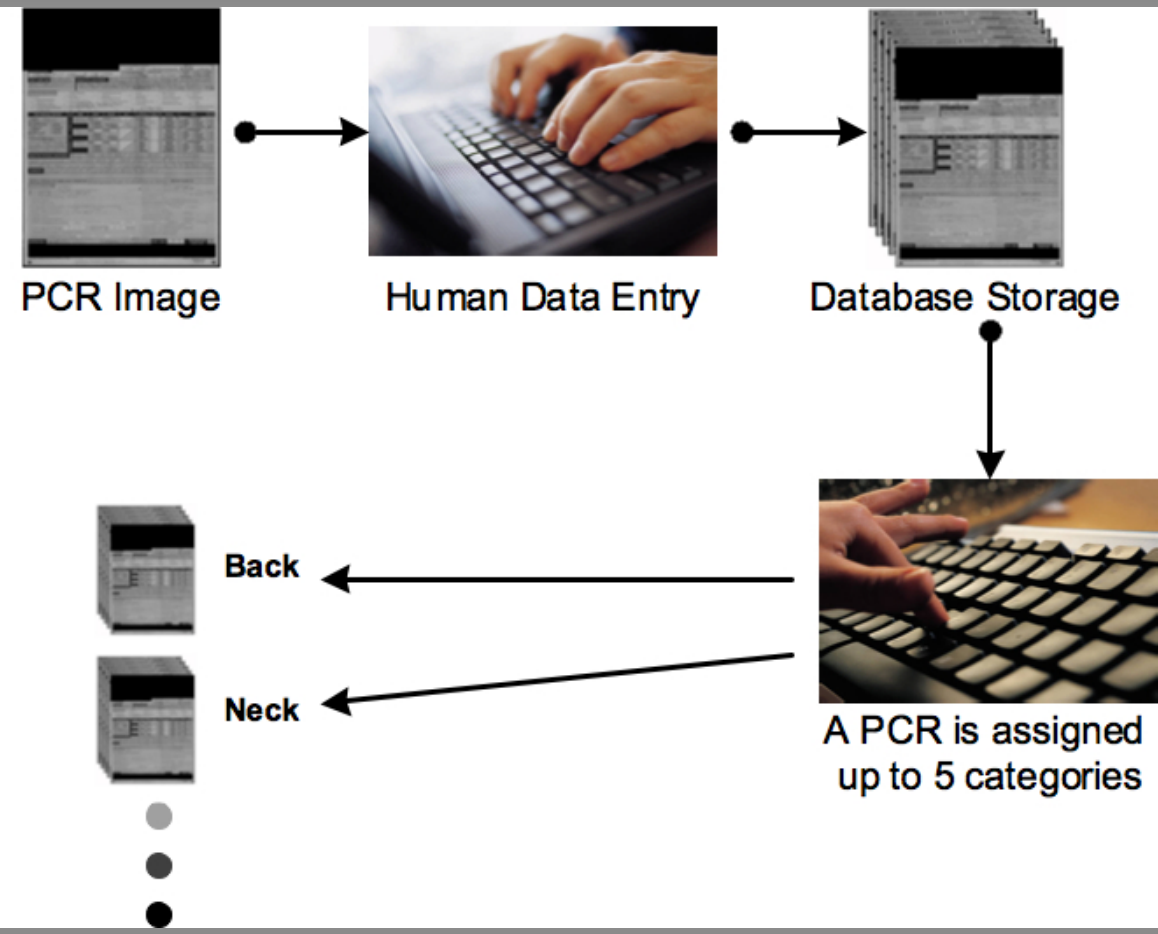
- Circulatory/Cardiovascular
- Digestive
- Endocrine
- Excretory
- Immune
- Integumentary
- Musculoskeletal
- Nervous
- Reproductive
- Respiratory

4 Extremities/Joint Locations

- Arms/Shoulders/Elbows
- Feet/Ankles/Toes
- Hands/Wrists/Fingers
- Legs/Knees

4 General Locations

- Fluid/Chemical Imbalance
- Full Body
- Hospital Transfer/Transport
- Senses



6 Body Range Locations

- | | |
|----------------------|----------------------|
| Abdomen | Head |
| Back/Thoracic/Lumbar | Neck/Cervical |
| Chest | Pelvic/Sacrum/Coccyx |

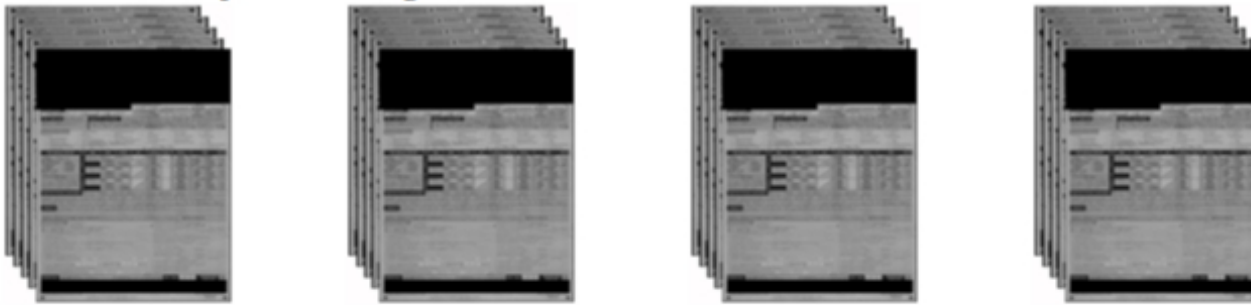
PCR-Category Relationship

Circulatory

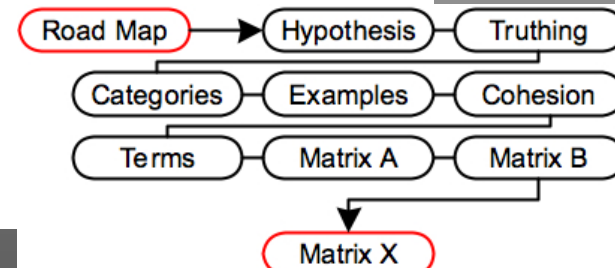
Digestive

Head

Neck

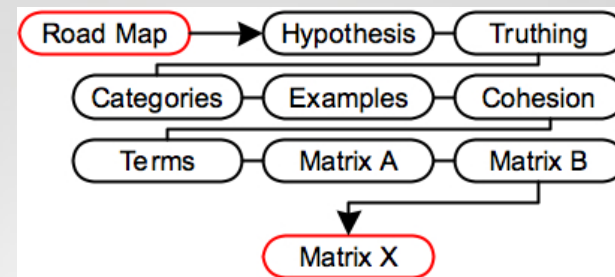


**PCR's are organized by category.
Each category has it's own lexicon (overlap can occur).**



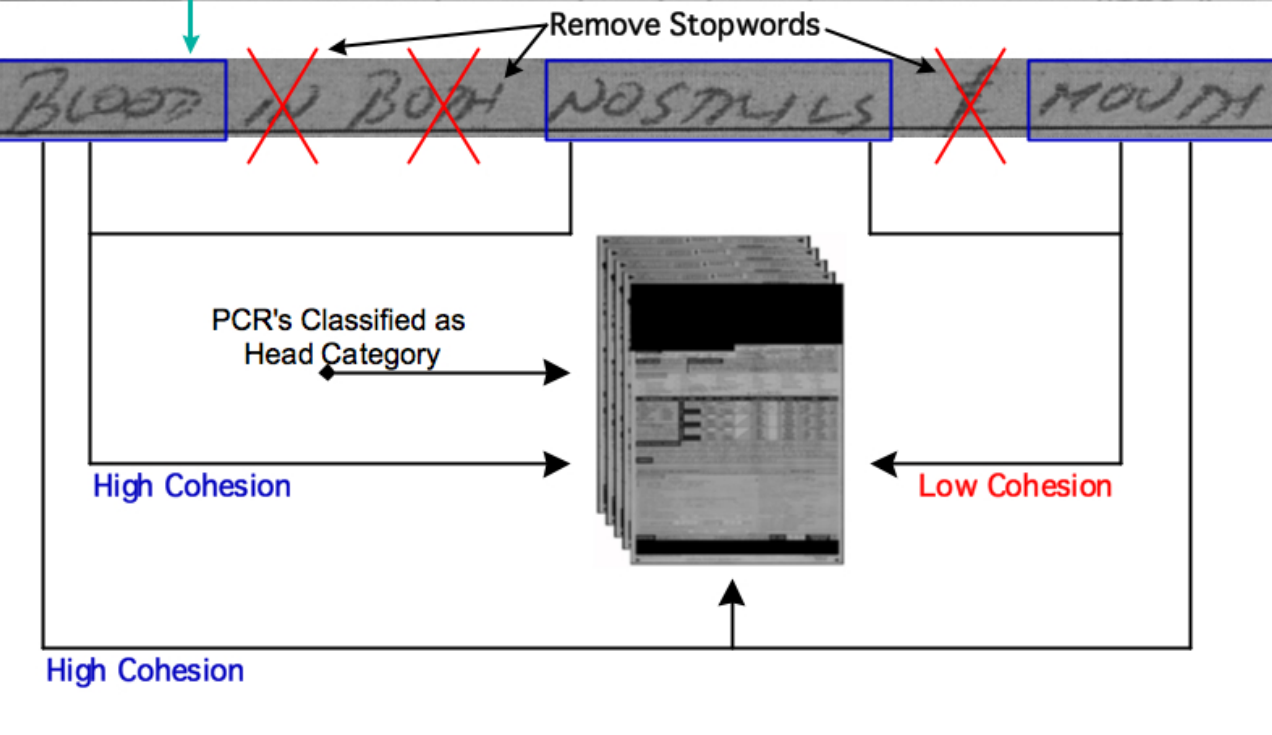
Modeled Category Examples

- ◆ *Example 1:* A patient treated for an emergency pregnancy would be considered the *Reproductive System* category.
- ◆ *Example 2:* A conscious and breathing patient treated for gun shot wounds to the abdominal region would be considered *Circulatory/Cardiovascular System*, due to potential loss of blood.



Phrase Extraction using Cohesion

OBJECTIVE PHYSICAL ASSESSMENT FOUND SUPINE ON GROUND - PE → CAOx3 - PEARL - SPAN LAL ON NOSE
 BLOOD IN BOTH NOSTRILS & MOUTH - LUNGS CTA - BOLT - FIRM PULSATING MASS
 MOVING ALL 4 EXT - SPINAL PRECAUTIONS, → AMP & ROUTINE BLS/ALS
 COMMENTS Tx 2 IV ENROUTE BLOOD GLUCOSE @ 95MG/DL, Tx 3 Don 100%
 A) SYNCOPE FACIAL INJURIES, POSS. AAA
 MEDICAL REASON FOR AMBULANCE TRANSPORT O₂/SPINAL/LAC/IV/EVAL/Tx → REPORT GIVEN TO KAREN - RN.



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

PELVIS

1860	2.49	PAIN HIP
144	0.34	HIP JVD
112	0.39	PAIN CHANGE
275	0.81	HIP FX
110	0.37	HIP CHANGE
163	0.40	JVD PAIN
106	0.40	CAOX3 PAIN
202	0.50	PAIN JVD
213	0.55	PAIN LEG
205	0.42	CHEST HIP

Phrases determined for each category (word order preserved)

$f(w_a, w_b)$ frequency of two words co-occurring

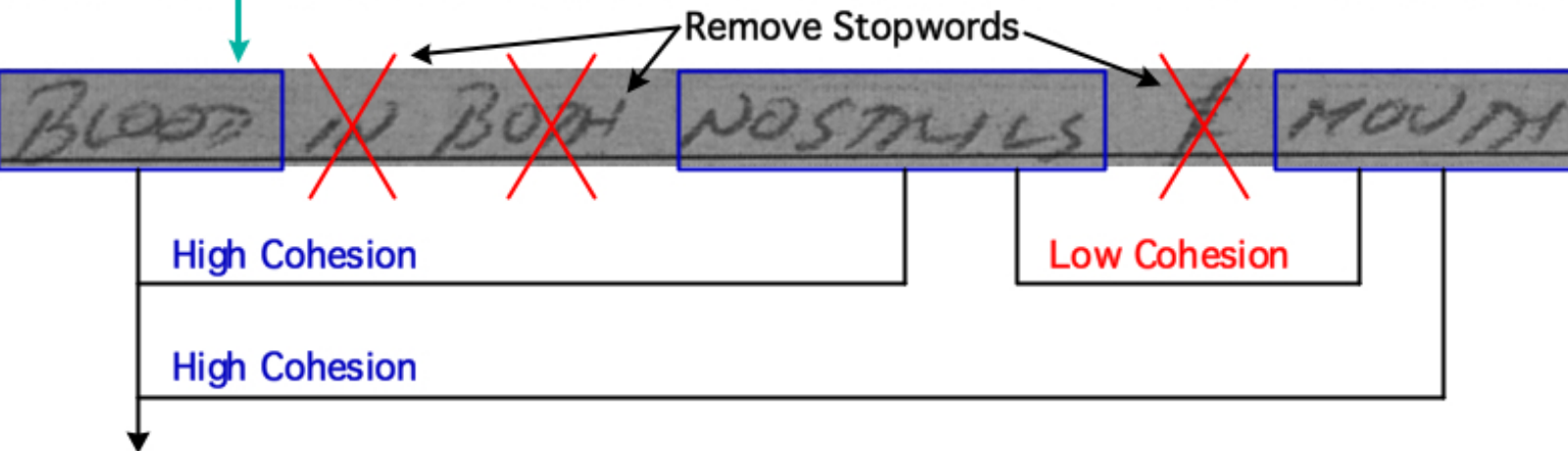
$f(w_a)$ independent word frequency of w_a

z a constant ($z=2$ in this research)

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

Term Construction

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

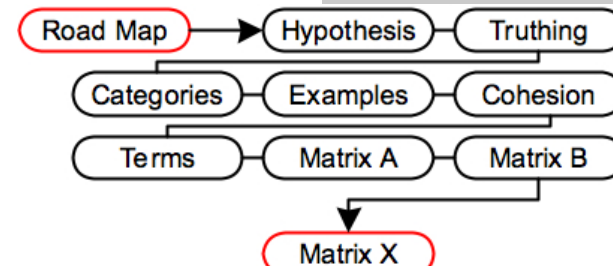


Generate All Possible Character Combinations (Remove Duplicates)

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

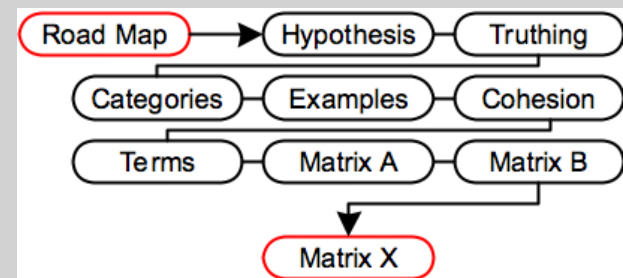
Pairing All High Cohesion Phrases and Construct nm-gram Terms

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



Term Encoding Strategies

- ◆ Terms can be constructed in different ways:
 - ◆ No Spatial Information (NSI)
 - ◆ $*C_1*C_2*$
 - ◆ BLOOD $\rightarrow *L*D*$
 - ◆ Exact Spatial Information (ESI)
 - ◆ xC_1yC_2z
 - ◆ BLOOD $\rightarrow 1L2D0$
 - ◆ Approximate Spatial Information (ASI)
 - ◆ $x_{a-b}C_1y_{c-d}C_2z_{e-f}$
 - ◆ Encode a range of lengths to a code
 - ◆ Code 0: indicate no characters
 - ◆ Code 1: indicates 1 or 2 characters
 - ◆ Code 2: indicates at least 3 characters
 - ◆ BLOOD $\rightarrow 1L1D0$



Combinatorial Analysis (1/2)

Equation

$$C(n) = \left(\left(\sum_{i=1}^{n-1} \binom{n-i}{1} \right) + n \right) = \left(\left(\binom{n}{2} (n-1) \right) + n \right)$$

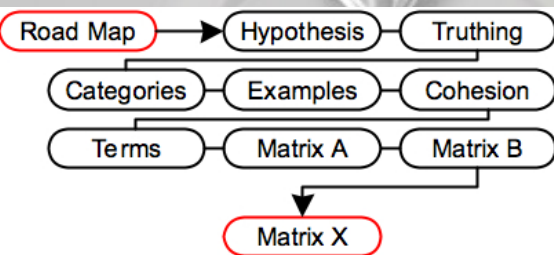
Input: Word of length n s.t. n is at least 1

Output: Quantity of terms to be generated

$$P(a,b) = C(a) \cdot C(b)$$

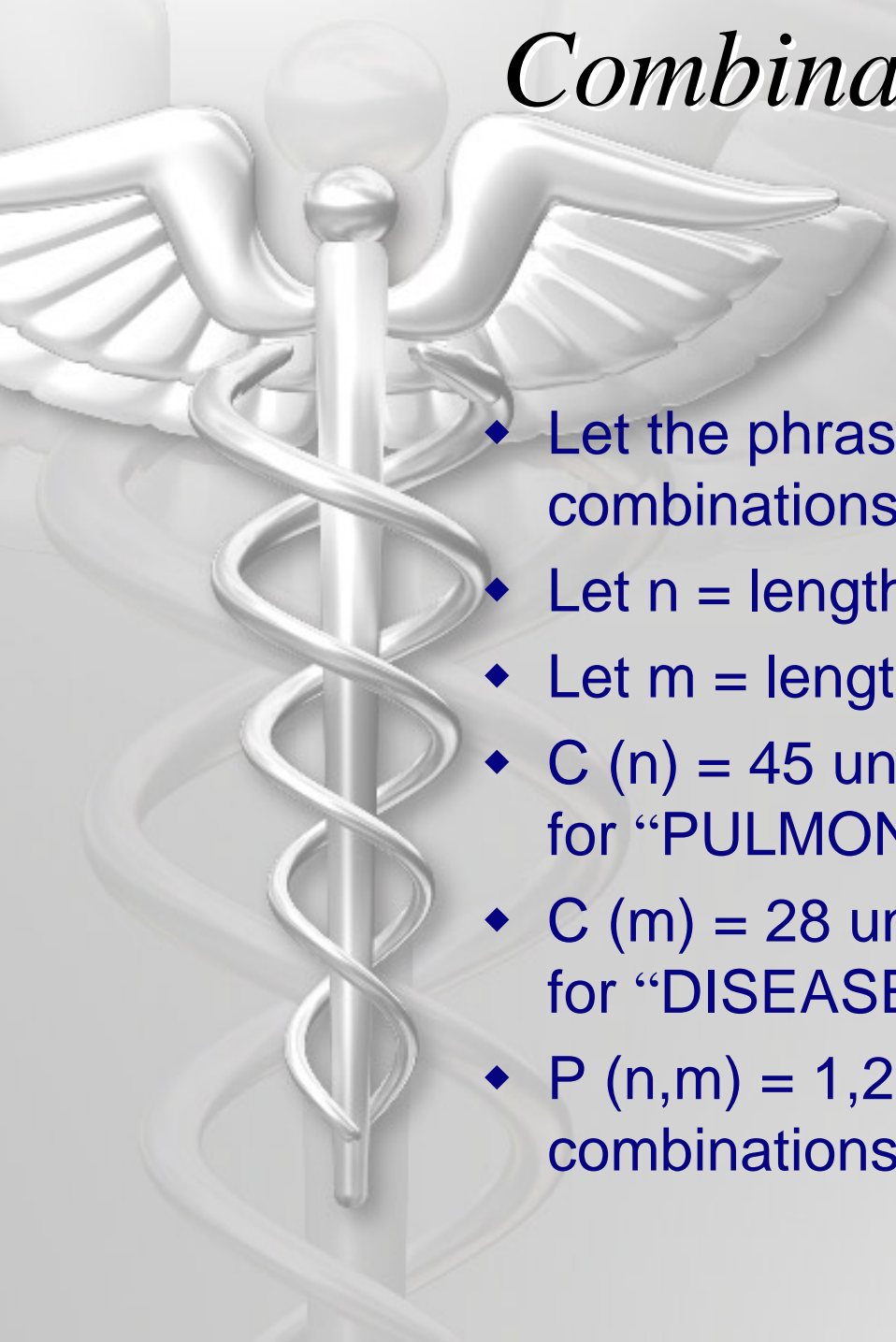
Input: Two word phrase a and b

Output: Quantity of terms to be generated



Combinatorial Analysis (2/2)

Example

- 
- ◆ Let the phrase to evaluate uni/bi-gram combinations be “*PULMONARY DISEASE*”
 - ◆ Let $n = \text{length}(\text{“PULMONARY”}) = 9$
 - ◆ Let $m = \text{length}(\text{“DISEASE”}) = 7$
 - ◆ $C(n) = 45$ uni-gram + bi-gram combinations for “PULMONARY”
 - ◆ $C(m) = 28$ uni-gram + bi-gram combinations for “DISEASE”
 - ◆ $P(n,m) = 1,260$ uni-gram + bi-gram phrase combinations for “*PULMONARY DISEASE*”

LDWR Character Extraction

Chest Pain/Pressure

- Phrases extracted from Circulatory/Cardiovascular-System category.
- Blue indicates successfully extracted top-choice character.
- Red indicates unsuccessfully extracted top-choice character.

CHEST PAIN/PRESSURE (CP) IMAGES		CP BI-GRAMS			
(ID: 48)		*H*U* *B*I*	(ID: 342-2)		*C*S* *A*N*
(ID: 65)		*C*A* *W*T*	(ID: 407-1)		*C*M* *P*A*
(ID: 120)		*E*S* *T*M*	(ID: 407-2)		*H*O* *A*W*
(ID: 183)		*E*O* *P*A*	(ID: 473)		*C*F* *L*A*
(ID: 185)		*C*O* *A*N*	(ID: 643)		*C*H* *A*W*
(ID: 189)		*C*W* *R*S*	(ID: 695-1)		*L*S* *L*R*
(ID: 199)		*F*F* *P*O*	(ID: 98)		*C*S* *S*S*
(ID: 276)		*L*F* *L*M*	(ID: 606)		*C*W* *D*M*
(ID: 342-1)		*C*S* *A*N*	(ID: 695-2)		*C*F* *W*L*

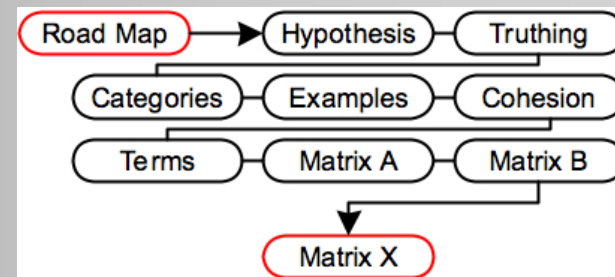
Term-Category Matrix Construction

Construct a Term-Category Frequency Matrix called A

	Digestive System	Circulatory System	Immune System	Head	Neck	Back	Pelvis	Legs	Arms
(bl, no)	0	0	0	32	0	0	0	0	0
(bl, ns)	6	0	0	28	0	28	0	0	0
(bl, nt)	0	0	0	28	0	0	0	0	0
(bl, nr)	0	4	0	31	0	0	0	14	0
(bl, mo)	0	0	0	33	0	55	0	0	0
(bl, mu)	0	0	0	33	0	0	0	41	0
(bl, mt)	0	5	0	42	0	18	0	0	0
(bl, mh)	0	0	0	28	0	0	0	0	0
(bo, no)	9	0	0	28	0	12	0	0	0
(bo, ns)	0	0	0	32	0	38	0	0	0

(bl, mt) occurred 42 times across all PCR's under Head

- Terms on rows.
- Categories on columns.
- Terms generated from high cohesive phrases under category.
- The Term-Category matrix is large.



Matrix Normalization

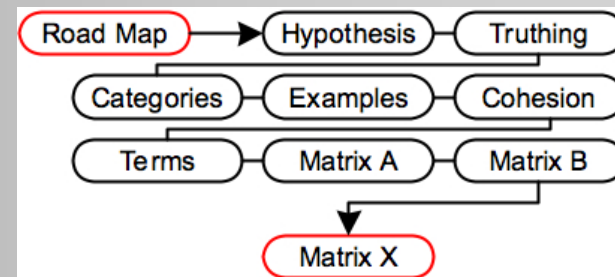
Produce Matrix B by Normalizing Matrix A

	Digestive System	Circulatory System	Immune System	Head	Neck	Back	Pelvis	Legs	Arms
(bl, no)	0	0	0	1.0	0	0	0	0	0
(bl, ns)	.15	0	0	.70	0	.70	0	0	0
(bl, nt)	0	0	0	1.0	0	0	0	0	0
(bl, nr)	0	.12	0	.91	0	0	0	.41	0
(bl, mo)	0	0	0	.51	0	.86	0	0	0
(bl, mu)	0	0	0	.63	0	0	0	.78	0
(bl, mt)	0	.11	0	.91	0	.39	0	0	0
(bl, mh)	0	0	0	1.0	0	0	0	0	0
(bo, no)	.28	0	0	.88	0	.38	0	0	0
(bo, ns)	0	0	0	.64	0	.76	0	0	0

Matrix values are now normalized.
Term Vector's are now all of unit length.

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

- A is the input matrix containing raw frequencies.
- B is the output matrix with normalized frequencies.
- (t,c) is a (term, category) coordinate within a matrix.



Matrix Weighting

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

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

The IDF is used to improve term discrimination ability.

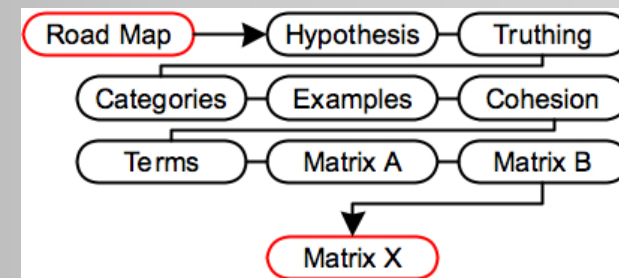
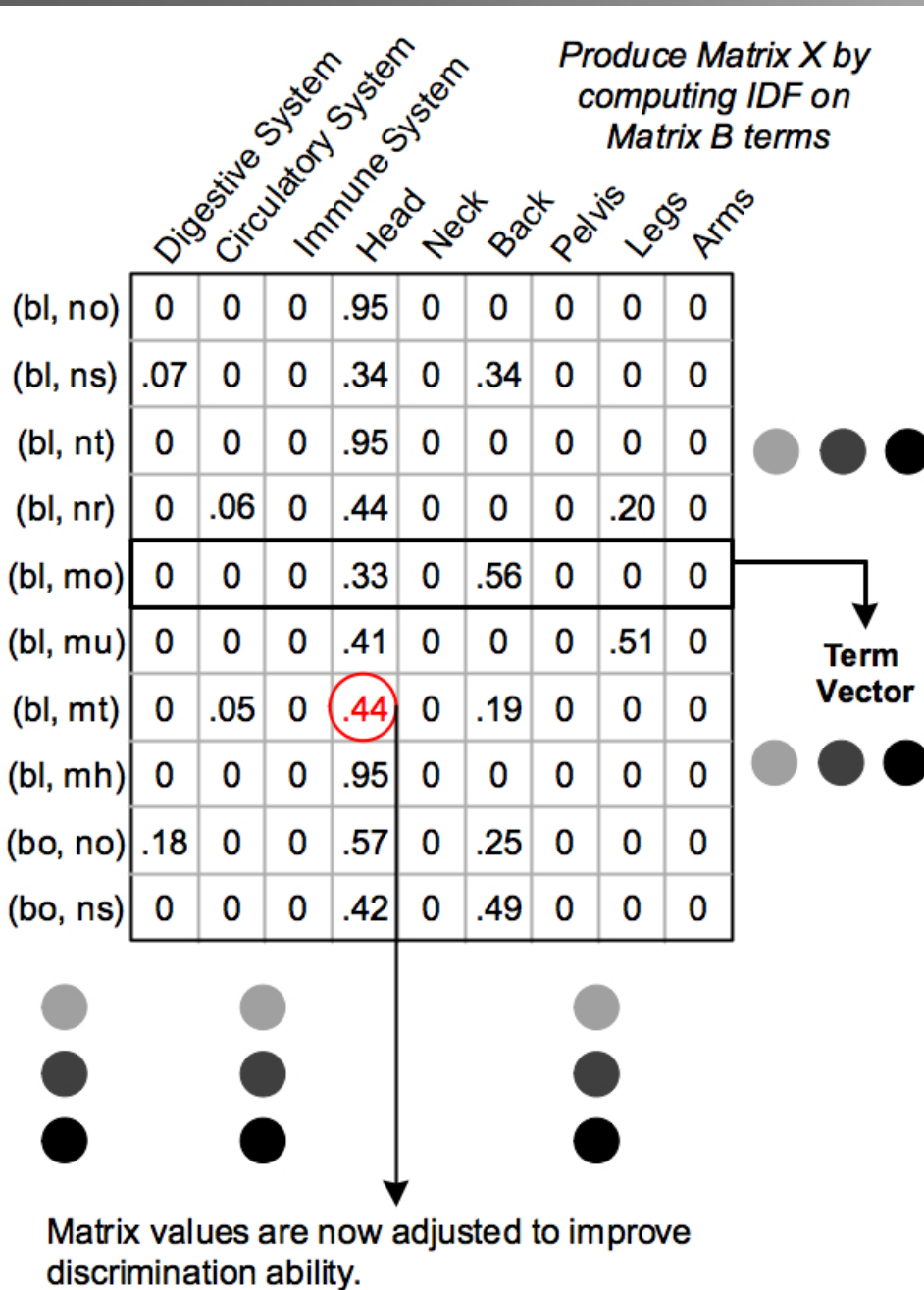
IDF(t) computes IDF on term t

c(t) = number of categories containing term t


B is previous normalized/weighted matrix

Pair (t,c) represents a matrix coordinate

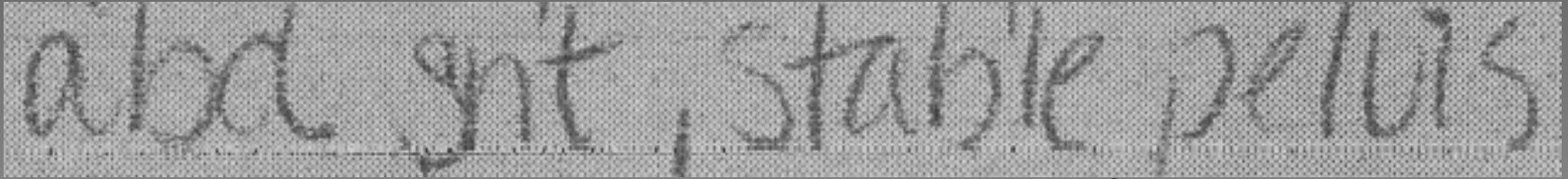
e.g. (bl, mt) occurred in 3 categories.



TOC

- 
- ◆ Introduction
 - ◆ Prior Work
 - ◆ Binarization
 - ◆ Topic Categorization: Training
 - ◆ **Topic Categorization: Testing**
 - ◆ Results
 - ◆ Applications

Text Image Extraction

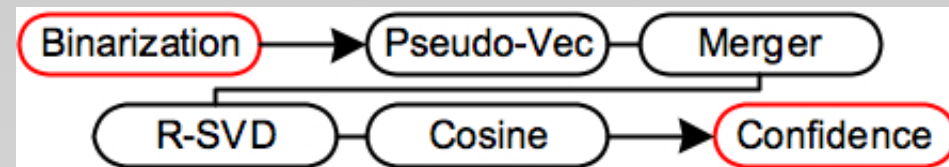


Binarization and post-processing



Recognizer Confusion:

abd	“a ? q”
snt	“t”
stable	“t”, “a”, “b”, “e ? c”
pelvis	“e ? c”, “l”, “v ? u”



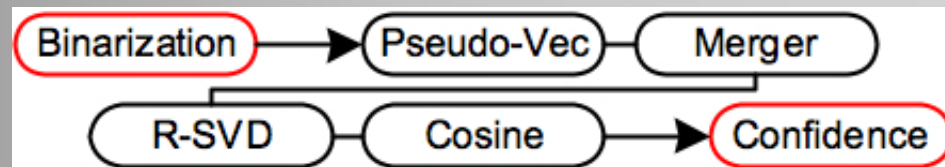
Pseudo-Category Vector

Pseudo-Category Vector

(bl, no)	0
(bl, ns)	2
(bl, nt)	1
(bl, nr)	1
(bl, mo)	1
(bl, mu)	0
(bl, mt)	2
(bl, mh)	0
(bo, no)	3
(bo, ns)	1

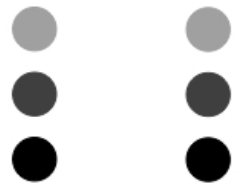
The LDWR recognition engine found the nm-gram sequence (bo, no) three times on a single PCR.

- ◆ A vector Q of size $(|T| \times 1)$ is constructed.
- ◆ Each cell in the vector represents the frequency of the uni/bi-gram sequence generated from the LDWR.

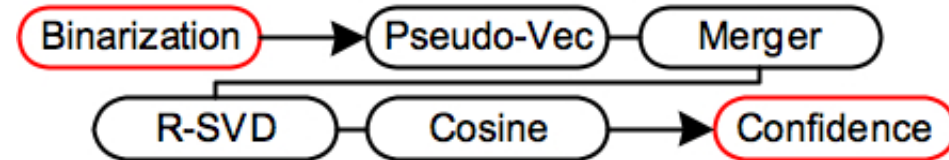


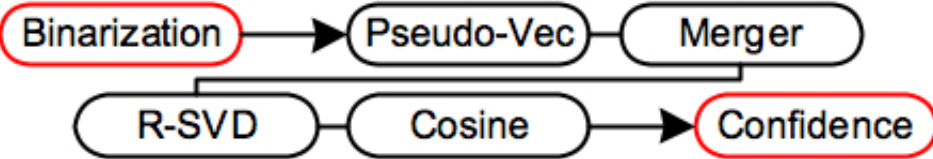
Vector Merger

	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

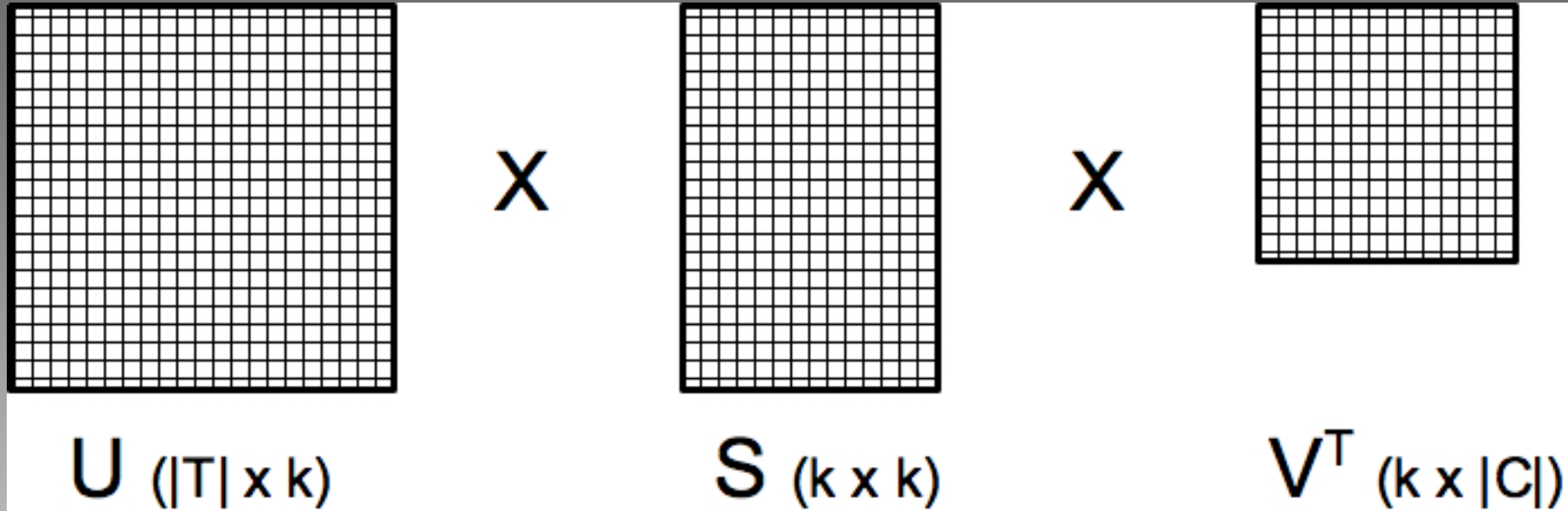

 Vector Q is now added to the original normalized/weighted matrix X

- Q is merged with the original matrix X.
- The purpose is to produce a vector V_q which represents Q in k -dimensional space.





SVD



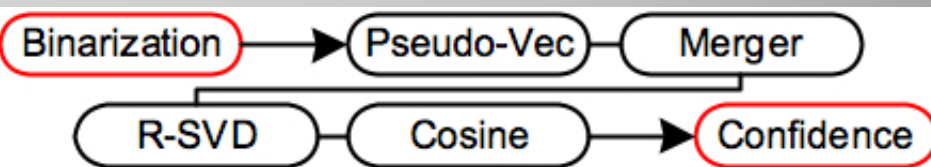
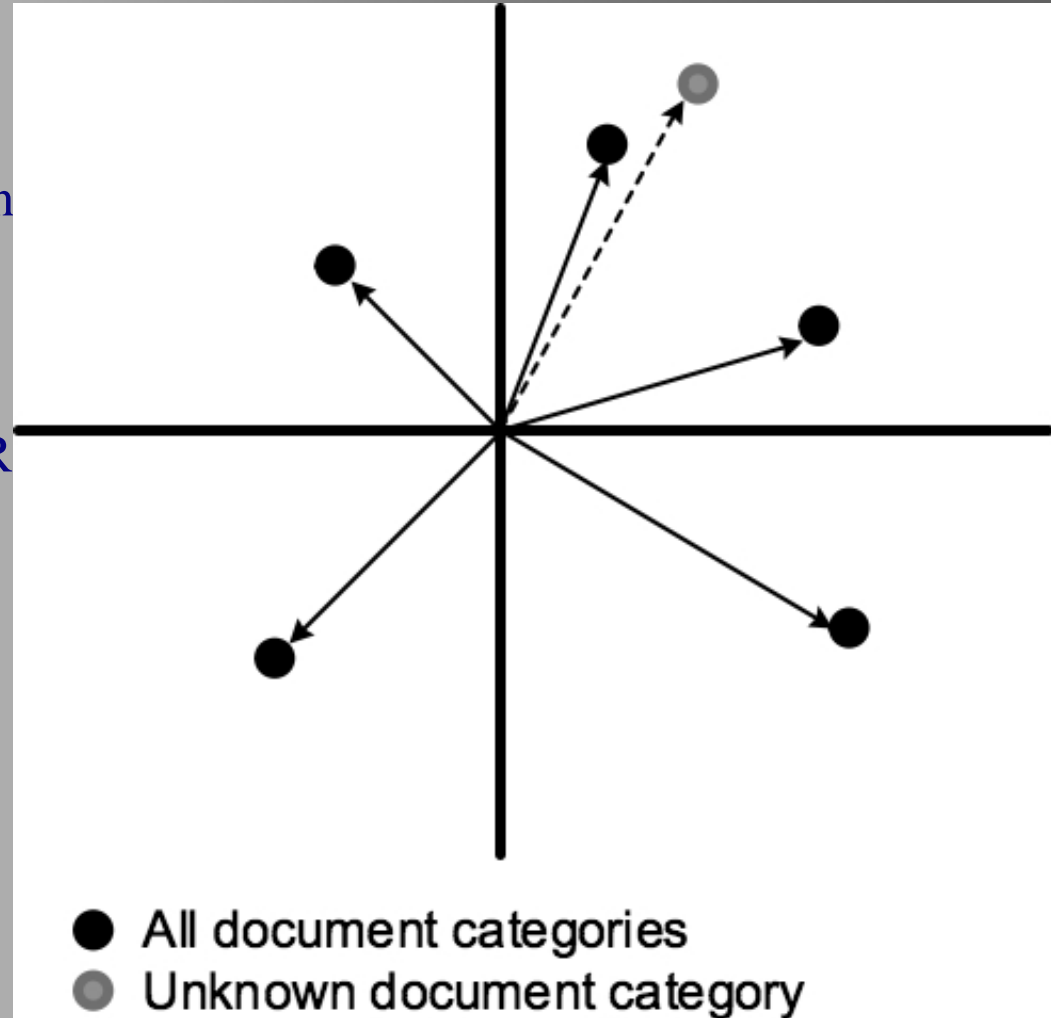
$$X = U \cdot S \cdot V^T$$

- ◆ X is a matrix which is decomposed into 3 matrices
- ◆ U is a $T \times k$ matrix representing term vectors
- ◆ S is a $k \times k$ matrix
- ◆ V is a $k \times C$ matrix representing category vectors
- ◆ The correctly scaled pseudo-category vector is now produced by computing $V_q \cdot S_r$

Category Querying

- ◆ Vector x : correctly scaled pseudo-category
- ◆ Vector y : A single category vector in $V_r \cdot S_r$
- ◆ The value z is the score between x and each y .
 - ◆ z is the ‘closeness’ of the LDWR data to a category

$$z = \cos(x, y) = \frac{x \cdot y^T}{\sqrt{\sum_{i=1}^n x_i^2 \cdot \sum_{i=1}^n y_i^2}}$$



(Simplified Dimensional Space)

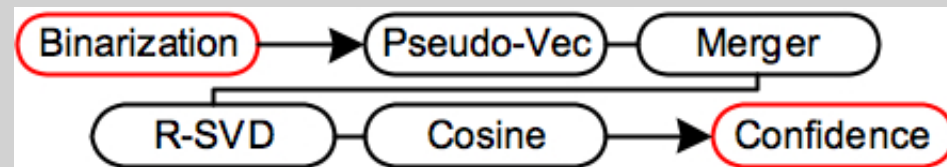
Sigmoid Fitting

- ◆ Lucent Technologies showed that the use of a sigmoid function instead of the cosine score can reduce error rate.
- ◆ The cosine score is mapped to a sigmoid function.
- ◆ Least Squares Regression line is computed as follows:


$$a = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} \quad b = \frac{1}{n} \left(\sum_{i=1}^n y_i - a \sum_{i=1}^n x_i \right)$$

- ◆ The final *confidence* score is produced by mapping the regression line to the sigmoid function:

$$\text{confidence}(a, b, z) = \frac{1}{\left(1 + e^{-(az+b)} \right)}$$

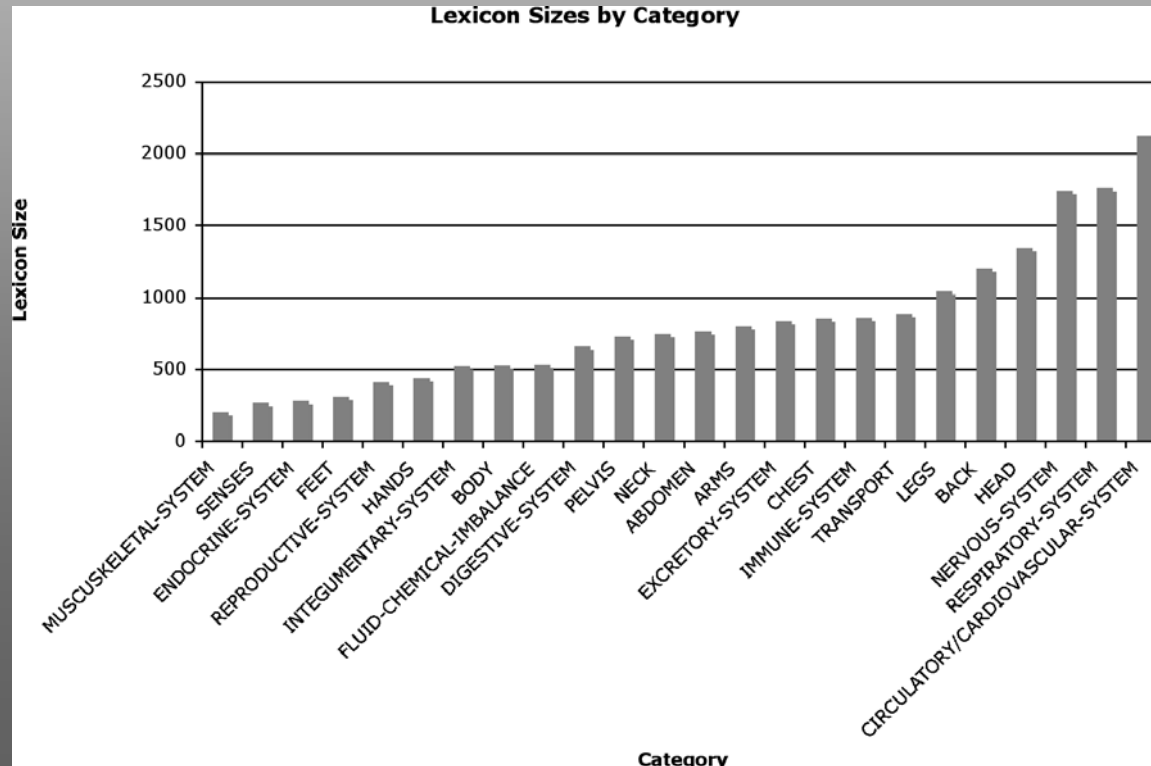


TOC

- 
- ◆ Introduction
 - ◆ Prior Work
 - ◆ Binarization
 - ◆ Topic Categorization: Training
 - ◆ Topic Categorization: Testing
 - ◆ **Results**
 - ◆ Applications

Lexicon Organization

UNIQUE LEXICON SIZE FROM TRAINING SET:	5,628
TOTAL WORDS ACROSS LEXICONS (OVERLAPPING):	8,156
AVERAGE CATEGORIES PER FORM	1.4
MAXIMUM CATEGORIES PER FORM	5
TOTAL TRAINING SET PCR'S:	750
TOTAL TESTING SET PCR'S:	62
TOTAL TRAINING SET WORDS:	42,226
TOTAL TEST SET WORDS:	3,089



Handwriting Recognition(HR) Rates

	CL	CLT	AL	ALT	SL	SLT	RL	RLT
ACC	76.34%	76.92%	63.52%	66.59%	70.51%	71.51%	70.70%	71.06%
ERR	71.93%	69.95%	57.24%	47.12%	62.26%	59.44%	62.04%	59.45%
RAW	23.31%	25.32%	32.31%	41.73%	30.30%	32.73%	30.62%	32.63%
LS	5,628	8,156	1,193	1,246	2,514	2,620	2,401	2,463
!L	n/a	n/a	23.89%	8.02%	16.06%	10.46%	16.61%	12.23%
!HL	n/a	n/a	33.33%	97.98%	48.19%	73.99%	46.59%	62.96%

ACC: accept rate

ERR: error rate

RAW: raw recognition rate

LS: lexicon size for experiment

!L: percentage of truther word not in reduced lexicon

!HL: percentage from !L set in which a human could not decipher the word

CL: complete lexicon

AL: category oracle machine

SL: synthetic term generation

RL: reduced lexicon approach

“-T” indicates test deck lexicon inserted into reduced lexicon.

HR Interpretation

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

HR Performance based on Lexicon Reduction

METRIC	VALUE
Reduction Accuracy (α)	0.33
Reduction Degree (ρ)	0.83
Reduction Efficacy (η)	0.06
Lexicon Density (words)	1.07 → 0.87
Lexicon Density (uni/bi-grams)	0.50 → 0.78

Word present after reduction

Words removed from lexicon

Effectiveness Rating $\eta = \Delta_{LDWR} \times \alpha^{(1-\rho)}$

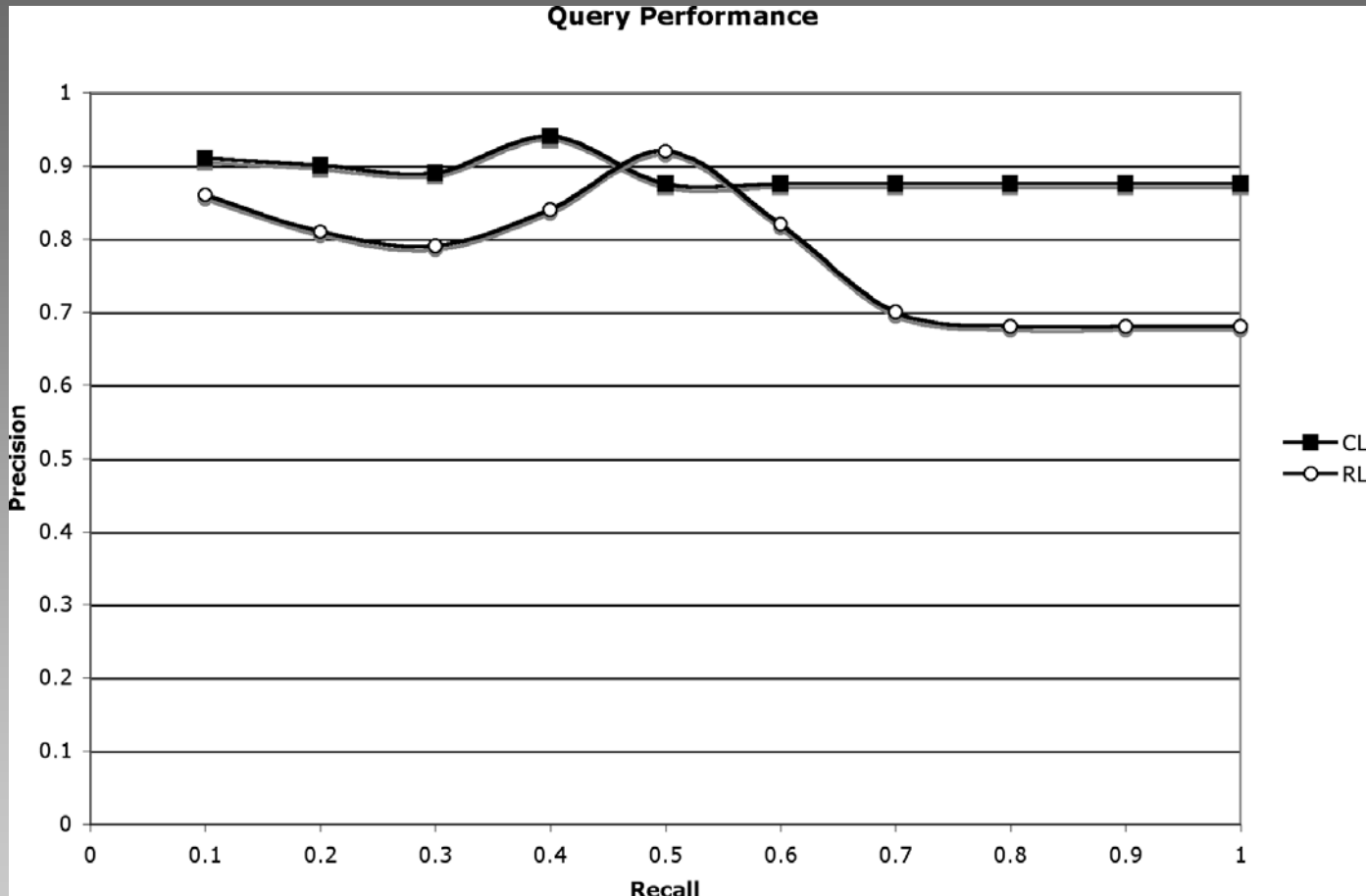
Lexicon becomes less confusing

N-Gram's are more common

Lexicon Reduction Performance

Retrieval Performance

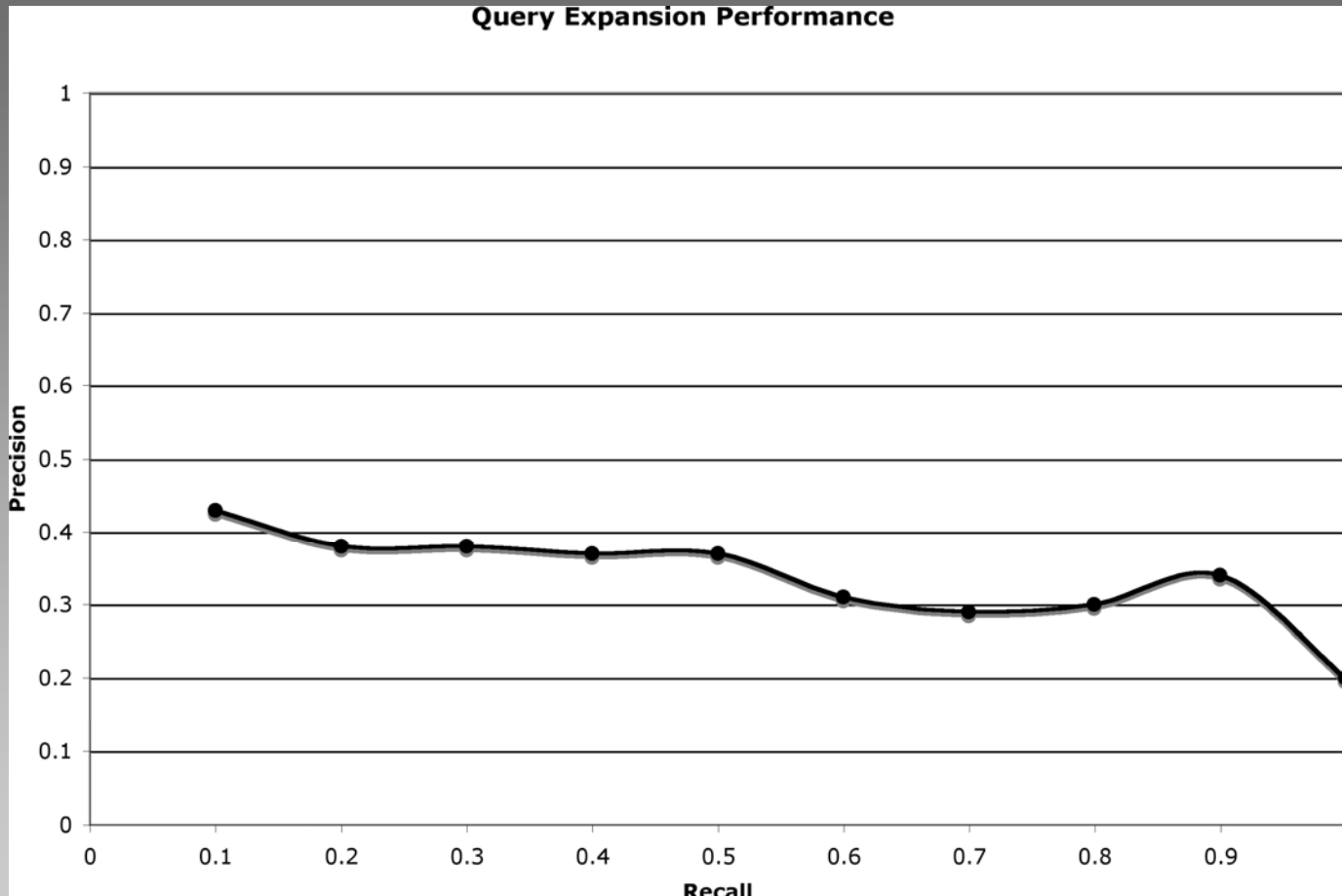
Cohesive Phrase Search



- ◆ Querying 800 forms with 1,250 cohesive phrases from isolated deck of 200 forms:
- ◆ CL returns 0 forms 73% of the time
- ◆ RL returns 0 forms 23% of the time


Retrieval Performance

Query Expansion Search

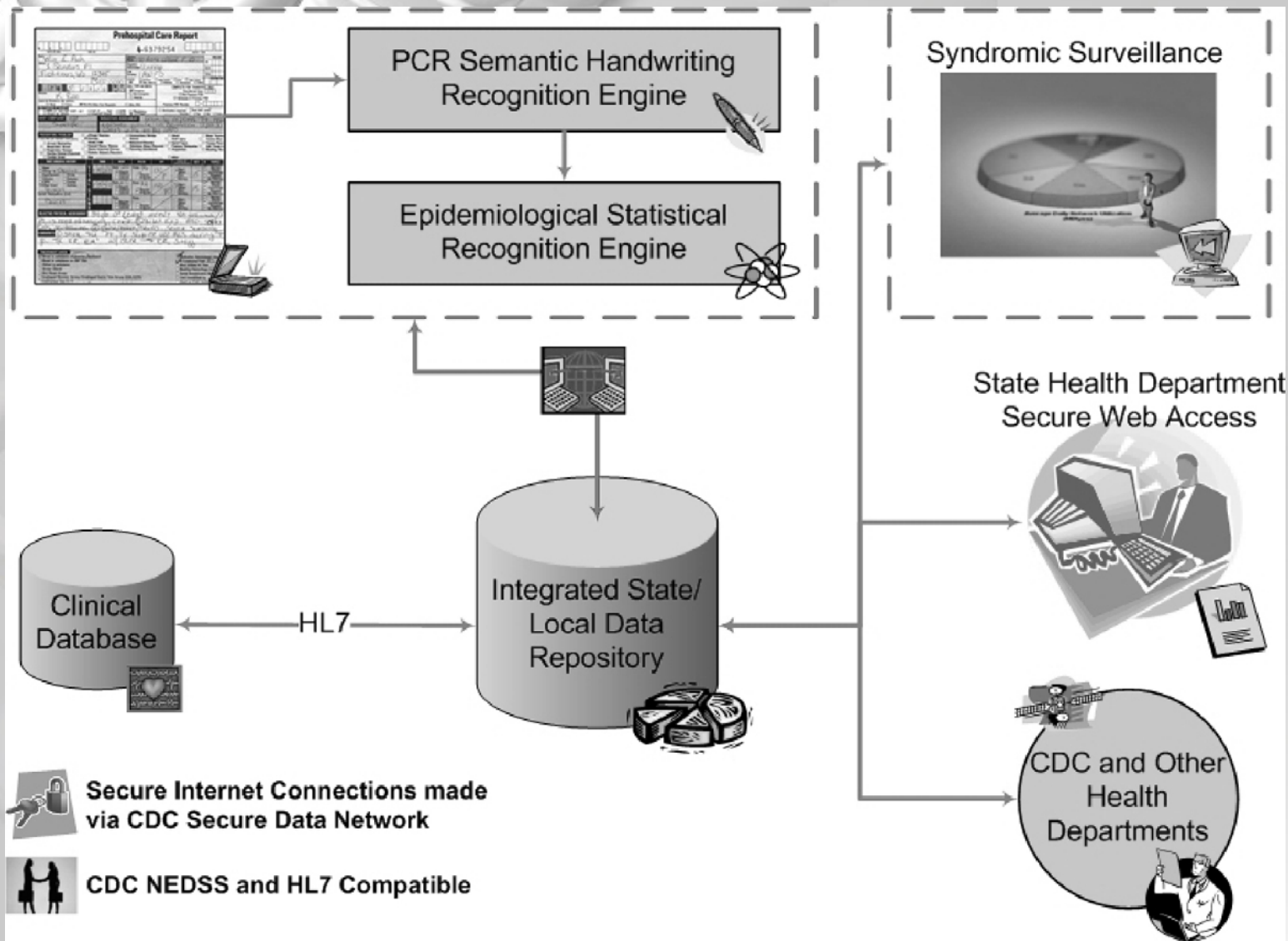


- ◆ Cohesive phrase is decomposed into uni/bi-gram components
- ◆ Search performed by matching only uni/bi-gram components
- ◆ Poor performance

TOC

- 
- ◆ Introduction
 - ◆ Prior Work
 - ◆ Binarization
 - ◆ Topic Categorization: Training
 - ◆ Topic Categorization: Testing
 - ◆ Results
 - ◆ **Applications**

Health Surveillance



Contributions

- ✓ The first binarization and post-processing strategy on carbon forms.
- ✓ The first application of recognition of handwritten medical forms.
- ✓ The first search engine using handwritten forms.
- ✓ Construction of the first data set of actual handwritten emergency medical documents for use in document analysis research.
- ✓ A paradigm showing a mapping between character encodings to a topic categorization used for lexicon reduction. This strategy is reusable for other lexicon driven handwriting recognizers that are based on character segmentation.
- ✓ New metrics for measuring the performance of lexicon reduction systems.
- ✓ Compatibility with standard information protocols used by Health Level 7 (HL7) and the Center for Disease Control (CDC).
- ✓ A framework for automated, centralized and secure health surveillance network.
- ✓ An advanced software system with diverse visual interfaces and command-line execution modes.