

Text Extraction from Handwritten Medical Form Images

February 13, 2006

Robert Milewski
Venu Govindaraju

Outline

I. Project Overview

II. Document Image Analysis Challenges

- Binarization of carbon copy forms
- Unconstrained Handwriting Recognition (OHR)

III. Preview: IR of Handwritten Medical Forms

Project Overview

• Pre Hospital Care Report

WNY: 250,000 filed a year

NYC: 50,000 filed in a day

• OHR issues

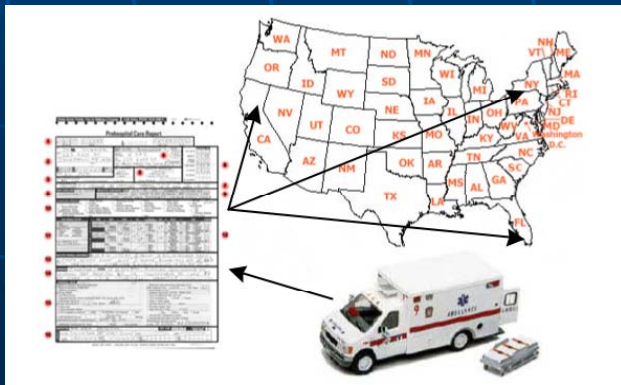
Loosely constrained writing style

Large lexicons

Heterogeneous data

6,700 carbon forms stored at 300 DPI

1000 PCR forms ground truthed



Prehospital Care Report

4-6379254

Name: Win E. Pooh Agency Name: 100 Acre Wood F.D. MILEAGE: END CALL REC'D:

Address: Fictitious, Wb 12345 Dispatch Information: Unresp ENROUTE:

Call Location: 555-0000 Call: 1 A.W.F.D. TOTAL: ARRIVED AT SCENE:

CALL TYPE AS REQ'D: Emergency Non-Emergency Stand-by Transferred from: No Previous PCR Unknown if Previous PCR

PROCESSED BY: K. Roco PREVIOUS PCR NUMBER:

CARE IN PROGRESS ON ARRIVAL: None Citizen PD/Other First Responder Other EMS

MECHANISM OF INJURY: MVA (seat belt used →) Fall of feet GSW Machinery Extrication required Seat belt used? Yes No Unknown

CHIEF COMPLAINT: SP SUBJECTIVE ASSESSMENT: family reports Pt. had syncopal episode while in bathroom upon arrival Pt. LOA x3 w/ 102 on pt arrival

PRESENTING PROBLEM: Allergic Reaction Unconscious/Unresp. Shock Major Trauma OB/GYN Syncope Stroke/CVA Seizure Head Injury Trauma-Blunt Burns Airway Obstruction General Illness/Malaise Behavioral Disorder Spinal Injury Trauma-Penetrating Environmental Respiratory Arrest Gastro-Intestinal Distress Substance Abuse (Potential) Fracture/Dislocation Soft Tissue Injury Heat Respiratory Distress Cardiac Related (Potential) Poisoning (Accidental) Amputation Bleeding/Hemorrhage Cold Cardiac Arrest Diabetic Related (Potential) Other Obvious Death Hazardous Materials

| PAST MEDICAL HISTORY | TIME | RESP | PULSE | B.P. | LEVEL OF CONSCIOUSNESS | GCS | R | PUPILS | L | SKIN | STATUS |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------|-----------------|-----------------|------------|----------------------------------------------------------------------------------------------------------|-----------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------|--------|
| <input type="checkbox"/> Allergy to <u>Denies</u> | <u>13:20</u> | Rate: <u>20</u> | Rate: <u>76</u> | <u>130</u> | <input checked="" type="checkbox"/> Alert <input type="checkbox"/> Pain <input type="checkbox"/> Unresp. | <u>15</u> | <input type="checkbox"/> Normal <input type="checkbox"/> Dilated <input type="checkbox"/> Constricted <input type="checkbox"/> Sluggish <input type="checkbox"/> No Reaction | <input type="checkbox"/> Cool <input type="checkbox"/> Warm <input type="checkbox"/> Moist <input type="checkbox"/> Dry | <input type="checkbox"/> Pale <input type="checkbox"/> Cyanotic <input type="checkbox"/> Flushed <input type="checkbox"/> Jaundiced | <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> | |
| <input type="checkbox"/> Hypertension <input type="checkbox"/> Stroke <input type="checkbox"/> Diabetes <input type="checkbox"/> COPD <input type="checkbox"/> Other (List) <u>Smokes</u> | <u>13:30</u> | Rate: <u>20</u> | Rate: <u>72</u> | <u>130</u> | <input checked="" type="checkbox"/> Alert <input type="checkbox"/> Pain <input type="checkbox"/> Unresp. | <u>15</u> | <input type="checkbox"/> Normal <input type="checkbox"/> Dilated <input type="checkbox"/> Constricted <input type="checkbox"/> Sluggish <input type="checkbox"/> No Reaction | <input type="checkbox"/> Cool <input type="checkbox"/> Warm <input type="checkbox"/> Moist <input type="checkbox"/> Dry | <input type="checkbox"/> Pale <input type="checkbox"/> Cyanotic <input type="checkbox"/> Flushed <input type="checkbox"/> Jaundiced | <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> | |
| <input type="checkbox"/> Seizures <input type="checkbox"/> Cardiac <input type="checkbox"/> Asthma | | Rate: <u>20</u> | Rate: <u>72</u> | <u>130</u> | <input checked="" type="checkbox"/> Alert <input type="checkbox"/> Pain <input type="checkbox"/> Unresp. | <u>15</u> | <input type="checkbox"/> Normal <input type="checkbox"/> Dilated <input type="checkbox"/> Constricted <input type="checkbox"/> Sluggish <input type="checkbox"/> No Reaction | <input type="checkbox"/> Cool <input type="checkbox"/> Warm <input type="checkbox"/> Moist <input type="checkbox"/> Dry | <input type="checkbox"/> Pale <input type="checkbox"/> Cyanotic <input type="checkbox"/> Flushed <input type="checkbox"/> Jaundiced | <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> | |
| Current Medications (List) <u>Denies</u> | | Rate: <u>20</u> | Rate: <u>72</u> | <u>130</u> | <input checked="" type="checkbox"/> Alert <input type="checkbox"/> Pain <input type="checkbox"/> Unresp. | <u>15</u> | <input type="checkbox"/> Normal <input type="checkbox"/> Dilated <input type="checkbox"/> Constricted <input type="checkbox"/> Sluggish <input type="checkbox"/> No Reaction | <input type="checkbox"/> Cool <input type="checkbox"/> Warm <input type="checkbox"/> Moist <input type="checkbox"/> Dry | <input type="checkbox"/> Pale <input type="checkbox"/> Cyanotic <input type="checkbox"/> Flushed <input type="checkbox"/> Jaundiced | <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> | |

OBJECTIVE PHYSICAL ASSESSMENT: 84 y/o M COA x3 HEENT NO trauma / ABN / PEARL Pt is hard of hearing. Chest: @ Bilat. Exp. ABD: soft x4 Exts: no trauma @ pulse/motion/NEURO. Some sweating noted

COMMENTS: Dstick 76 Pt Tx supine w/ ALS during Tx. P. TX ER RM w/ care → ER Staff

TREATMENT GIVEN: Moved to ambulance on stretcher/backboard Moved to ambulance on stair chair Walked to ambulance Airway Cleared Oral/Nasal Airway Esophageal Obturator Airway/Esophageal Gastric Tube Airway (EOA/EGTA) Endotracheal Tube (E/T) Oxygen Administered @ 100 L.P.M. Method NRB Suction Used Artificial Ventilation Method C.P.R. in progress on arrival by: Citizen PD/Other First Responder Other C.P.R. Started @ Time 13:20 Time from Arrest 10 Minutes AED/ECG Monitored (Attach Tracing) (Rhythm(s)) NSR Defibrillation/Cardioversion No. Times Manual Semi-automatic Medication Administered (Use Continuation Form) IV Established Fluid NS RVD Cath. Gauge 210 Most Inflated @ Time Bleeding/Hemorrhage Controlled (Method Used:) Spinal Immobilization Neck and Back Limbs Immobilized by Fixation Traction Heat or (Cold) Applied Warming Included @ Time Method Restraints Applied, Type Baby Delivered @ Time In County Alive Stillborn Male Female Transported in Trendelenburg position Transported in left lateral recumbent position Transported with head elevated Other Vitals TX

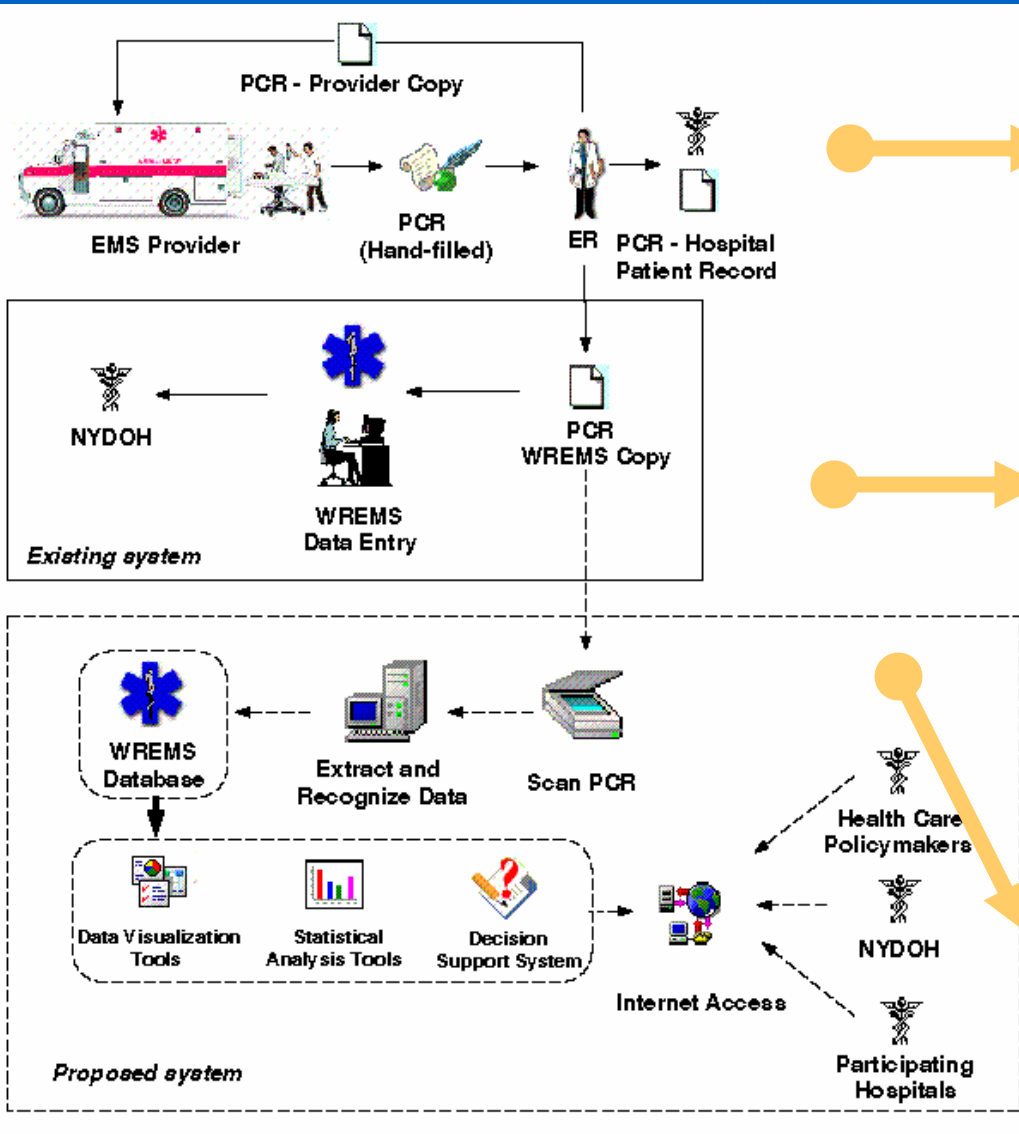
DISPOSITION (See list) Admitted Discharged Deceased Other Vitals TX

IN CHARGE: TIGGER DRIVER'S NAME: Eyore

EMT # AEMT # EMT # AEMT #

AGENCY COPY/WHITE RESEARCH COPY/YELLOW HOSPITAL PATIENT RECORD COPY/PINK

Current PCR Life Cycle

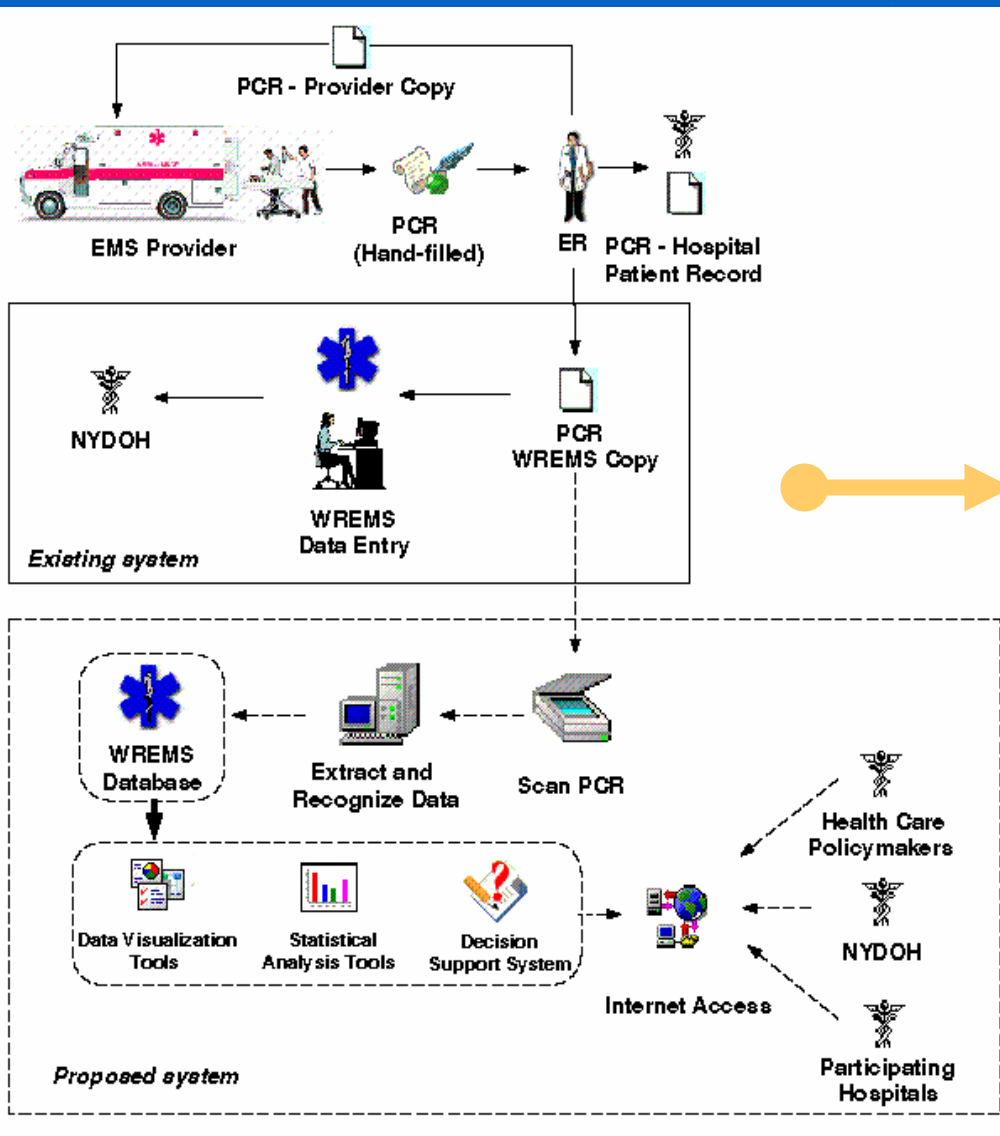


- Ambulance rescues patient
- Patient information hand filled in emergency environment
- Form passed to ER

- **Transcription**
- Data Dissemination
 - Aggregated copies (Carbon)
 - Data extracted
 - Data Analysis

- Scan carbon copy
- Information Extraction
- Data Analysis

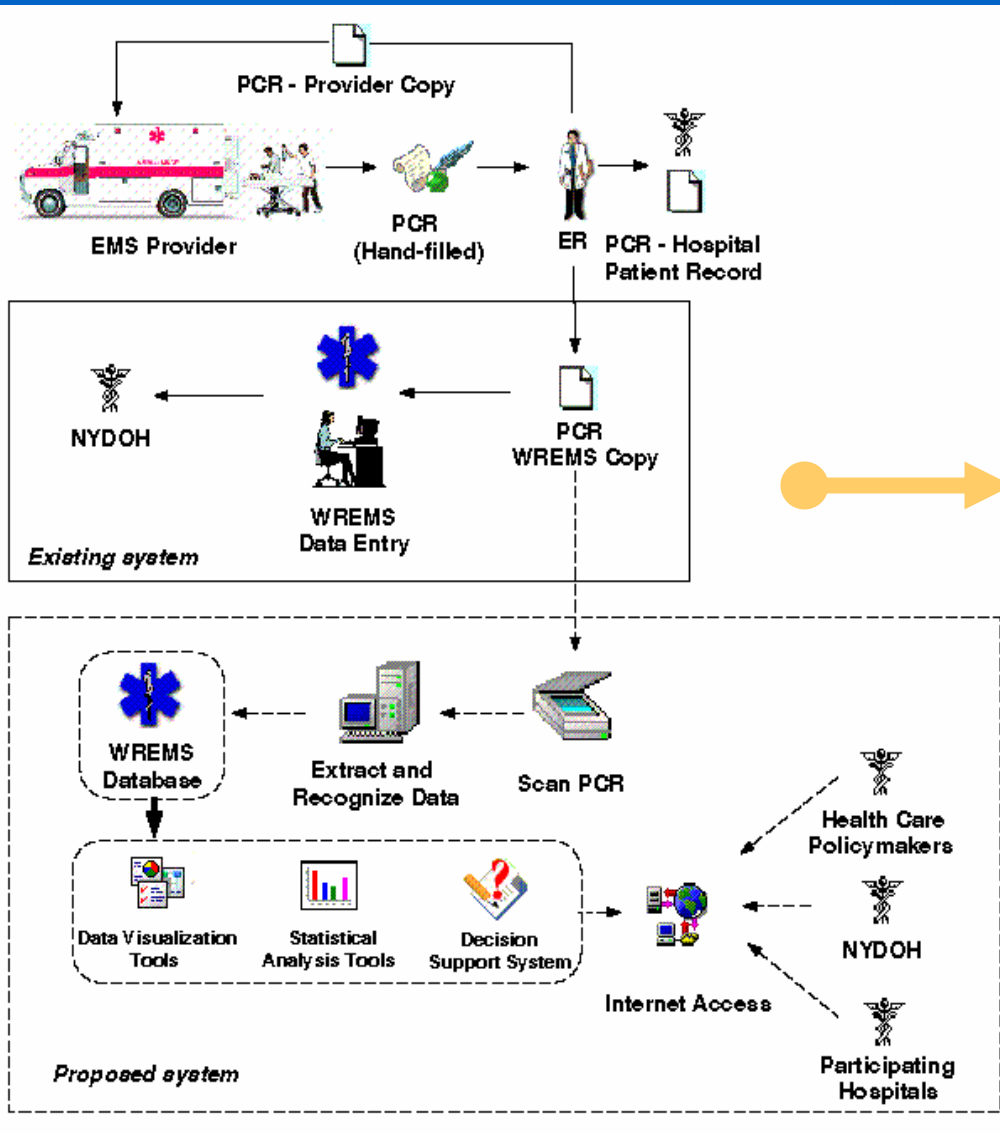
Problems



Major problems

1. No real-time event correlation
2. Labor intensive transcriptions

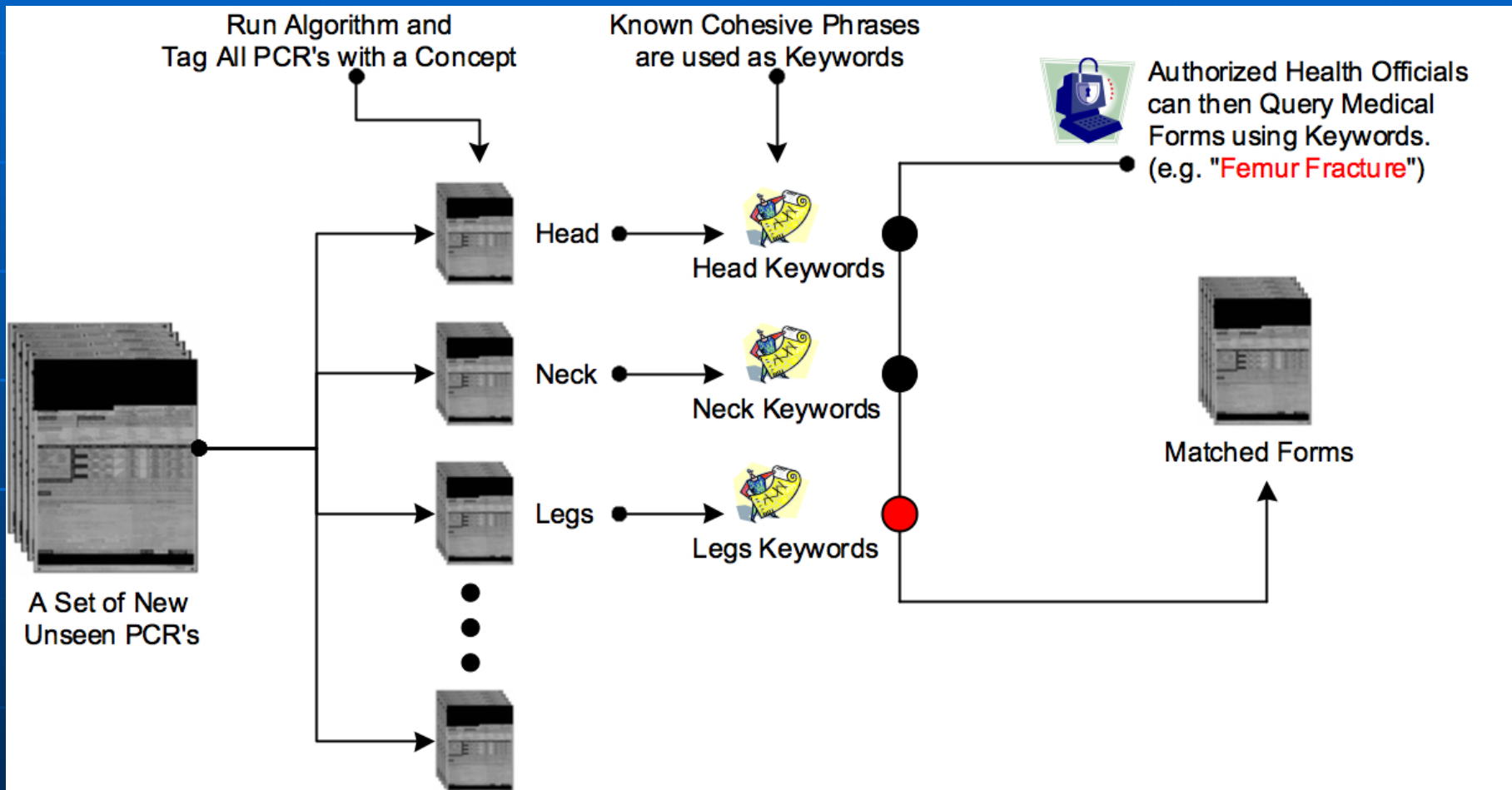
Proposed Solution



DIA Solution
OHR

Application

Medical Form Search Engine



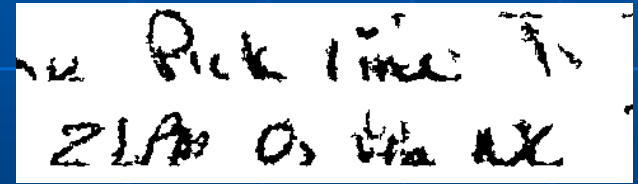
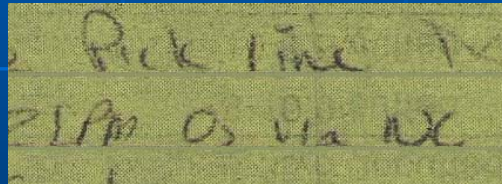
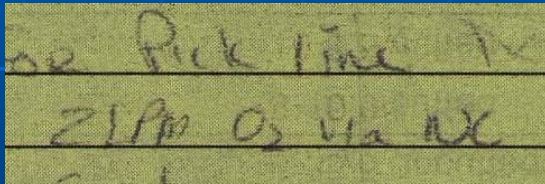
DIA Tasks

- **Form Registration/Classification**
 - PCR form templates are stored in a database.
- **Preprocessing**
 - ✓ **Binarization**
 - Form dropout (line removal etc.)
 - Word Segmentation
- **Zone Detection/Segmentation**
- **OHR**
 - ✓ **Word Recognition**
 - Keyword Spotting
- ✓ **Indexing and Retrieval**

Binarization

Challenges

- Yellow color media
- Noisy carbon mesh
- Uneven carbon imprint
- Poor, shaky writing (emergency context)



| | |
|------------------------------------------------|--------------------------------------|
| OBJECTIVE PHYSICAL ASSESSMENT | (+) ARC = CAOx3, Hee |
| | bilat, Abdsnt xy, Pelvis Stable, (+) |
| | Pt was upset and was weeping |
| COMMENTS | Pt + x = stretcher → Amb → |
| | staff |
| MEDICAL REASON FOR AMBULANCE TRANSPORT: | Dec for be |

Previous Work- Application Driven

- Gatos et al, PR 06
- Shi & Govindaraju, ICPR 04
- Wolf & Doerman, ICPR 02
- Sauvola & Niblack, PR 00
- Wu & Manmatha, SPIE 98
- Yang & Yan, PR 98
- Liu & Srihari, PAMI 97
- Trier & Taxt, PAMI 95
- Degraded handwritten documents
- Palm leaf manuscripts (historical, handwritten)
- Low Quality Text, MRF Model
- Degraded Documents
- Machine printed text on complex backgrounds
- Engineering drawings images in poor illumination
- Machine printed mail piece images
- Map Images; Comparison survey Niblack + Yanowitz/Bruckstein

Approaches

■ Global Thresholding

- Assume separation between foreground and background
- Cannot handle shadows, low contrast, smear
- Otsu's Method, Kittler et al, Kapur et al

■ Adaptive Thresholding

- Niblack: $T_{\text{window}} = \mu + k\sigma^2$

window = 1-2 chars; k = -0.2

- leaves residual noise in background

- Sauvola: $T_{\text{window}} = \mu + 1 - k(1 - \sigma^2/R)$

k=0.5; R=128

■ Background Modeling

- Shi & Govindaraju (ICPR 04)
- Gatos et al (PR 2006)

Domain Attributes

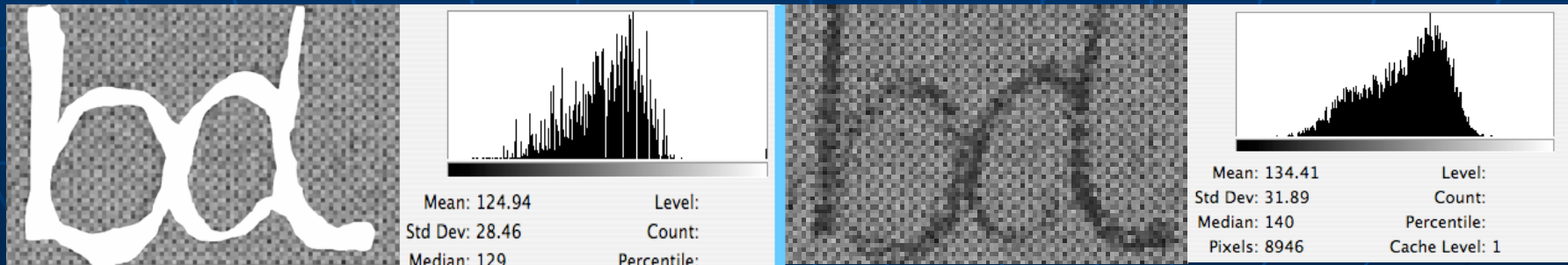
- ☑ Text occupies separable range in the histogram
- ☑ Prior knowledge of number of peaks in the histogram
- ☑ Text images contain highly structured stroke units

■ Handwriting Attributes

- Irregular strokes
- Character stroke width consistent
- Non-uniform intensity of strokes

■ Carbon Mesh Attributes

- Unevenly Textured ; Smears/Smudges
- Carbon intensity *may* be same as stroke intensity



Global Histogram Analysis

OBJECTIVE PHYSICAL ASSESSMENT

75 yo male PAU x 3, HEENT pt wears glasses neck. ⊖TD, ⊖JVD
Chest = expansion, clear sounds abd gut, stable pelvis, PUS x 4 ↑↓ extremities
Ⓜ arm head trap, prior IV site, 1° vitals 1° assessment enroute, tx to Buffalo Mercy

COMMENTS

without incident. PT walked from table to stretcher -

Pixels less than threshold (1 peak)



OBJECTIVE PHYSICAL ASSESSMENT

75 yo male PAU x 3, HEENT pt wears glasses neck. ⊖TD, ⊖JVD
Chest = expansion, clear sounds abd gut, stable pelvis, PUS x 4 ↑↓ extremities
Ⓜ arm head trap, prior IV site, 1° vitals 1° assessment enroute, tx to Buffalo Mercy

COMMENTS

without incident. PT walked from table to stretcher -

Pixels greater than threshold (1 peak)
Background does not have a separate peak



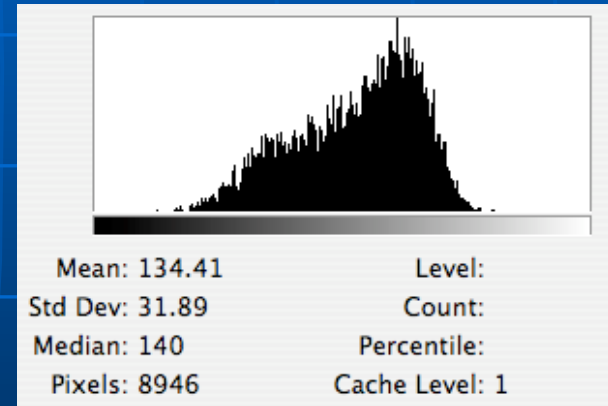
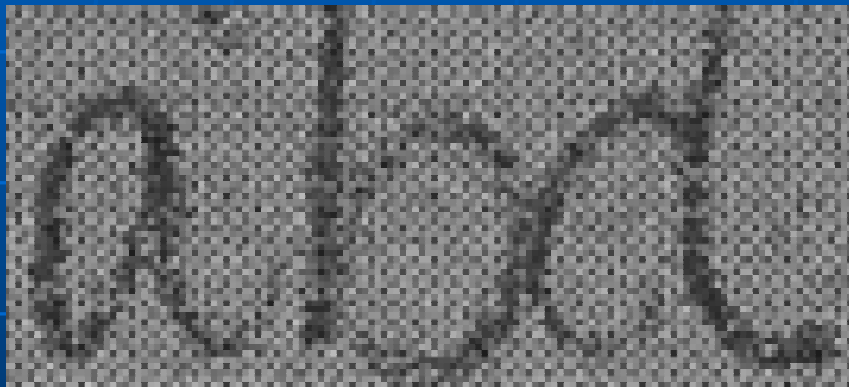
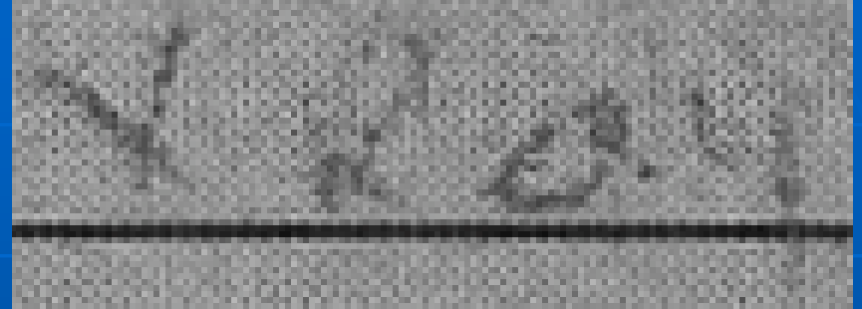
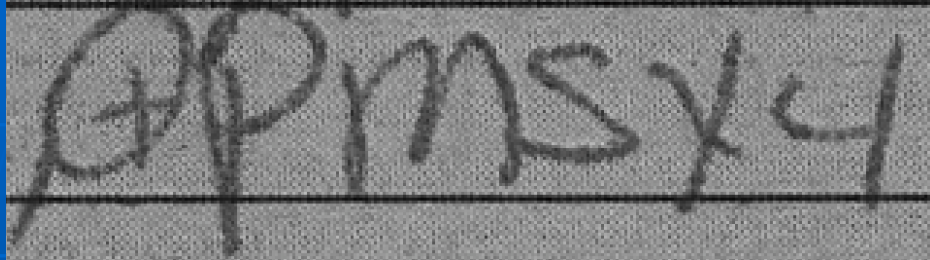
OBJECTIVE PHYSICAL ASSESSMENT

75 yo male PAU x 3, HEENT pt wears glasses neck. ⊖TD, ⊖JVD
Chest = expansion, clear sounds abd gut, stable pelvis, PUS x 4 ↑↓ extremities
Ⓜ arm head trap, prior IV site, 1° vitals 1° assessment enroute, tx to Buffalo Mercy

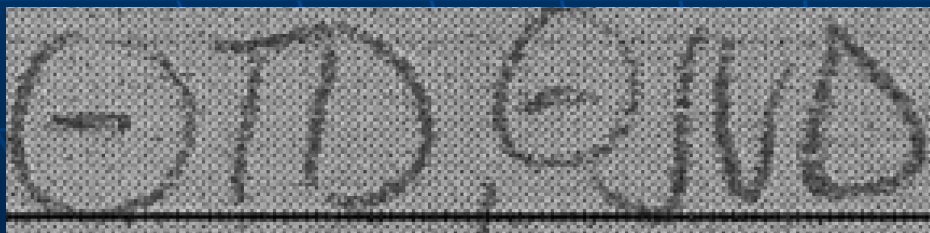
COMMENTS

without incident. PT walked from table to stretcher -

Local Histogram Analysis



Locally- single peak histogram



Locally Adaptive Method

| | |
|------------------------------------------------|-----------------------------------------------------------------------------------|
| OBJECTIVE PHYSICAL ASSESSMENT | APC= CAAx3, Hce bilat. Abdant X4, Delus State G H was upset and was wearing |
| COMMENTS | H was → Struck → Amb Staff |
| MEDICAL REASON FOR AMBULANCE TRANSPORT: | Delus |

Otsu's Global Thresholding Method

| | |
|------------------------------------------------|-----------------------------------------------------------------------------------|
| OBJECTIVE PHYSICAL ASSESSMENT | APC= CAAx3, Hce bilat. Abdant X4, Delus State G H was upset and was wearing |
| COMMENTS | H was → Struck → Amb Staff |
| MEDICAL REASON FOR AMBULANCE TRANSPORT: | Delus |

Background Modeling

OBJECTIVE PHYSICAL ASSESSMENT (+) ARC's CAX's, Heer
bilat, Abdsnt xy, Pelvis Stable, (+)
Pt was upset and was weeping

COMMENTS Pt + x → stretcher → Amb →
rest in O/P staff

MEDICAL REASON FOR AMBULANCE TRANSPORT: pleur tube

[This section contains a large area of blurred text, likely representing a continuation of the medical notes or a redacted portion of the document.]

Background Substraction

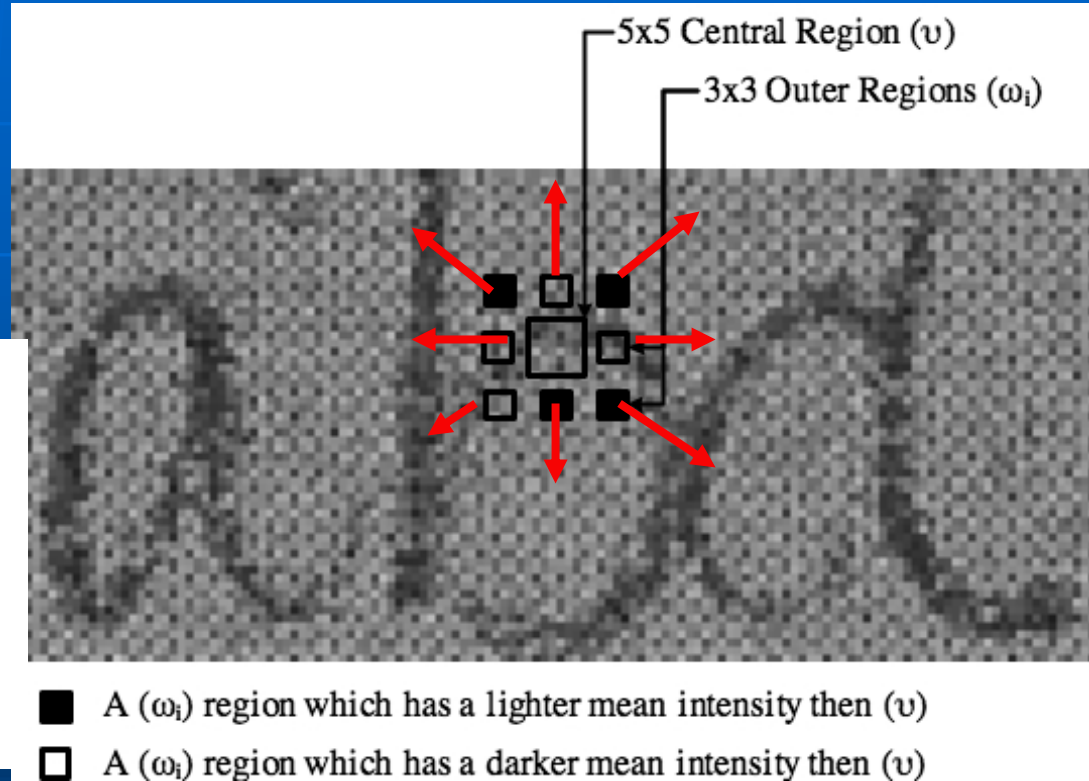
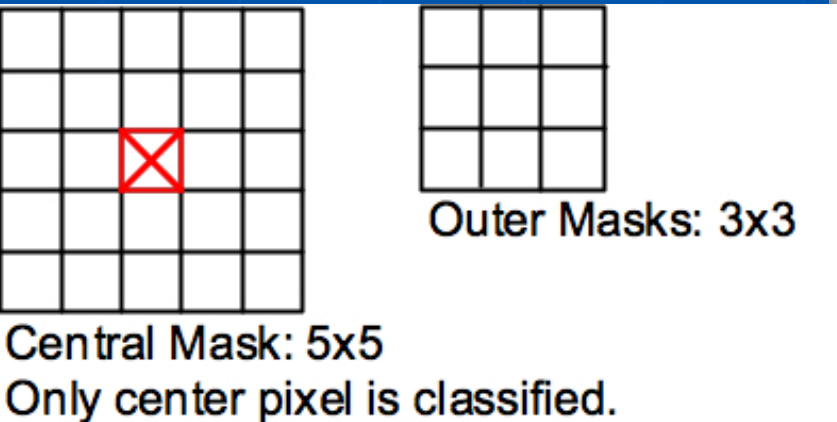
| | |
|------------------------------------------------|------------------------------------------------------------------------------------------|
| OBJECTIVE PHYSICAL ASSESSMENT | ARC = CAOx3, Hec bilat. Abdsnt xy, Pelvis Stable, (+) Pt was upset and was weeping |
| COMMENTS | Pt + x → stretcher → Amb with 2 staff |
| MEDICAL REASON FOR AMBULANCE TRANSPORT: | pelvic |

Additional Postprocessing

| | |
|------------------------------------------------|------------------------------------------------------------------------------------------|
| OBJECTIVE PHYSICAL ASSESSMENT | ARC = CAOx3, Hec bilat. Abdsnt xy, Pelvis Stable, (+) Pt was upset and was weeping |
| COMMENTS | Pt + x → stretcher → Amb with 2 staff |
| MEDICAL REASON FOR AMBULANCE TRANSPORT: | pelvic |

Algorithm

- Central mask size fn of stroke width (local)
- Thresholds are adaptive



Step 1: A central mask v slides across image

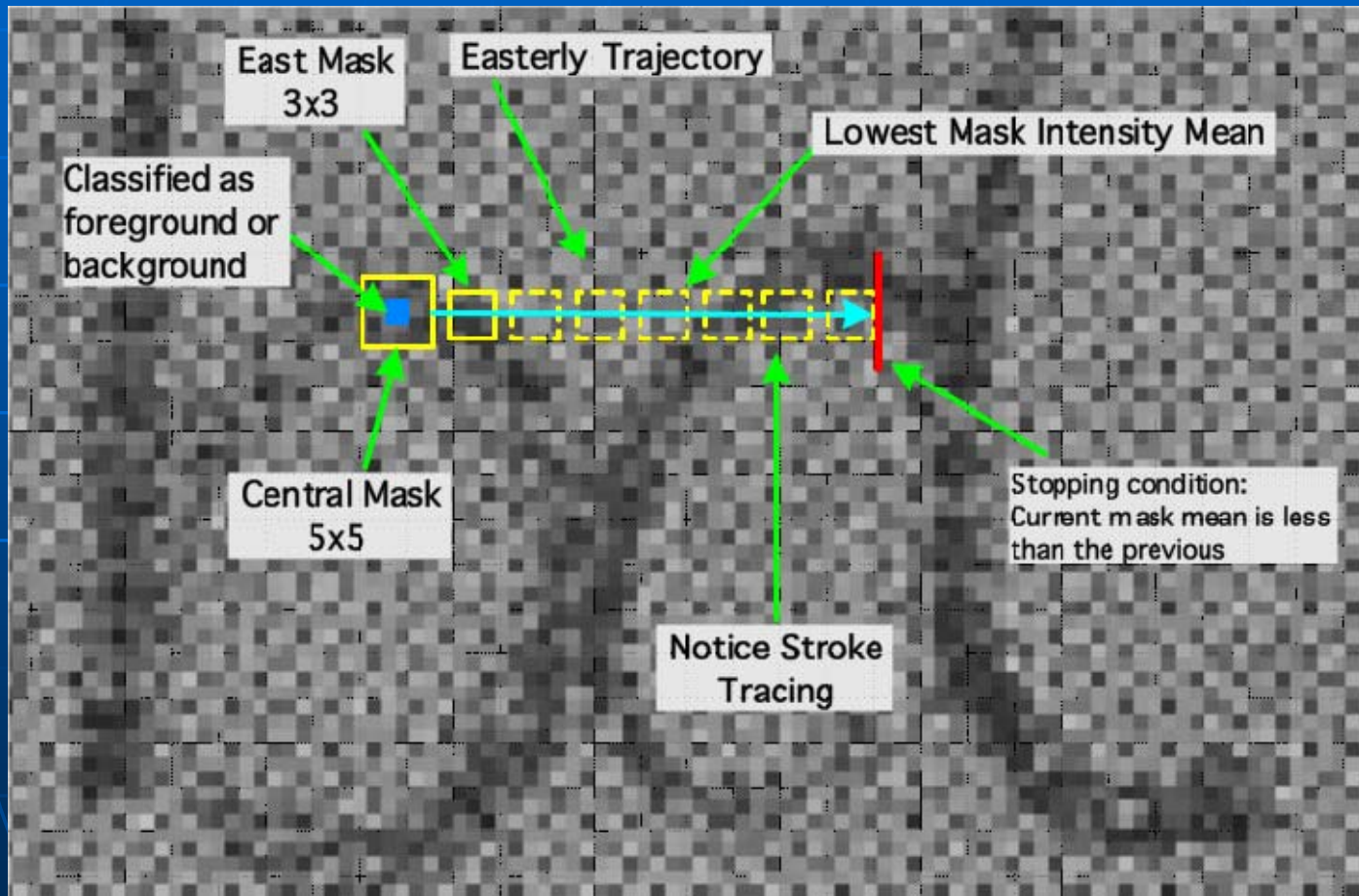
Step 2: At each central mask location, 8 scouting masks ω are launched in 8 directions

Step 3: Scouting masks search along trajectories for the lowest average mask intensity

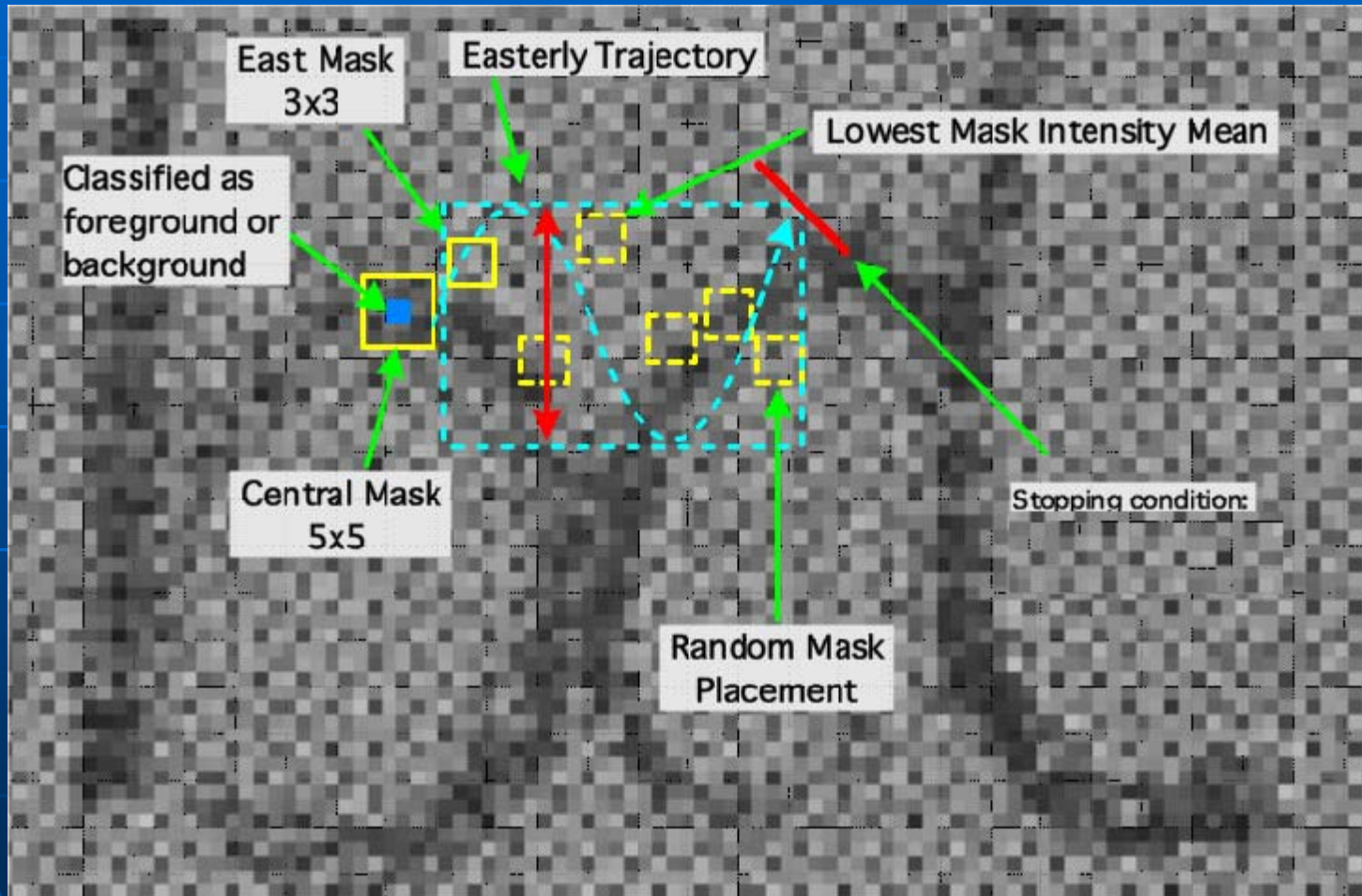
Let $\omega_i^* = \min_j \omega_{I,j}$ be the min of the means in the trajectory of ω_i

Step 4: If majority (of the 8 masks) satisfy $M(\omega_i^*) - M(v) \geq \kappa (=10)$
label **central pixel** as foreground

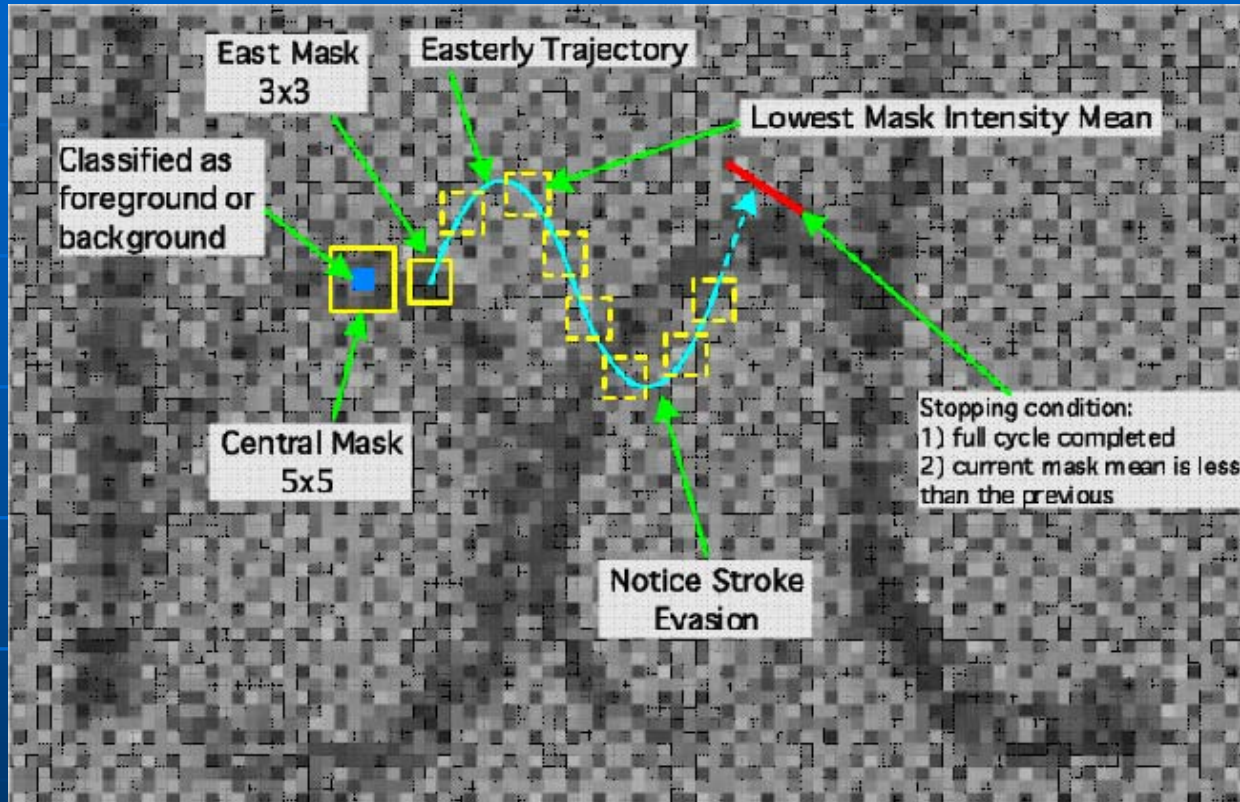
Linear Scouting



Random Scouting

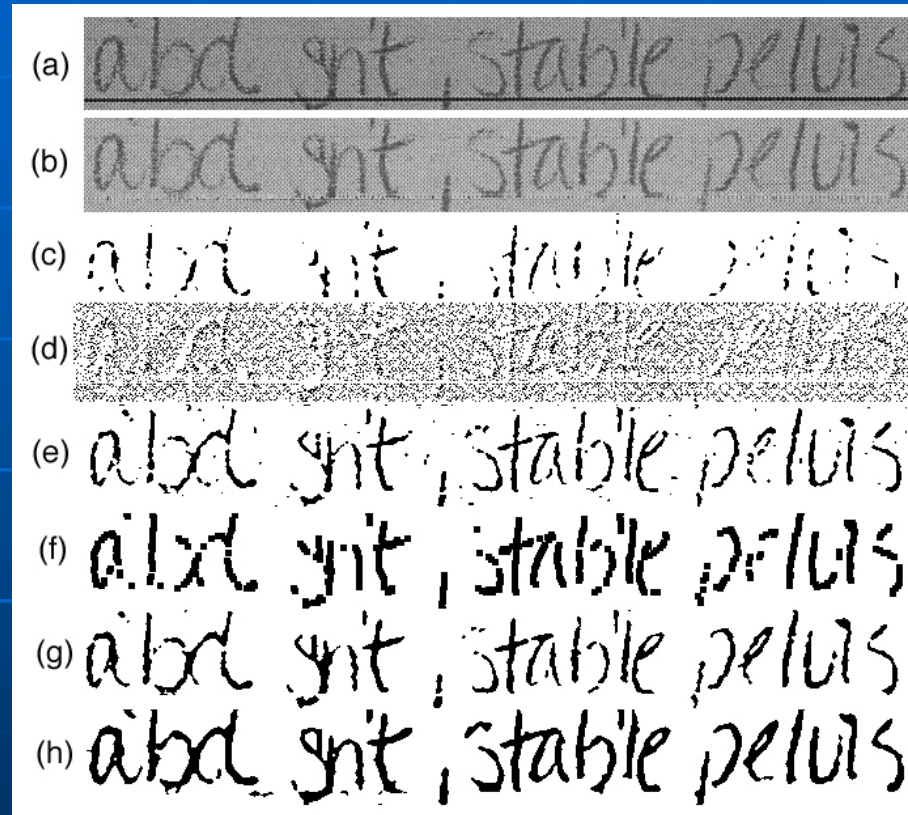
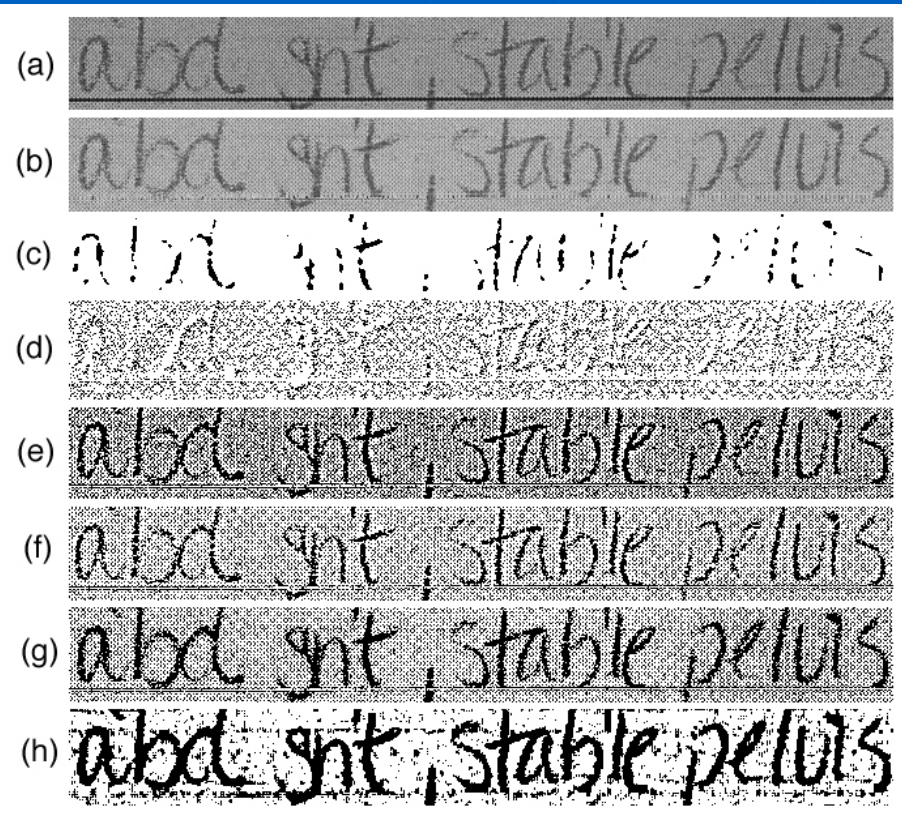


Compound Mask



- $y = 2\phi \sin(1/2x)$:
 - ϕ : estimated stroke width in pixels ($\phi=5$)
 - 2ϕ : to ensure sampling from both sides of a stroke
- Implement as a static compound mask

Binarization + Postprocessing



W
K
N
S
O
Z

a) Original

b) Form Dropout

c) (W) Wu/Manmatha

d) (K) Kamel/Zhao

e) (N) Niblack

f) (S) Sauvola

g) (O) Otsu

h) (Z) Our Method

Postprocessing: Median filtering and noise removal

Experimental Results

Test set: 30 PCR's; 1,440 word images

Binarization

| LEX | W | K | N | S | O | Z |
|-----|------|------|------|-------|-------|-------|
| 100 | 3.1% | 4.6% | < 1% | 19.4% | 35.3% | 56.5% |
| 1K | < 1% | 1.5% | 0.0% | 10.1% | 17.1% | 26.9% |
| 4K | 0.0% | 0.0% | 0.0% | 4.4% | 11.2% | 20.3% |

With Post Processing

| LEX | W | K | N | S | O | Z |
|-----|------|------|-------|-------|-------|-------|
| 100 | 3.1% | 4.6% | 46.9% | 48.3% | 52.6% | 59.1% |
| 1K | < 1% | 1.5% | 26.0% | 26.4% | 36.1% | 40.1% |
| 4K | 0.0% | 0.0% | 17.8% | 12.7% | 21.2% | 25.3% |

(W) Wu/Manmatha
(N) Niblack
(O) Otsu

(K) Kamel/Zhao
(S) Sauvola
(Z) Our Method