# CLEAT:

## A <u>CL</u>assification, <u>E</u>nhancement and <u>A</u>nalysis <u>T</u>oolkit
## for
## Heterogeneous Document Image Collections

*David Doermann*
*University of Maryland*
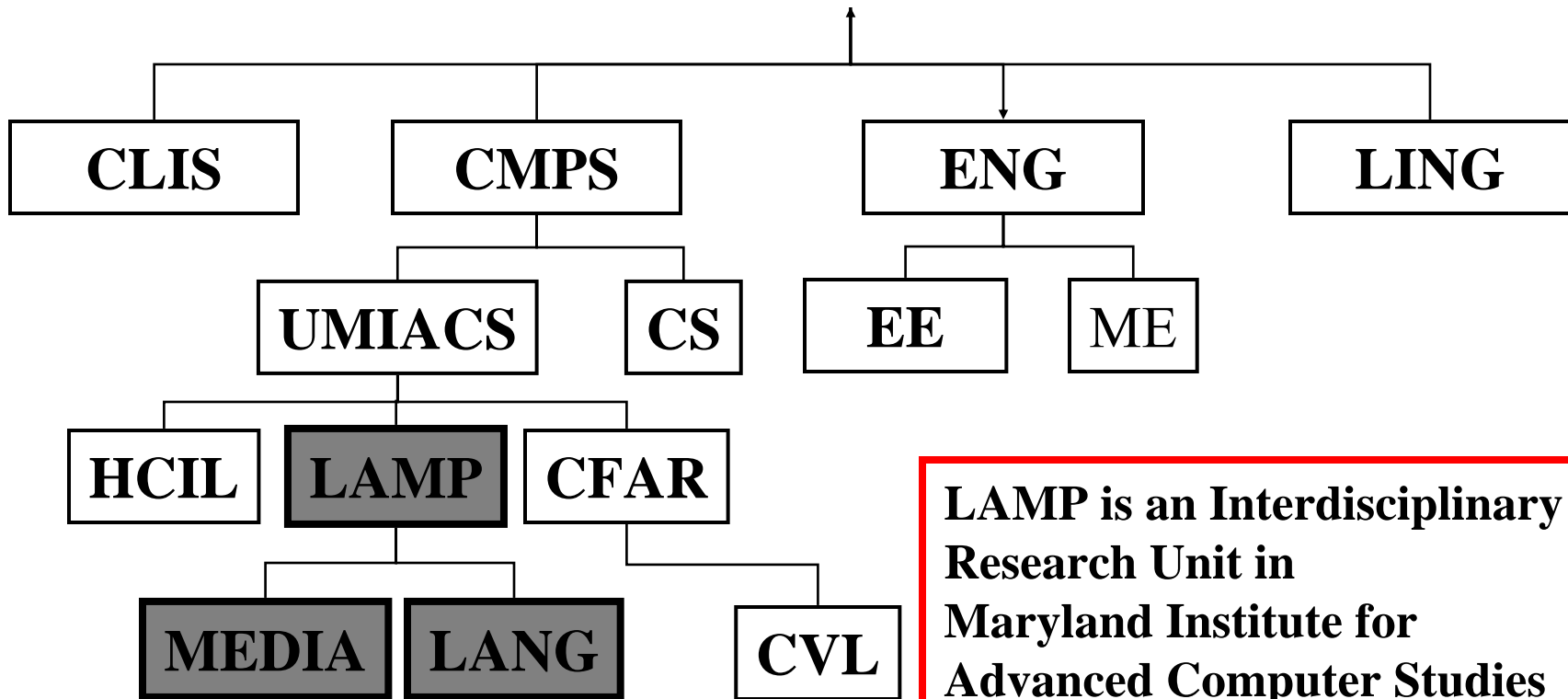
# Who are we?
# LAMP History

- Began in 1996 with a focus on documents
- Produced 9 PhD (2 more expected in 2007)
- Over 200 scientific publications
- Almost 50 Students (Undergrad-Graduate)
- Numerous Technology Transfer Opportunities

# Mission

To conduct research and education in  analysis and processing of multimedia information sources including documents, images and video, to develop natural language tools for real world applications, and to foster collaboration in these areas between researchers at the university and representatives of government agencies and industry

CLIS

CMPS

ENG

LING

UMIACS

CS

EE

ME

HCIL

LAMP

CFAR

MEDIA

LANG

CVL

LAMP is an Interdisciplinary Research Unit in Maryland Institute for Advanced Computer Studies (UMIACS)

# Outreach

- Bi-Annual SDIUT Conference
  - Soon to be included in Google Books Project
- Host of workshops and short courses
- Editorial Office of IJDAR
- Data Collection and Evaluations
- LAMP Seminar Series
- Chairing Program Committee for ICDAR 2007
- Organizing Arabic OCR competition at ICDAR'07

# Research Focal Areas

- Document image analysis
  - Providing fundamental tools for the enhancement summarization, navigation, indexing and retrieval in document image databases

- Content based video analysis
  - Providing access to video content through extraction, structure representation, classification, visualization and indexing

- In General
  - Ability to access large heterogeneous collections of material
  - Adaptable systems – OCR, MT
  - Low density to resource poor languages
  - Enhancing low quality input – document images, OCR

# Intelligence Value Estimation

- How can we take **large, noisy, unstructured, heterogeneous** collections of image and video data to:
  - Mine the nuggets?
  - Bubble the important things to the top?
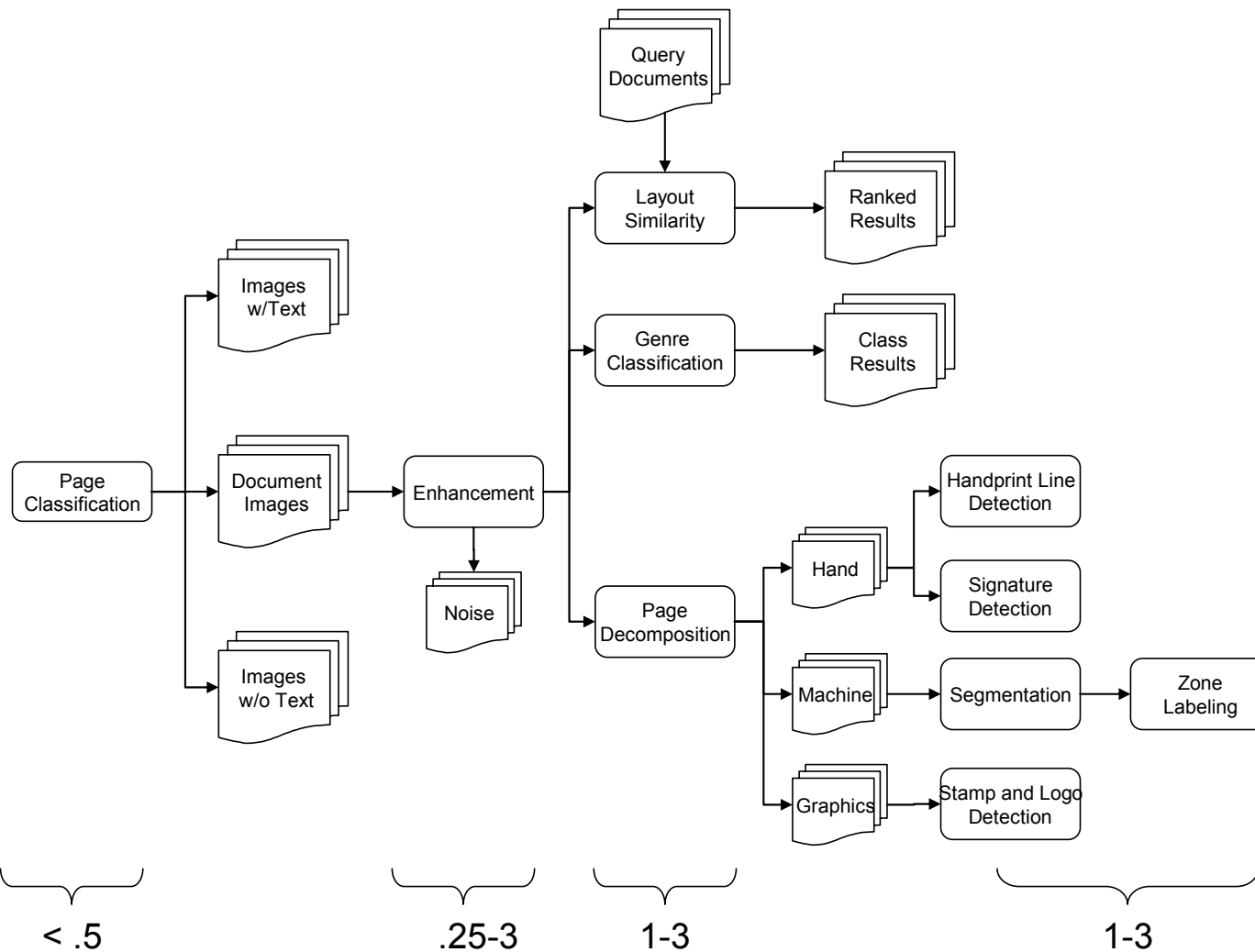  - Provide tools for Information Discovery?

# Challenges We Face….

- An overwhelming number of documents
  - Only a small fraction will ever be seen
- Huge variations in types, qualities and "value"
- Documents value diminishes with time
- We need to bring relevant documents to the top of the stack

# Approach

- Build robustness to noise into algorithms
  - Train noise as its own class
  - Integration of recognition and segmentation

- Provide mid level tools to organize collections
  - Genre Classification
  - Logo, Stamp and Signature Detection/Recognition

- Focus on Ranking rather then "conversion"
  - Page Layout Similarity

- Provide tools necessary for efficient research and evaluation
  - Datasets
  - GEDI – Groundtruth and Evaluation

# Project Overview



Target Processing Speed in Seconds

< .5          .25-3    1-3              1-3

# Task Objectives

Task 1:     Data Collection
Task 2:     Ground Truthing
Task 3:     Evaluation Framework
Task 4:     Evaluation and Visualization Tool

Task 5:     Page Classification Module
Task 6:     Enhancement Module
Task 7:     Layout Analysis Module
Task 8:     Content Labeling module

Task 9:     Evaluation
Task 10:    Training

# Performance Goals

| Task | Performance Goal |
|------|------------------|
| Page Classification | 80% precision across all three classes |
| Enhancement | 10-30% increase in accuracy of downstream processes – segmentation, detection |
| Layer Separation | 90% coverage at the pixel level |
| Segmentation (Print and Hand) | 85% using implementation of existing methods |
| Logo and Stamp Detection | 75% precision at 85% recall |
| Signature Detection | 75% precision at 85% recall |

# GEDI – Java Interface

# Data Collection and Evaluation

| Type | Number |
|---|---|
| Class 1: Traditional Document Images | 9000 |
| Class 2: Camera captured, Text in Scene, and Color documents | 500 |
| Class 3: Non-document Images | 500 |
| **Genre** | **Number** |
| Forms, Drawing, Tables | 1000 |
| Business Documents, Memos, Letters | 2500 |
| Journal and Conference Papers, Articles | 2500 |
| Newsletters, Flyers | 1000 |
| Structured Documents – phone books, dictionaries | 1000 |
| Handwritten | 1000 |
| Foreign Language – handwritten and machine printed | 1000 |
| Highly Degraded | 500 |
| Mixed Annotation | 2000 |

15



4



German - Basic

Bim, Bam, Bum - Ein Glockenton

fliegt durch die Nacht, als

hätt' er Vogelflügel, er fliegt

in römischer Kirchentracht

wohl über Tal und Hügel. Er



German - Digits/Symbols, Uppercase

**HABEN SIE WEITER GEDRÜCKT?**

2900 + 458 = 3131 - 227?

{498, 814, 4687}

(371) 422-2878: 80011

*897* [943] [3882]

52,9%! $762 #70:00



FFX 115c

MELLO ... TO RUN FOR MAYOR OF FREMONT

Fremont Councilman Gary Mello closed the door Thursday on a mayoral bid, increasing the chances that Mayor Bill Ball will be re-elected to a second term in November.

Mello announced in March that he was considering a possible run for mayor, but has said subsequently on several occasions he wouldn't rule out the would challenge Ball.

"After a long and difficult decision-making process, I have decided not to run for mayor of Fremont, at this time," Mello said in a written statement.

Mello's council seat does not expire until 1993.

In a telephone interview, the 34-year old little insurance company executive said he did not want to devote more time to city business at the expense of his family and job. He currently spends about 30 hours a week on city-related business, and that being mayor would mean at least 10 hours per week.



FFX 121J

GENEROSITY FINANCES MERCY MISSION TO FARMWORKERS

A silent disaster has stricken the people of California's Central Valley, and two Fremont city Hall staffers are spearheading a drive to help.

The situation is dark.

"People are dying or close to death," warned Fremont Mayor Bill Ball.

# New Data

- 25,000 pages ground truthed to the zone level

- Sampled from the Tobacco Litigation Corpus of 49 Million pages

# 25,000 pages ground truthed

|  | DOCS | PAGES |  | DOCS | PAGES |
|---|---|---|---|---|---|
| dt_calendar | 44 | 90 | dt_email | 973 | 1151 |
| dt_photograph | 227 | 461 | co_tables | 1049 | 1980 |
| dt_questionnaire | 188 | 461 | dt_form | 1582 | 2265 |
| dt_bibliography | 175 | 530 | co_foreign | 1669 | 2300 |
| dt_periodical | 479 | 693 | dt_notes | 2288 | 2925 |
| dt_list | 405 | 710 | co_illegible | 2598 | 3983 |
| dt_advertisement | 519 | 894 | dt_graphic | 2061 | 4307 |
| dt_newspaper | 688 | 921 | dt_letter | 3145 | 4601 |
| co_fax | 830 | 1150 | dt_report | 2213 | 4604 |
| co_drawings | 638 | 1150 | dt_memo | 2762 | 4611 |
|  |  |  | co_handwritten | 4894 | 6903 |
|  |  |  | co_marginalia | 10665 | 17251 |

# DocLib Architecture

- **Efficient Technology Transfer**
  - software compatibility
  - balance of academia, governemnt, and industry needs
  - common framework for document processing

- **Scalability**
  - rapid prototyping of new methods
  - simple algorithm comparison

- **Robustness and Stability**
  - high quality standards
  - platform-independence
  - accommodation of frequently changing requirements

# DocLib Status

- Core DocLib components matured and stable (in use by a variety of government installations)\

- Addons being integrated/implemented, primarily by developers

- Freely available to government researchers

- Core supported on Solaris, Linux and Windows

# Core vs Add-ons

- Core components are loosely defined as necessary building blocks for ANY document analysis process

- Addons are tools and applications for specific types of analysis

*We try to put as few constraints on the representations as possible.*

# Image Factory



**Design Factors:**

- ➤ **Image Type objects are static/singleton objects created on startup**
- ➤ **DLImageFactory is a static/singleton object**
- ➤ **Image Type objects registers itself with the DLImageFactory during startup**
- ➤ **DLImageFactory keeps a list of supported Image objects as each image type calls the register function**
- ➤ **Additional image types can be plugged into DOCLIB without modifying existing DOCLIB code.**

# DocLib Architecture

**DocLib´s architecture rests on two pillars:**

**DLImage:**
➢ **Image Processing**

**DLDocument:**
➢ **Document Processing**

DLImage ⟷ DLDoc

**e.g.**
➢ **image rotation**
➢ **image deskewing**
➢ **image conversions**
➢ **cc calculation**
➢ **shape drawing**

**e.g.**
➢ **page segmentation**
➢ **text line extraction**
➢ **logo detection**
➢ **XML input/output**
➢ **page layout analysis**

# Document Hierarchy

# Recent Modules

- Thinning
- Rotation
- Deskewing

- XML i/o
- Degradation
- OCR Scansoft interface (Windows)
- Docstrum

- Logo detection
- Signature processing

- LogoDetect
- TokenMatch
- Machine vs. Handwritten
- Jargon
- Text Line Detection

# XML Output Extension

```xml
<?xml version="1.0" encoding="UTF-8" ?>
<!-- GEDI is developed at Language and Media Processing Laboratory,
     University of Maryland.   -->
<GEDI xmlns="http://lamp.cfar.umd.edu/GEDI" version="1.0">
  <USER name="Elena" date="Sun, 14 Oct 2007 8:28 PM" />
  <DL_DOCUMENT src="aaa27e00.tif" docTag="xml" NrOfPages="2">
    <DL_PAGE gedi_type="DL_PAGE" src="aaa27e00.tif" pageID="1«
        width="2560" height="3296">
      <DL_ZONE gedi_type="STAMP" id="None" col="1174" row="495"
          width="447" height="132" />
      <DL_ZONE gedi_type="LOGO" id="None" col="274" row="569"
          width="346" height="159" contents="" />
      <DL_ZONE gedi_type="MACHINEPRINT" id="None" col="647"
          row="626" width="1372" height="105" contents="" />
      <DL_ZONE gedi_type="MACHINEPRINT" id="None" col="2410"
          row="2479" width="511" height="110" orientation="-
1.6295521495106193" contents="" />
    </DL_PAGE>
  </DL_DOCUMENT>
```

# Technical Presentations

- Page Segmentation (and rule line separation)

- Logo Detection and Recognition
- Signature Detection
- Stamp Detection

- Document ID/Script ID
- Page Layout Similarity

- Video Research
  - Tracking and Analysis of People
  - Video Content Classification

# Page Layer Segmentation
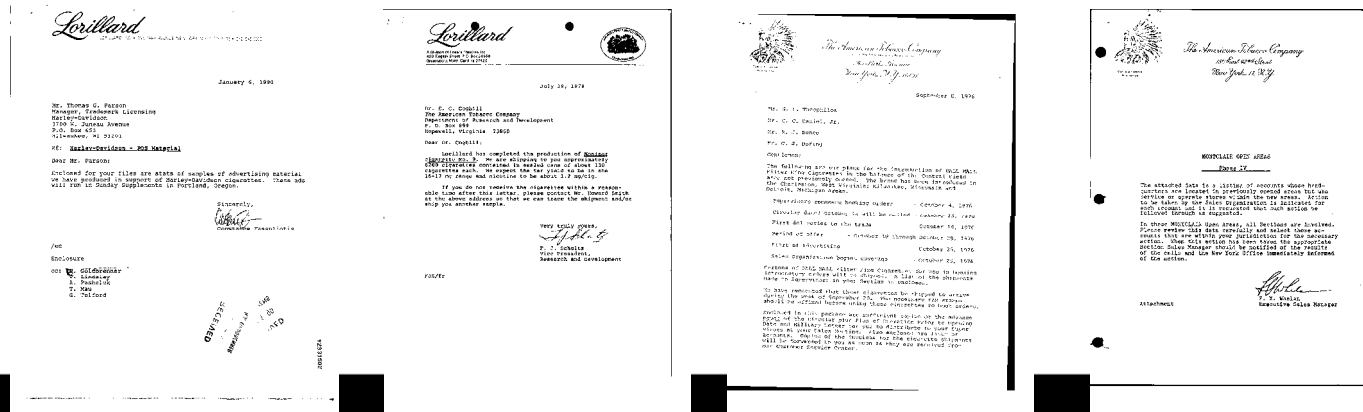
- Document image generation model
  - A document consists many layers, such as handwriting, machine printed text, background patterns, tables, figures, noise, etc.

# Motivation

- Document analysis has been viewed as a solved problem in clean, well-constrained documents.

- However, the performance degrades significantly when a small amount of noise is introduced.

- We further separate handwriting from machine printed text.

# Page Segmentation for Noisy Documents



\* Docstrum page segmentation technique is used

# Overview of Our Approach

- Segment the document to word level using connected component based, bottom-up approach.

- Classify each segmented block into noise, handwriting or printed text, based on extracted features and the Fisher classifier.

- Using MRF (Markov Random Field) to refine the classification result.

# Feature Extraction and Selection

- We extracted 140 features and 31 of them are selected to train the

| | Usage description | Dimensio | Selected |
|---|---|---|---|
| Structural | Region size, connected components | 18 | 9 |
| Gabor filter | Stroke orientation | 16 | 4 |
| Run-length histogram | Stroke length | 20 | 5 |
| Crossing counts histogram | Stroke complexity | 10 | 6 |
| Co-occurrence | Texture | 16 | 2 |
| $2 \times 2$ gram | Texture | 60 | 5 |
| Total | | 140 | 31 |

# Classification Results with Fisher Classifier

Printed text
Handwriting
Noise

# Using Context

- The results are reasonable with a few mis-classification due to the overlapping of different classes in the feature space.

- Context can be used to refine classification results further
  - Words of printed text tend to lie on the same line.
  - Noise block are likely to overlap each other.

- This kind of local dependency among neighboring components can be described with the Markov Random Field (MRF).

# Clique Definition

- Low level MRF is defined on regular lattice (pixel)

- Our high level MRF is defined on a graph.
  - After defining the connection between word blocks, a graph is generated.
  - Neighborhood of MRF is defined on the graph.

- Clique $C_p$ for printed text

| Left | Center | Right |
|------|--------|-------|

$O_v$  $\rightarrow D_h \leftarrow\!-$

- Clique $C_v$ for Noise

# MRF Postprocessing Example

Printed text
Handwriting
Noise



Before MRF-based postprocessing

After MRF-based postprocessing

# Evaluation

- Data Collection
  - 318 documents provided by the tobacco industry.
  - 94 documents of testing, the other for training.

| | #Total | Percentage | Before Post-processing | | After Post-processing | |
|---|---|---|---|---|---|---|
| | | | Accuracy | Precision | Accuracy | Precision |
| Printed Words | 19,227 | 66.9% | 95.9% | 99.5% | 98.0% | 99.7% |
| Handwritten Words | 701 | 2.4% | 93.2% | 62.9% | 93.0% | 83.3% |
| Noise Blocks | 8,802 | 30.7% | 96.8% | 93.0% | 98.6% | 96.0% |
| Total | 28,730 | 100% | 96.1% | N/A | 98.1% | N/A |

# Application to Page Segmentation



Before enhancement

After enhancement

# Rule Line Detection Example



**Model-based line detection result**

**After rule line removal**

# Evaluation

- Database
    - 168 Arabic documents with a total of 3,870 groundtruthed lines.
    - 100 images for the training of the HMM model, 68 images for the testing.
- Quantitative evaluation (evaluation metrics are discussed in the paper in detail).

QUANTITATIVE EVALUATION OF THE RULE LINE DETECTION RESULT.

| | Groundtruthed Lines | Detected Lines | Correct | Partial Correct | Missed | False Alarm |
|---|---|---|---|---|---|---|
| Training Set | 2,274 | 2,319 | 2,212 (97.3%) | 56 (2.5%) | 6 (0.3%) | 51 (2.2%) |
| Test Set | 1,596 | 1,631 | 1,545 (96.8%) | 49 (3.0%) | 2 (0.1%) | 37 (2.3%) |

# Technical Presentations

- Page Segmentation (and rule line separation)

- Signature Detection
- Logo Detection and Recognition
- Stamp Detection
- Document ID/Script ID

**Metadata Extraction**

- Page Layout Similarity

- Video Research
  – Tracking and Analysis of People
  – Video Content Classification

# Problem Statement

Given a large heterogeneous document image database, we are facing a few very challenging problems

- How can we retrieve documents authored or approved by a specific individual in unconstrained settings?



- How can we retrieve documents originating from an organization?

# Motivation

- Signatures and logos provide exciting new dimensions for document image mining

- Solution to these problems are also important in document analysis systems in a range of application domains

  - Signature verification and identification

  - Business process automation

# Our Tasks

- Two problems are of fundamental interest to general content-based image retrieval

  – Detection and segmentation

  – Matching

    • Representation

    • Similarity measures

    • Matching algorithms

| Signature Detection & Segmentation | Extract Pointset & Compute Shape Descriptor | Compute Shape Distance |
|---|---|---|
| | Solve for Point Correspondences | |
| | Solve for Non-rigid Transformations | |

# Overview of our approach

- We treat a signature as a global symbol. Rather than focusing on local features that typically have large variations, our approach aims to capture the structural saliency of a signature by searching over multiple scales

- We consider identifying salient structure and grouping its parts in two separate steps

- Two keys questions we addressed are:

  - How to effectively model off-line signature production under reasonable assumptions without its temporal information

  - What to effectively measure the structural saliency of signatures under such production model

# Evaluation

- We used two large collections of real world documents—Tobacco-800 and University of Maryland Arabic datasets.

- Using document context, our multi-scale signature detector achieves 92.8% and 86.6% detection rates for the Tobacco-800 and Maryland Arabic datasets, at 0.3 false-positives per image.



(a)                    (b)

ROC curves for (a) Tobacco-800 dataset and (b) Maryland Arabic dataset.

# Evaluation



Examples of detected signatures from Tobacco-800 and their saliency maps.

# Evaluation



Examples of detected signatures from Maryland Arabic dataset and their saliency maps.

# Evaluation



(a)

(b)

Examples of (a) falsely alarms (b) missed signatures

# Shape representation



(a)  (b)  (c)

(d)  (e)  (f)

Shape contexts [Belongie *et al.*, 2002]  and local-neighborhood-graph [Zheng and Doermann, 2006] constructed from detected and segmented signatures.

# Shape matching



(a)          (b)

(d)          (e)

(g)                                    (h)

Illustration of signature matching using shape contexts and local-neighborhood-graph

# Shape matching



Illustration of signature matching using shape contexts and local-neighborhood-graph

# Shape matching evaluation

A query with eight relevant signature instances



Top eight retrieved in the ranked list



(1)    (2)    (3)    (4)

(5)    (6)    (7)    (8)

Relevant instance outside the top eight in the ranked list



(10)

A signature query example. Among the total of eight relevant signature instances, seven appear in the top eight of the 460-element ranked list, giving an average precision of 94.2%, and an R-Precision of 87.5%. The irrelevant signature that is ranked among the top eight is highlighted with a blue bounding box.

# Signature matching results

**Table 1: Signature retrieval result using different similarity measures.**

| Similarity measures | Mean average precision | Mean R-precision |
|---|---|---|
| $D_{sc}$ | 66.9% | 62.8% |
| $D_{af}$ | 61.3% | 57.0% |
| $D_{be}$ | 59.8% | 55.6% |
| $D_{re}$ | 52.5% | 48.3% |
| $D_{sc} + D_{be}$ | 78.7% | 74.3% |
| $D_{sc} + D_{af} + D_{sc} + D_{re}$ | 84.5% | 80.8% |

**Table 2: Signature retrieval result using multiple instances of signatures from the same person in each query.**

| Number of instances | Mean average precision | Mean R-precision |
|---|---|---|
| One | 84.5% | 80.8% |
| Two | 88.6% | 85.2% |
| Three | 91.3% | 88.1% |

# Logo Detection and Recognition

– enables identification of the source of documents from a given organization

– Most studies assume good logo detection and segmentation is available

- Challenges

– Detection is required for any prior to extraction

– Extraction is required for any shape based matching/recognition  process

# Challenges

- Extremely large intra-class variations among logos
- Continuum between graphics, logos and text

# Challenges

- Diverse document layouts, scanning qualities, and image degradations on real document datasets

# Claim #1

- Documents exists where *spatial* segmentation of Logos, Signatures and Stamps is not an option!

# Claim #2

- Considering the more general problem of **Detection** *(as opposed to segmentation->classification)* allows us to integrate identification and extraction, and possibly recognition

- The concept has successfully been applied to:

  – *Guangyu Zhu, Yefeng Zheng, David Doermann and Stefan Jaeger. Multi-scale Structural Saliency for **Signature Detection**. (CVPR 2007).*

  – *Guangyu Zhu, Stefan Jaeger and David Doermann. A Robust **Stamp Detection** Framework on Degraded Documents. SPIE 2006.*

# Multiscale Detection

- Each logo candidate region is further classified at successively finer image scales by a cascade of simple classifiers



- The overall classifier is a strong learner, even if each individual classifier is in fact a weak learner

# Feature selection and extraction

- How can we explore document context for logo detection?



**Clustering result of logo positions using *k*-means (*k* = 3)**

$$D_c(P) = \min_{i \in \{1,2,\cdots,k\}} (|p_x - c_x^{\cdot}| + \lambda|p_y - c_y^{\cdot}|)$$

- We define context distance as

| Context Distance | Area | Symmetry |
|---|---|---|
| Spatial Density | Aspect Ratio | Text Uniformity |

# Evaluation

- We use tobacco-800, a large public dataset that consists of 1290 real-world documents (full dataset 49 million pages)

- Use accuracy and precision as evaluation metrics

$$\text{Accuracy} = \frac{\text{\# of correctly detected logos}}{\text{\# of logos in groundtruth}} \qquad \text{Precision} = \frac{\text{\# of correctly detected logos}}{\text{\# of detected logos}}$$

- Detection is at least > 75% and < 125% pixel are overlap (determined from shape matching approach – Zhang. PAMI 2006)

**Summary of logo detection performance on the Tobacco-800 dataset**

|  | Accuracy | Precision |
|---|---|---|
| Improved spatial density [9] | 39.3% | 32.1% |
| Fisher classifier only, *i.e.*, \|S\| = 1 | 59.2% | 41.7% |
| Multi-scale approach with \|S\| = 2 | 57.0% | 68.1% |
| Multi-scale approach with \|S\| = 3 | 84.2% | 73.5% |

# Evaluation



**Examples of correctly detected logos from Tobacco-800**

# Evaluation



**(a) Over/under-segmented logos**



**(b) Non logos**

Examples of incorrectly detected logos



Examples of missed logos

# Challenges in stamp detection

- Unique characteristics of stamps
    - Unstable and unpredictable patterns in documents
    - Outliers and occlusions are typical
    - Much lower spatial density compared to logo
    - Stamp instances appear as weaker regions within a full spectrum of background – text, figures, tables, watermark
    - Not generally valid to assume its location within the source

# Our stamp detection approach

**Image Scaling & Conversion**

**Gaussian Smoothing**

**Extract Edge Strength**

**Extract Edge Orientation**

**Remove Junctions**

**Link Connected Edges**

**Filter Connected Edges**

**Select Edge Pairs**

**Parameter Estimation on $\{x_o, y_o, area\}$ by Voting**

**Verification**

Edge Extraction

Construct and effectively constraint the feature Space

Obtain stamp parameters

# Ellipse detection method using pairs of edges

$t_1(x, y) = 0$

$E_1$

$l(x, y) = 0$

$E_2$

$t_2(x, y) = 0$

parabola

circle

ellipse

hyperbola

The quadratic function $f(x, y)$ represents the family of 2nd-order curves that pass points $E_1$ and $E_2$ and tangent to lines $t_1(x, y)$ and $t_2(x, y)$.

# Demo



Region of a sample image



Strength of edge gradient

# Demo



Strong edges



Orientation of edge gradient

# Demo



Top 10 candidates in the 3-D parameter space in ellipse center and area, i.e. ($x_o$, $y_o$, $area$)

(68, 238, 11313), score = [5485509]
(56, 202, 6464),   score = [501958]
(52, 226, 8080),   score = [431456]
(72, 206, 8080),   score = [352608]
(84, 266, 6464),   score = [278291]
(84, 210, 6464),   score = [260775]
(44, 222, 8080),   score = [247448]
(28, 270, 3232),   score = [241991]
(40, 202, 4848),   score = [224263]
(76, 230, 9696),   score = [215384]

# Demo

# Demo

# Demo

Capability to detect multiple stamp instances

# Demo

Capability to detect stamp instances in diverse backgrounds

# Software releases

- Signature detection and logo detection code are released as Doclib add-on modules

- Production test on 32,000+ documents

- Signature matching and logo matching expected

# Demo

# Experiment

| Test Databases | Total Images | Images with The Retrieved Stamp |
|---|---|---|
| Database 1 | 436 | 92 (Circular) |
| Database 2 | 193 | 68 (Elliptic) |
| Database 3 | 287 | 102 (Rectangular) |

# Script and ImageID

- ScriptID
  - Given a set of handwritten document images, identify the scripts.
  - Dataset: UMD handwritten dataset + Arabic dataset

- ImageID
  - Given an arbitrary image, identify that it is
    - document image
    - image with text
    - Image w/o text
  - Dataset: ~3700 images from Internet.

# The Observation

朱雀桥边野草花，
乌衣巷口夕阳斜。
旧时王谢堂前燕，
飞入寻常百姓家。

옷을 깔끔하게 차
| 들어 있는 것
어릴 것 같은데
입고 있었다.

स्वायत्त विनियोग
जो राष्ट्रीय आय
वर्तमान आर्थिक
न्धारित होता है।

# The Observation (con't)

- The relationship of connected edges could be used for description;

- The dominant descriptors for different scripts could be different;

- The statistics of the descriptors could be used for discriminating different scripts.

# The descriptor

- Fit edges to small lines

- Adjacent lines: encode the relative coordinates w.r.t pivot point.

  – C / Z shape

  – Y shape

Yu et al, Object Detection Using Shape Codebook, BMVC 2007

# The codebook for the descriptor

- The advantage of the codebook
  - Generic
  - Quantization -> fast

- generate the codebook
  - A large dataset
  - Extract descriptor
  - Cluster the descriptor

# The implementation

- Given a document image
  - Preprocessing
    - Binarize if necessary
    - Skeletonize
    - Clean the image using mathematical morphology.
  - Extract descriptors
    - Extract line segments
    - Compute shape descriptors
    - Quantize the shape descriptors and compute their histogram.
  - Train and classify

# Result

- Confusion matrix (experimental result, july 2007)

|  | Arabic | Chinese | Hindi | Korean |
|---|---|---|---|---|
| Arabic | 11 (74%) | 1 | 2 | 1 |
| Chinese | 0 | 10 (77%) | 0 | 3 |
| Hindi | 1 | 1 | 10 (83%) | 0 |
| Korean | 1 | 3 | 0 | 9 (70%) |

# Failed examples

Arabic

Chinese

# Failure example (Korean)

# Image ID

- Determine which class:
  - Text, Image w/text, or image
- Adopt different vision modules
  - For different categories we can adopt different strategy in computer vision
- Improve efficiency
  - Use the category as prior.
- Speedup OCR module in real world environment.

# The Challenge

- Images are arbitrary
  - Appearance model cannot be used for the classification.
  - We use the same shape descriptor because the code book is generic.

- Ambiguity
  - "images / text vs images", e.g., Coke can.
  - "doc vs images / text", e.g. "publication cover" usually has figures.

# Dataset for ImageID

- Collected form Internet, through search using different keywords

- Manual inspection, removal of duplicate images.

| Page Classification Datasets (Google Image) | |
|---|---|
| Document | 797 |
| Image with Text | 1695 |
| Non-Document | 1275 |
| Total | 3767 |

google_cd_cover_
0.tif

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_c...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_c...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_c...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_c...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_c...

google_apple_0.tif    google_apple_1.tif    google_apple_2.tif    google_apple_5.tif    google_apple_9.tif    google_apple_1...    google_apple_1...    google_apple

google_apple_1...    google_apple_1...    google_apple_1...    google_apple_2...    google_apple_2...    google_apple_2...    google_apple_2...    google_apple

google_apple_2...    google_apple_2...    google_apple_2...    google_apple_3...    google_apple_3...    google_apple_3...    google_apple_3...    google_apple

google_apple_4...    google_apple_4...    google_apple_4...    google_apple_4...    google_apple_4...    google_apple_4...    google_apple_4...    google_apple

google_apple_5...    google_apple_5...    google_apple_5...    google_apple_5...    google_apple_5...    google_apple_5...    google_apple_5...    google_apple

# LBP: Local Binary Pattern

- Define
  - Texture: Joint distribution of center pos $g_c$ given neighbor sampling $g_p$ (p=0,..., P-1)

$$T = t(g_c, g_0, ..., g_{P-1})$$

- Example



$(P=4,R=1.0)$  $(P=8,R=1.0)$  $(P=12,R=1.5)$  $(P=16,R=2.0)$  $(P=24,R=3.0)$

# The performance

- Confusion matrix

|  | Doc | Image w/ | Non doc |
|---|---|---|---|
| Doc | 0.8557 | 0.1340 | 0.0103 |
| Image w/ | 0.1725 | 0.6011 | 0.2264 |
| Non doc | 0.0444 | 0.1422 | 0.8133 |

# The Module

- Input
  - Training: an text file contains a list of training images.
  - Testing: a filename to an image.

- Output
  - Training: an SVM classifier (model.txt)
  - Testing: XML format  (JEDI readable) for corresponding input image.

- Performance
  - 700 seconds for 3000 images
  - Similar speed for every image
  - No exceptions and memory leaks

# Results – classified as documents

# Images w/ text

# Images

# Technical Presentations

- Page Segmentation (and rule line separation)

- Signature Detection
- Logo Detection and Recognition
- Stamp Detection
- Document ID/Script ID

- Page Layout Similarity

<span style="color:red">Document Ranking</span>

- Video Research
  - Tracking and Analysis of People
  - Video Content Classification

# Motivation

- *In a large collection of documents (forms, academic papers, handwritten letters, checks, receipts, etc.), most times people need to handle only those with some specific layout.*

- ***Drawback*** *of our previous system for document ranking based on layout : training is restarted from beginning each time a new layout comes*

- ***Reason***: *we do not give an explicit definition of layout,  the system learns no concept of layout, but image content.*

- ***Proposal***: Let the system itself figure out important dissimilarities for layout classification.

# Layout Examples



1C    2C    1r2C    3C    2c_asym

2c2c_asym    class1    class2    class3    class5

# Document Representation
## -- Building blocks

- Text lines extracted by TB library (endpoint coordinates, line orientations)

# Quadrilaterals from text line pairs

- A document := {Quadrilaterals}



- Merits:
  - Text line properties (length, orientation) are defined implicitly by their relative contribution to the quadrilateral shape

- Drawbacks:
  - $O(n) \rightarrow O(n^2)$

# Quadrilateral Shape Vector

- 5D shape vector



$L_1$, $L_2$: text lines

$L_4$, $L_5$: diagonals

$L_3$: midpoints connection line

- – Vector uniquely defines the quadrilateral shape
- – Text line correspondence guaranteed
- – Efficient clustering

- Document represented this way is translation and 180° rotation invariant

# Dictionary of Quadrilaterals

- We need to establish correspondences between quadrilaterals so that documents comparison can break down into quadrilateral comparison.

- Clustering in 5D space using range search, each quadrilateral cluster is regarded as a word in the dictionary

- Need a rich dictionary to avoid too many unknowns in a query

- From 101 documents, we built a dictionary with 976 words

# Score a query document

- Each document has a signature *S* like

| 1 | 0 | 0 | 1 | 0 | … … | 1 | 1 |
|---|---|---|---|---|------|---|---|

- Each layout class has a relaxed signature *RS* averaged from training samples. (consistency)

| 0.9 | 0.1 | 0.12 | 1 | 0.07 | … … | 0.875 |
|-----|-----|------|---|------|------|-------|

- Each classifier has a performance value *P* on validation set. (discriminativity)

| 0.75 | 0.8 | 0.66 | 0.55 | 0.7 | … … | 0.6 |
|------|-----|------|------|-----|------|-----|

- Score of a query against layout class i

$$\text{Score}_i = \sum_k F(S_k, RS_{i,k}) * P_k$$

$$C = \text{argmax}_i \ \text{Score}_i$$

# Evaluation Scheme

- **Mean Average Precision (MAP)**
  - $P_i = (\sum_{i \le j} P_j) / (\sum_{i \le j} 1)$

- **Average Relevance Rank (ARR)**
  - $I = (\sum (R_i - (N_t+1)/2)) / (N * N_t)$

    $R_i$ : rank of one wanted testing document.

    N : testing size

    $N_t$: wanted testing size

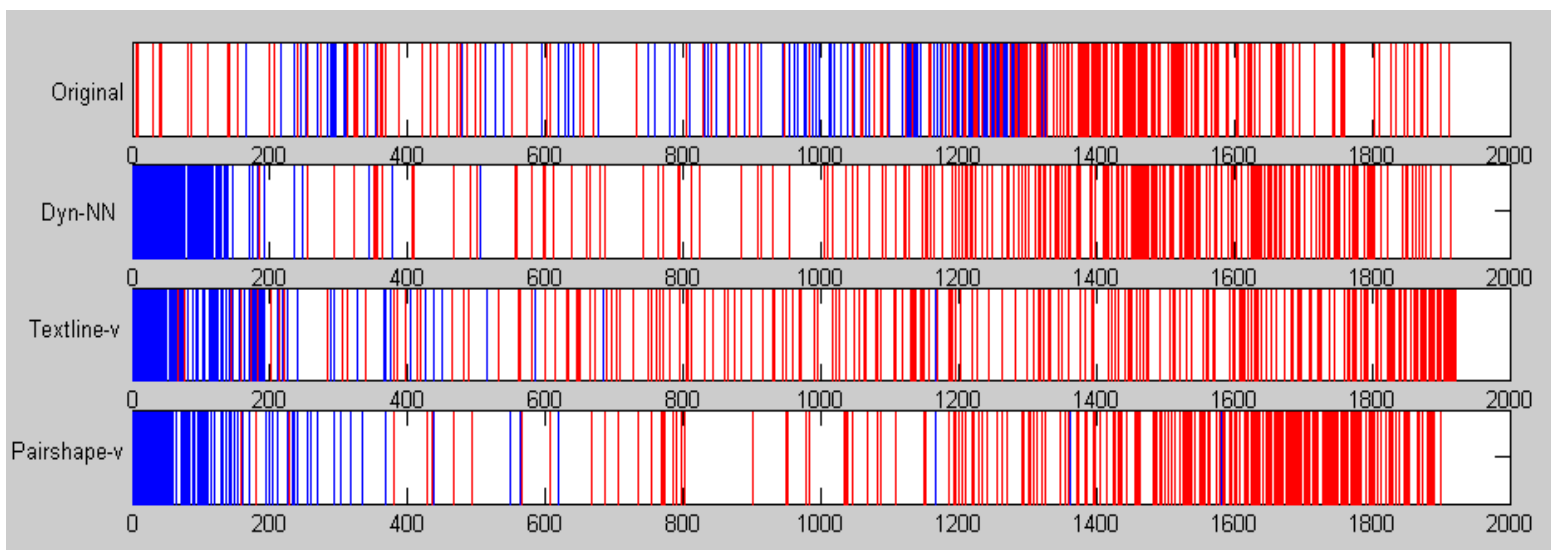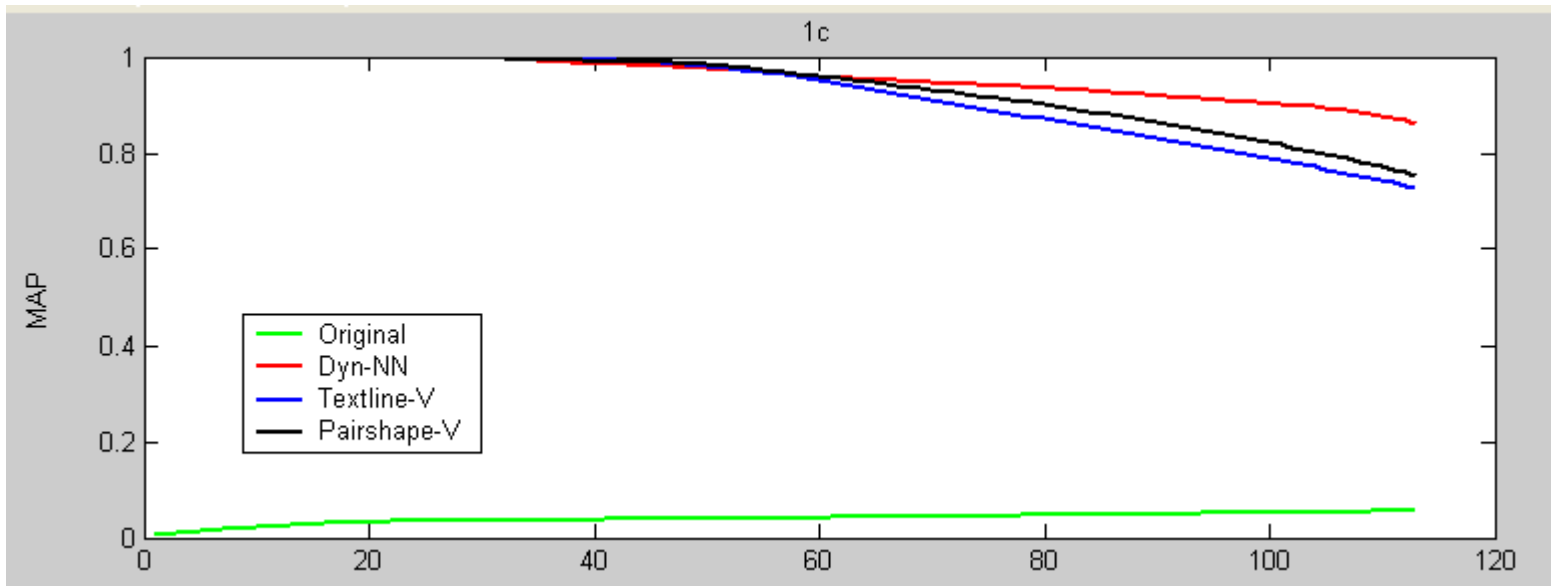  - $I \in [0, 1-N_t/N)$, the lower the better

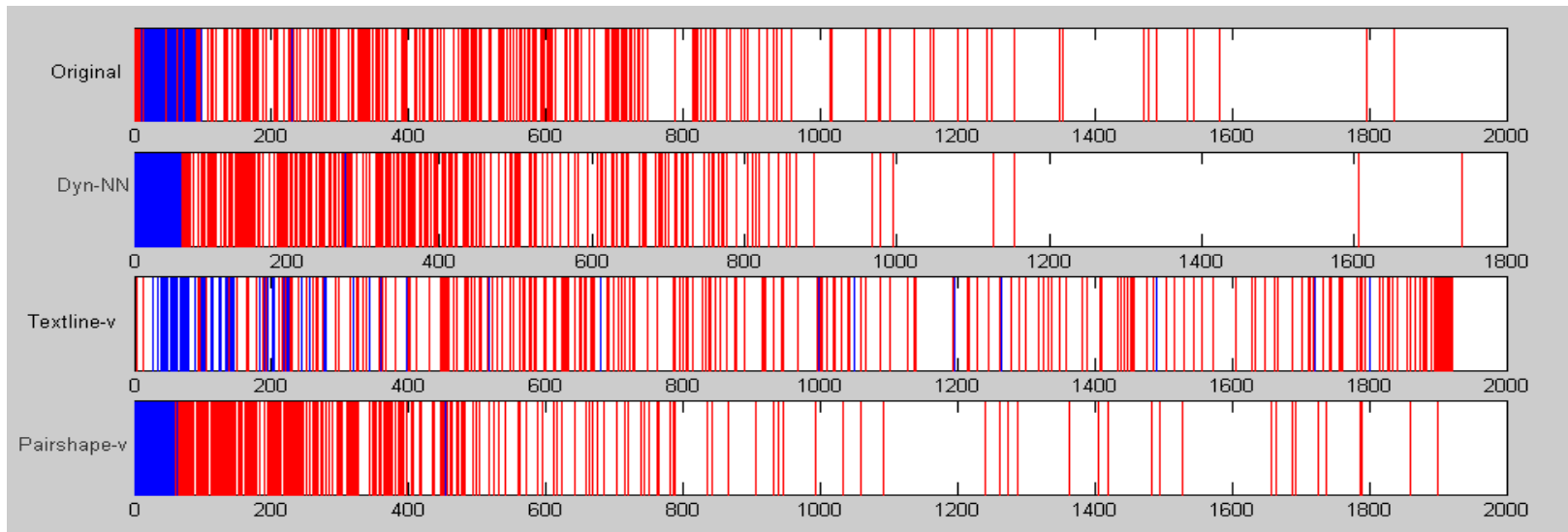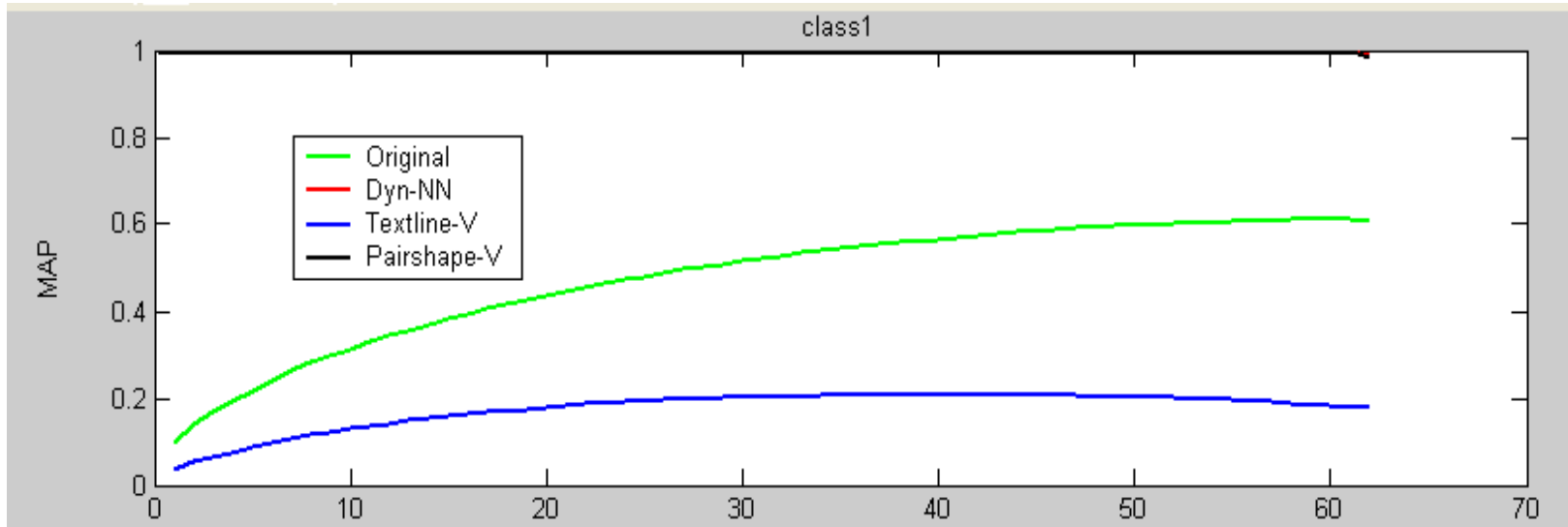# Experimental Results
## --Confusion Matrix

| | 1c | 2c | 1r2c | 3c | 2c_asym | 2c2c_asym | class 1 | class 2 | class3 | class4 |
|---|---|---|---|---|---|---|---|---|---|---|
| **1c** (113) | 87 | 8 | 16 | | 2 | | | | | |
| **2c** (144) | | 133 | 4 | 1 | | 5 | 1 | | | |
| **1r2c** (431) | 9 | 168 | 246 | | | 8 | | | | |
| **3c** (23) | | | | 23 | | | | | | |
| **2c_asym** (6) | | | | | 3 | 3 | | | | |
| **2c2c_asym** (45) | | 1 | | | | 44 | | | | |
| **Class1** (62) | | | | | | | 62 | | | |
| **Class2** (264) | 3 | | | | | 2 | 3 | 230 | 2 | 24 |
| **Class3** (121) | 1 | | | 1 | | | 13 | 2 | 101 | 3 |
| **Class4** | | | | 1 | | 1 | 17 | 27 | 7 | 52 |

# Experiments – ARR Results

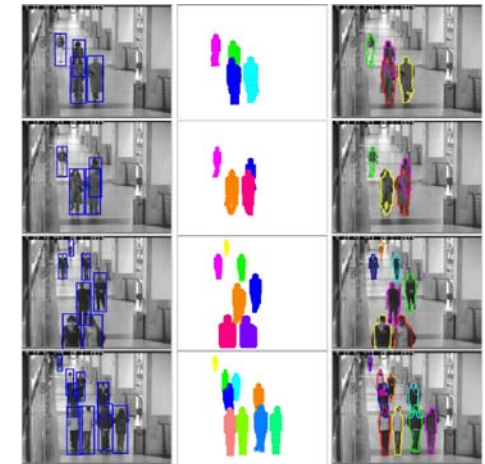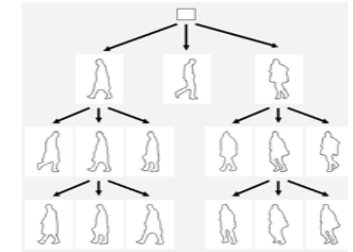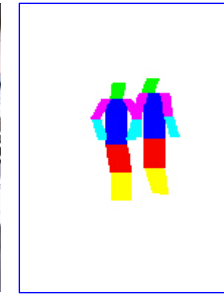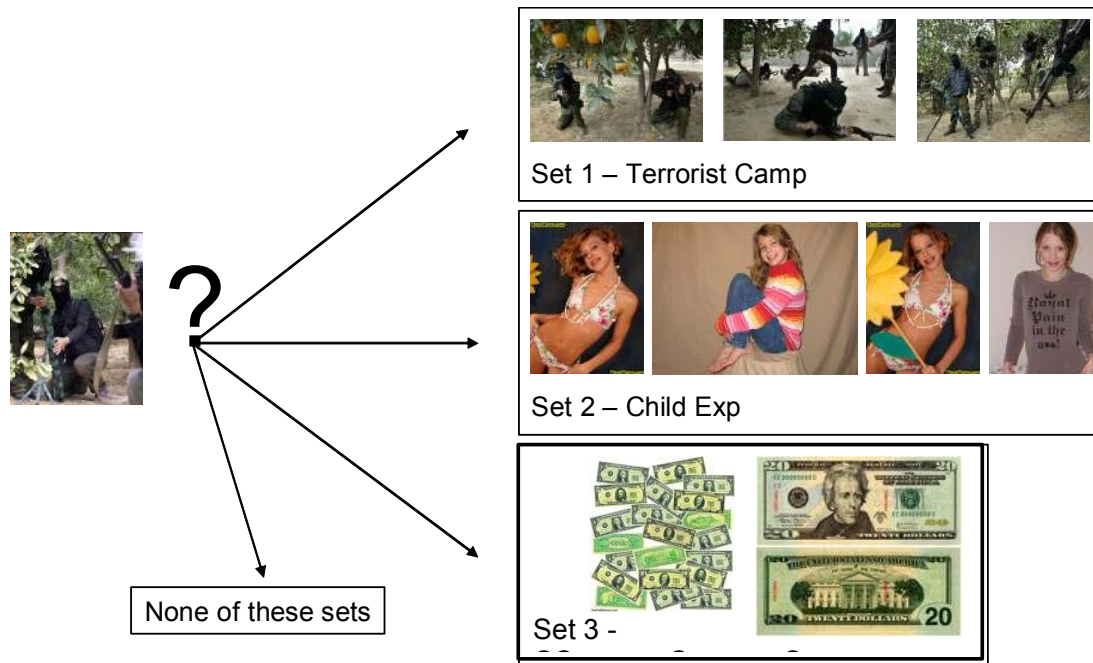|  | Original | Dyn-NN | Text-V | Pair_V |
|---|---|---|---|---|
| 1c | 0.450 | 0.011 | 0.038 | 0.043 |
| 2c | 0.062 | 0.010 | 0.324 | 0.087 |
| 3c | 0.028 | 0.0002 | 0.504 | 0.013 |
| 1r2c | 0.148 | 0.063 | 0.245 | 0.105 |
| 1r1r2c | 0.159 | 0.010 | 0.103 | 0.045 |
| 1r2c2c | 0.121 | 0.067 | 0.186 | 0.139 |
| 2c_asym | 0.137 | 0.025 | 0.360 | 0.039 |
| 2c2c_asym | 0.025 | 0.0002 | 0.097 | 0.010 |
| class1 | 0.009 | 0.002 | 0.133 | 0.003 |
| class2 | 0.398 | 0.011 | 0.004 | 0.075 |
| class3 | 0.160 | 0.026 | 0.146 | 0.090 |
| class5 | 0.302 | 0.056 | 0.103 | 0.085 |

class5

# Image and Video Research

- Surveillance Video
  - People Tracking
  - Appearance Modeling
  - Pose Estimation
- Partial Image Matching
  - Robust to changes in view point
  - Able to match partial images

# Forensic Image Search

- Consider a "search pack" which contains a "model" of a set of images of interest
- Hard Drive is searched and produces a report without revealing the search content



Set 1 – Terrorist Camp

Set 2 – Child Exp

Set 3 -

?

None of these sets
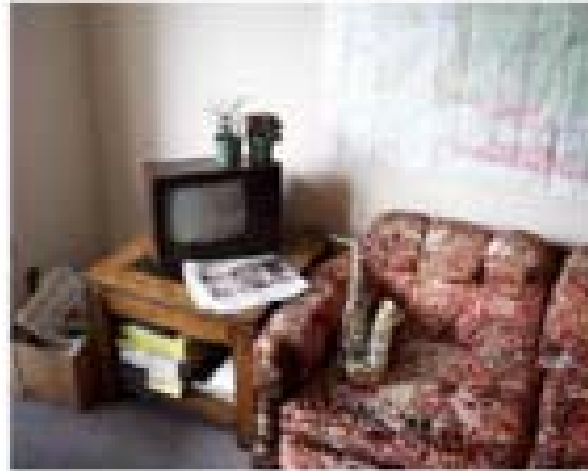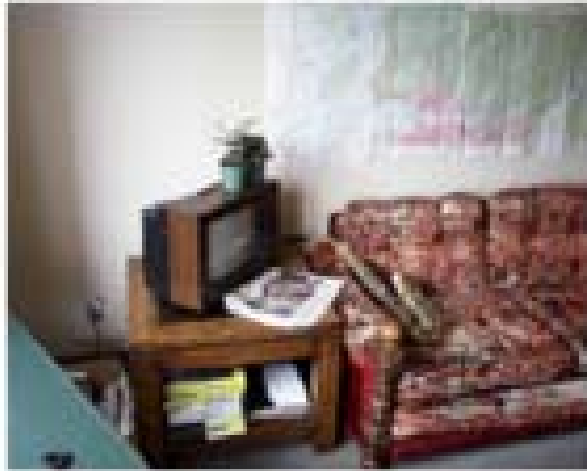
# High Speed Image Classification

- The purpose of this project
  - To create a new content based image retrieval (CBIR) algorithm that will remove some of limitations of the state of the art

- The task description
  - A user provides a set of training images belonging to several known categories (called SearchPak) and a set of test images.
  - For a test image, the user wants to know if it is similar to one of the SearchPak categories and otherwise classify it as "non-SearchPak image" or "junk image".

# Using what features?

- Histogram, correlogram of color, edge, texture…?
- A good feature: keypoint
  - A feature based on neighborhood edge histogram that is scale and rotation-invariant
  - Independent of color
  - Approach is called SIFT (*Scale Invariant Feature Transform*)
  - Captures salient visual information
- Groups of keypoints are powerful description of objects in images and video

# What Can we do with Keypoints?

- Searching (Video Google, Zisserman Oxford, K-means clustering)
- Mining (find most significant objects)
- Indexing (find anchor and cluster frames)
- Browsing
- Logo search
- Near-duplicate detection
- Face detection
- Building detection

# Summary

- Focused on Integration with DocLib framework
- Need Software engineering support
- Detailed evaluation and evaluation tools as part of Prototypes.

# Possible Research Extentions

- Increasing the speed of processing (software or hardware)

- Script independent word spotting

- Stamp and signature recognition

- Scene text recognition and super-resolution in video.

- Word level Script and Language ID