# CLEAT:

## A CLassification, Enhancement and Analysis Toolkit for Heterogeneous Document Image Collections
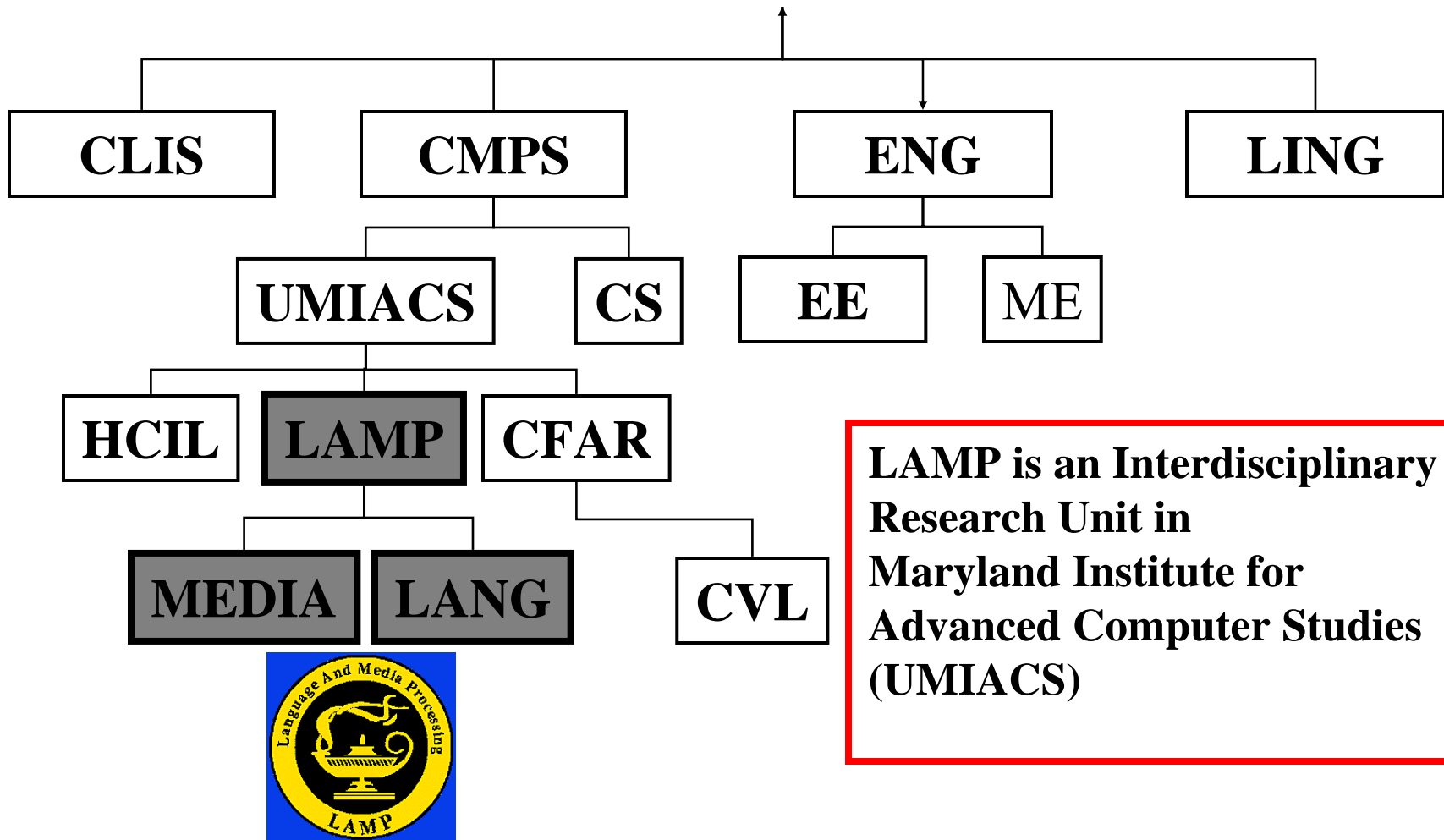
# LAMP History

- 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)

# 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

# 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
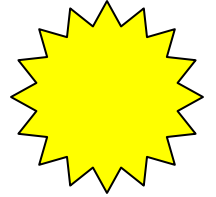- Organizing Arabic OCR competition at ICDAR'07

# Schedule

- Tuesday
  - AM: Project Overview and Status
  - PM: Logo Detection/Recognition

- Wednesday
  - AM: Font OCR, Word Level ScriptID
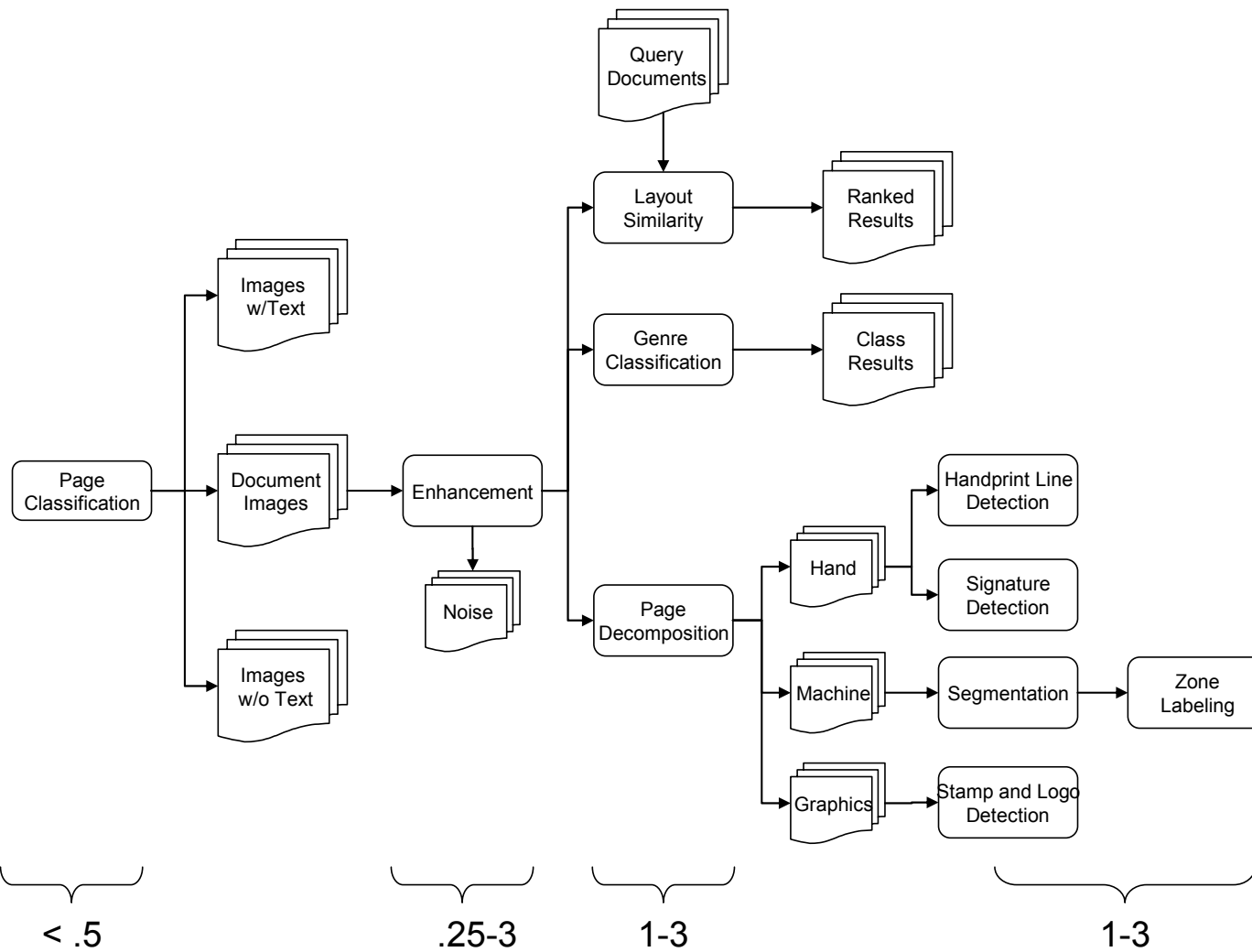
    Vision Related Research
  - PM: Review and Feedback

# Agenda
# Tuesday AM

- ## Project Overview
  - Introduction
  - Goals and Objectives

- ## Tools
  - GEDI Display Environment
  - Datasets

- ## DocLib and Algorithms
  - Technical Presentations

# Project Overview



Target Processing Speed in Seconds

# 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 |

# Phase I – March 2007

- Delivered complete CLEAT data collection.
- Provide ground truth for subset of data including signatures, stamps, logos, handwritten, and machine printed text.
- Provide document describing evaluation framework.

# Phase II – July 2007

- Deliver completed ground truthing and visualization tool for CLEAT metadata.
- Deliver Prototype version of CLEAT Software API Modules:
  - Document Image Enhancement,
  - Document Text/Image Text/Non-Text Discrimination,
  - Page Layout Similarity Ranking on CLEAT data,
  - Page Layer Segmentation and Zone Labeling, and
  - Content Labeling of Signatures, annotations, Stamps and Logos.
- Provide results of CLEAT API run on CLEAT datasets.
- Provide preliminary evaluation report.
- Provide basic API documentation

# Phase III

- Deliver Final version of CLEAT API.
- Provide training on use of CLEAT.
- Provide complete evaluation results on CLEAT data.
- Provide complete documentation of API.
- Provide feasibility report for system extensions.
- Provide a list of publications generated and planned as a result of this effort.

# WWW

- lamp.cfar.umd.edu/media/projects/cleat
- Contains
  - Summary
  - Proposal
  - Reports
  - Presentations
  - Milestones and Deliverables
  - Software
  - Data

# Agenda
# Tuesday AM

- Project Overview
  - Introduction
  - Goals and Objectives
- Tools
  - GEDI Display Environment
  - Datasets
- DocLib and Algorithms
  - Technical Presentations

# GEDI – Java Interface

# GEDI Features

- allows users to label and display rectangular zones in images

- supports user specified zone types

- handles type-specific attribute lists

- offers a graphical interface for editing and displaying zones

- enables users to create and distribute configuration files

- provides hotkeys for faster labeling

- can list multiple images in thumbnail views

- saves ground-truth and metadata as XML (compatible with DocLib)

# New Features

- Polygon and Oriented Boxes
- Scripting
- Text Alignment
- Multilingual support
- Additional Function Keys
- Bug Fixes

# Data Collection

**Datasets and Ground Truth –** A dataset containing examples of each class of document we process will be included. The dataset will contain a minimum of 5000 documents and be collected from a variety of sources, including the internet, existing training and testing datasets, public collections, project collections, and scanning. All ground truth will be provided in GEDI format and accompany the images

# 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 |

# Document Image Acquisition

- Sampling of Existing Databases
  - 20-25%
- Google Image Search
  - 60%
- Scanning hardcopy Document Images
  - 15-20%

# Document Genres

| Genre | |
|---|---|
| **Forms, Drawing, Tables et at.** | |
| **Forms** | 650 |
| **Drawing** | 80 |
| **Tables** | 100 |
| **Chemistry formulae** | 25 |
| | |
| **Math equations** | 165 |
| **Figures** | 40 |
| **Total** | 1060 |
| | |
| **Business documents and Memo letters** | |
| **Business documents** | 50 |
| | |
| **Business documents degraded** | 2700 |
| **Business documents with annotations** | 160 |
| **Memo letters** | 900 |
| **Total** | 3810 |
| | |
| **Journal and Conference Papers, Articles** | |
| **English** | 2785 |
| **German** | 360 |
| **Japanese** | 480 |
| **Total** | 3625 |

| Newsletters and Flyers | |
|---|---|
| **Google images** | 2400 |
| | |
| **Structured Documents** | |
| **Phonebook** | 229 |
| **Dictionaries (Chinese English, English Chinese)** | 1150 |
| **Yellowpage** | 80 |
| **Total** | 1459 |
| | |
| **Structured Documents** | |
| **Phonebook** | 229 |
| **Dictionaries (Chinese English, English Chinese)** | 1150 |
| **Yellowpage** | 80 |
| **Total** | 1459 |
| | |
| **Handwritten** | |
| **Chinese** | 146 |
| **Cyrillic** | 410 |
| **Japanese** | 47 |
| **Korean** | 80 |
| **Thai** | 319 |
| **Hindi** | 281 |
| | 1283 |

# Internet Downloads

| Genre | |
|---|---|
| **Figure** | |
| Good | 240 |
| Medium | 755 |
| Low | 548 |
| | |
| **Form** | |
| Good | 66 |
| Medium | 69 |
| Low | 32 |
| | |
| **Letter-Memo** | |
| Good | 55 |
| Medium | 88 |
| Low | 31 |
| | |
| **LIST** | |
| Good | 6 |
| Medium | 34 |
| Low | 11 |

| Newspaper | |
|---|---|
| Good | 22 |
| Medium | 37 |
| Low | 17 |
| | |
| **Publication Cover** | |
| Good | 130 |
| Medium | 425 |
| Low | 128 |
| | |
| **Receipt** | |
| Good | 10 |
| Medium | 50 |
| Low | 20 |
| | |
| **Screenshot** | |
| Good | 184 |
| Medium | 848 |
| Low | 566 |
| | |
| **Table** | |
| Good | 52 |
| Medium | 124 |
| Low | 42 |

# Maryland Datasets

- Collection of Free form Handwriting
  - Paid upto $1 for pages of native handwriting
  - Languages: Arabic, Amharic, Chinese, Korean, Japanese, Greek, Cyrillic, Hebrew, Thai, Burmese, and Hindi
  - Up to 1000 pages of each

以法治國的遺憾

李鵬基

　　趙紫陽同志不幸逝世了，噩訊傳來使人感到無比的悲痛。他為中國的改革開放和社會、經濟的發展作出了可讚頌的巨大貢獻，同期人民擁護他，讚揚他，愛戴他。然而由於"六·四"事件使他喪失了人身自由，十五年漫長的歲月過去了，但對他還是沒有作一個實事求是客觀的評價，使他到最終身，走綽不遺憾了！

　　現在不讀趙紫陽的"錯誤"性質為何，就在對他的處理方式而言，也是值得商榷的。遺憾以法治國的思想。使人撼與遺憾的是，方告有的領導人，一方面大讚特讚以法治國（嘴型也很重要），而另一方面卻又嘲諷以法治國，甚至踐踏以法治國。本來，趙紫陽是達此國家的重要領導人，曾是黨的總理、曾是黨的總書記，然而，未經任何法律手續予以審判，僅幾個老人的幾句話，就把趙紫陽長期敦禁起來，剝奪他的人身自由長達了十五年之久。這難道不是對以法治國的嘲弄和諷刺，是對以法

（1）

ที่ อเมริกา

5 กันยายน 2521

## 左ページ

Ⅱ. 途上国の社会経済的特質

1. 基本的様態 ～ 従属型二重経済

（本文の手書きメモ・図表）



## 右ページ

| | | |
|---|---|---|
| P85 ferment | n. | 불안, 동요, 정치적 동요 |
| usurp | v. | 강탈하다 |
| incessant | a. | 끊임없는 |
| shrink | v. | 줄어들다, 움츠리다, 기가 죽다 (p.p. shrunken) |
| rashly | ad. | 무모하게, 막무가내로, 성급하게 |
| massacre | v. | 학살하다 |
| flee | v. | 달아나다, 도피하다 (p. fled) |
| tribe | n. | 부족, 종족 |
| abandon | v. | 버리다, 포기하다 |
| nomadic | a. | 유목민의, 유목 생활을 하는 |
| *similize | v. | 동화하다, 동질으로 만들다 |
| cavalry | n. | 기병대, 기마부대 |
| ●refugee | n. | 망명자, 피난자, 난민 |
| perpetual | a. | 영구한, 부단한 |
| turmoil | n. | 혼란, 동요, 불안 |
| P87 undermine | v. | 약화시키다, 서서히 타락시키다 |
| *monastery | n. | 수도원, 사원, 사찰(원) |
| vast | a. | 광대한, 거대한, 엄청난 |
| *proportion | n. | 크기, 비례, 비율 |
| realm /relm/ | n. | 왕국, 영역 |
| bureaucracy | n. | 관료제, 관료정치, 관료주의 |
| *exert | v. | (힘·권력 등을) 행사하다, (영향을) 미치다 |

　　　< Taoism > 도교

| | | |
|---|---|---|
| Taoism | n. | 도교 |
| ●*calligraphy | n. | 서예, 서도, 달필, 필적 |
| conglomeration | n. | 응집, 집합 |

И. Бродский.

Под вечер он видит...

— 1 —

Под вечер он видит, застывши в дверях,
два всадника скачут в окрестных полях,
как будто по кругу, сквозь рощу и гарь,
и разно не могут друг друга догнать.
То бросив поводья, лошицу угав,
то снова в игле возбуждение привстав,
и быстро по светлому склону холма,
то в рощу опять, где сгущается тьма.

Два всадника скачут в вечерней грязи,
не только от дома, от сердца вблизи,
друг друга они окликают, зовут,
небесные рапе за рощу плывут.
И так никогда им на свете вдвоём
сквозь рощу и гарь, сквозь пустой водоём,
не скакать в виду стапулающих постов,
как будто меж нами не сотня кустов!

Стихи под эпиграфом.    И. Бродский.

„ То, что дозволено Юпитеру,
     не дозволено быку..."

Каждый пред Богом
                 наг.
Жалок,  наг
         и убог.
В каждой музыке  Бах,
В каждом из нас  Бог.
Ибо вечность — богам.
Бренность — удел быков...
Богово станет нам
Сумерками богов.
И надо небом рискнуть,
И, может быть, невпопад.
Ещё нас не раз распнут
И скажут потом: распад.

# Other "Documents"

# IBM Cross Pad Data

- 30 boxes, 30 writers producing 50-80 pages each
- 25000 pages total / 1 million words
- Most European Languages: German, French, Italian, English (UK), and Spanish
- Makeup: Characters (~8 boxes), Phrases, Freeform (1 box)
- Contracted with IBM to make the data public

# MELLO ~~DECIDES NOT~~ ~~T~~ 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 he said subsequently on several occasions the chances were slim he 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 40-year-old title insurance company executive said he did not want to devote more time to city business at the expense of his family and job. He said he currently spends about 30 hours a week on city-related ~~duties~~ and that being mayor would mean at least 10 hours per week.

---

**Bim, Bam, Bum - Ein Glockenton**

Bim, Bam, Bum - Ein Glockenton

fliegt durch die Nacht, als

fliegt durch die Nacht, als

hätt' er Vogelflügel, er fliegt

hätt' er Vogelflügel, er fliegt

in römischer Kirchentracht

in römischer Kirchentracht

wohl über Tal und Hügel. Er

wohl über Tal und Hügel. Er

---

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 dire.

"People are going to starve to death," warned Fremont Mayor Bill Ball.

The Fremont effort is focused on the area around Firebaugh, at California's largest olive processing plant and, according to some observers, the "heart of the California olive-growing area."

The severe winter wiped out the citrus crop, throwing unemployment, hovering between Fresno and Bakersfield, where unemployment, Fresno Week will not substantiate, will remain above those until November.

On May 5, a farmworker relief caravan from Fresno, San Francisco, El Cerrito, and Richmond delivered 12 tons of food, two tons of rice, two tons of beans and hundreds of pounds of sugar, bread, cooking oil and used clothing.

# New Data

- 25,000 pages ground truthed to the zone level
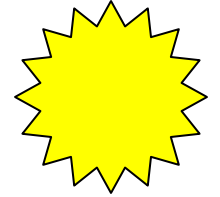- Sampled from the Tobacco Litigation Corpus of 49 Million pages

# Distribution (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 |

# Agenda
# Tuesday AM

- Project Overview
  – Introduction
  – Goals and Objectives
- Tools
  – GEDI Display Environment
  – Datasets

# DocLib and Algorithms
  – Technical Presentations

# 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>
```

# Agenda
# Tuesday AM

- Project Overview
  - Introduction
  - Goals and Objectives
- Tools
  - GEDI Display Environment
  - Datasets
- **DocLib and Algorithms**
  - Technical Presentations

# Technical Presentations

- Page Segmentation (and rule line separation)
- Page Layout Similarity
- Document ID/Script ID

This afternoon

- Logo Detection and Recognition
- Signature Detection
- Font OCR

# Technical Presentation Outline

- Overview of Problem
- Technical Approach
- Datasets
- Results
- Implementation and Software

# Examples of Drivers

- ScriptID [-x] [-l] filename
  - -x  ---  Write classification results into an xml file for each input image. It creates a new xml file if no associated xml file exists.
  - -l  ---  File containing the list of input images to execute
  - -h  ---  Show help at command line


- DocID [-x] [-l] filename
  - -x  ---  Write classification results into an xml file for each input image. It creates a new xml file if no associated xml file exists.
  - -l  ---  File containing the list of input images to execute
  - -h  ---  Show help at command line

# 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.

# Motivation

- Layer analysis and separation for general, heterogeneous documents, is a very hard problem.

- Handwritten documents are very important
  - Handwriting was developed a long time ago as a means to expand human memory and to facilitate communication.
  - We are continuing to produce handwritten documents.

# 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 |
| $2\times2$ 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

| 2 | | 1 |
|---|---|---|
| | 4 | Center |
| | | 3 |

$\rightarrow D_h \leftarrow$

# Clique Potential

- ## Clique potential:

$$V_p(c) = -\frac{P(x_l, x_c, x_r)}{\left(P(x_l)P(x_c)P(x_r)\right)^w} \qquad V_n(c) = -\frac{P(x_c, x_1, x_2, x_3, x_4)}{\left(P(x_c)P(x_1)P(x_2)P(x_3)P(x_4)\right)^w}$$

Probabilities are estimated from ground truth.

- ## Total energy of Gibbs distribution:

$$U(\underline{X}/\underline{Y}) = -w_s \sum_{s \in \Omega} P(x_s / y_s) + w_p \sum_{c \in C_p} V_p(c) + w_n \sum_{c \in C_n} V_n(c)$$

HCF (Highest Confidence First) method is used to minimize the energy function.

# MRF Postprocessing Example

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

# Background Pattern (Rule Line) Separation

- Many handwritten documents are produced on rule lined paper
- These lines should be detected and removed before we feed the text to an Optical Character Recognition (OCR) engine.



**An Arabic handwritten document on rule lined paper**

# Challenges

- Previous work
  - Hough transform
  - Projection based (Strip Projection [Chen98] and Skew Projection [Liu95])
  - Vectorization based (BAG [Jain96], SPV [Dori99], and DSCC [Zheng01])
- Challenges
  - The documents may be degraded with severely broken lines.
  - High accuracy and low false alarm rate are demanded.

# Challenges



**Original document**



**Line detection results using DSCC method**

- We propose a model-based method
  - Model the horizontal projection profile with an HMM model.
  - Under the model, lines are detected simultaneously.

# Preprocessing

- Text filtering
- Skew estimation and correction

# HMM Model for Parallel Lines

- Model the projection profile with an HMM model

  - The vertical position of lines $\{Y_i\}$ form a Markov Chain

  - We can not observe $\{Y_i\}$ directly, but projection profile

  - The gaps between neighboring lines are consistent on the same page



- Parameters of the model are estimated from ground truth.
- Viterbi algorithm is used to decode the model.

# 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%) |

# Comparison with Other Methods

- Hough transform
- DSCC
- Projection based methods

COMPARISON OF OUR MODEL-BASED METHOD WITH OTHER METHODS ON THE TEST SET (THERE ARE A TOTAL OF 1,596 GROUNDTRUTHED LINES).

| | Detected Lines | Correct | Partial Correct | Missed | False Alarm |
|---|---|---|---|---|---|
| Hough Transform | 1,588 | 1,299 (81.4%) | 60 (3.8%) | 237 (14.9%) | 229 (14.4%) |
| Projection Method | 1,577 | 1,310 (82.1%) | 112 (7.0%) | 174 (10.9%) | 155 (9.7%) |
| DSCC | 2,162 | 1,398 (87.6%) | 118 (7.4%) | 80 (5.0%) | 646 (40.5%) |
| Our Model-Based Method | 1,631 | **1,545 (96.8%)** | **49 (3.0%)** | **2 (0.1%)** | **37 (2.3%)** |

# Software

- Implemented as a set of Libraries
- Trainable with new data

# Technical Presentations

- Page Segmentation (and rule line separation)
- Page Layout Similarity
- Document ID/Script ID

This afternoon

- Logo Detection and Recognition
- Signature Detection
- Font OCR

# Multi-Class Classification using Document Layout

- Motivation

- Document Representation

- Random Chopping

- Feature Selection

- Score Function

- Experiments

- Summary and Future work

# 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

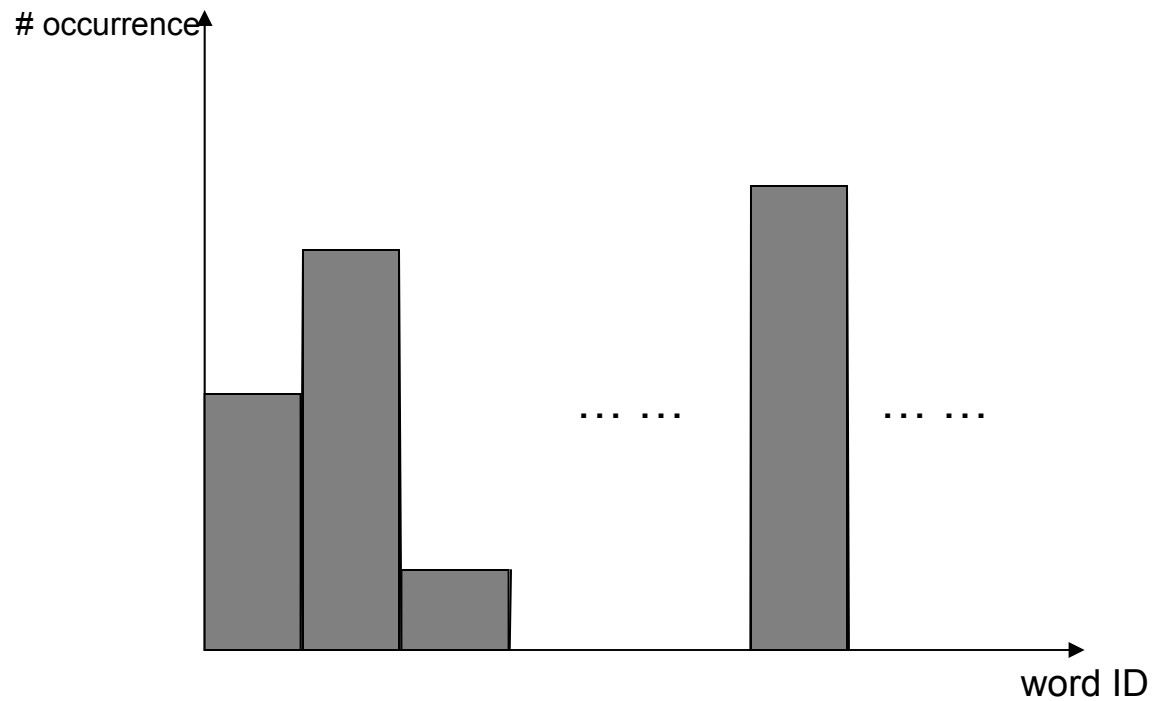- Document represented this way is translation and 180° rotation invariant

# 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

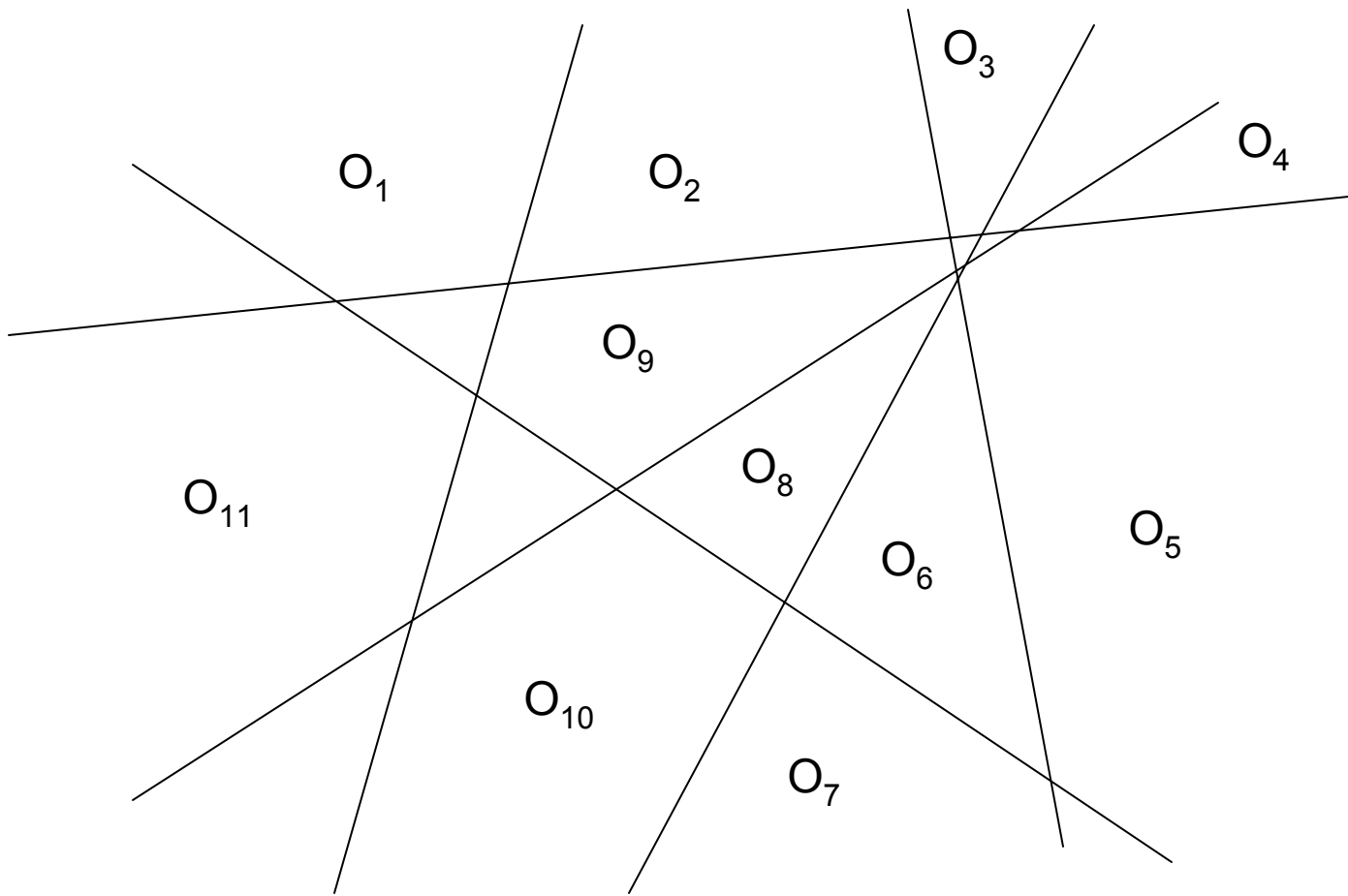- From 101 documents, we built a dictionary with 976 words

# Document Representation

# Random Chop – the idea

# Specifically

- For each layout class, we choose some training samples

- For i= 1 to NUM_CHOPS
  - Randomly chop layout classes into two classes
  - Validity checking of current chop
  - Feature Selection
  - Train a binary discriminative classifier using Logistic Regression on training samples
  - Evaluate the classifier on a validating set

# Feature Selection

- Document histogram vector lies in a very high dimension space

- Select subset of features that is relevant to the chopping in consideration

- CMIM criterion : Conditional Mutual Information Maximization

    $v(1) = \text{argmax } I(Y; X_i)$   $1 \leq i \leq N$

    $v(k+1) = \text{argmax } \{\min I(Y; X_i | X_{v(l)})\}$ , $1 \leq k \leq K$, $l \leq k$

- Stopping criteria:
    - Maximum number of selected features is reached
    - Information gain is lower than a threshold

# 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

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

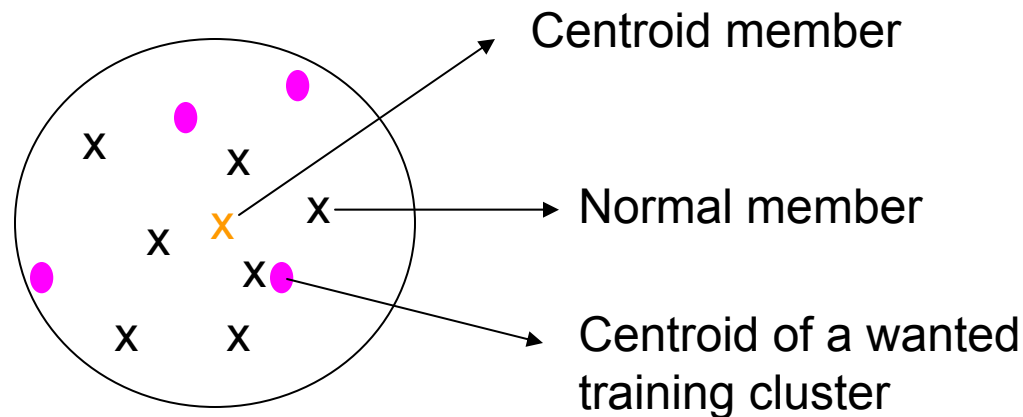$$C = \text{argmax}_i \ \text{Score}_i$$

# Scoring a testing document

- $S = (\sum_i N_i * W_i) / (\sum_i N_i)$

  $W_i$ : weight of the wanted training cluster which is the nearest neighbor within fixed range of testing cluster i.

  $N_i$ : size of cluster i.



Centroid member

Normal member

Centroid of a wanted training cluster

Original: neighbor with the highest weight;

Lamp: nearest neighbor

# 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

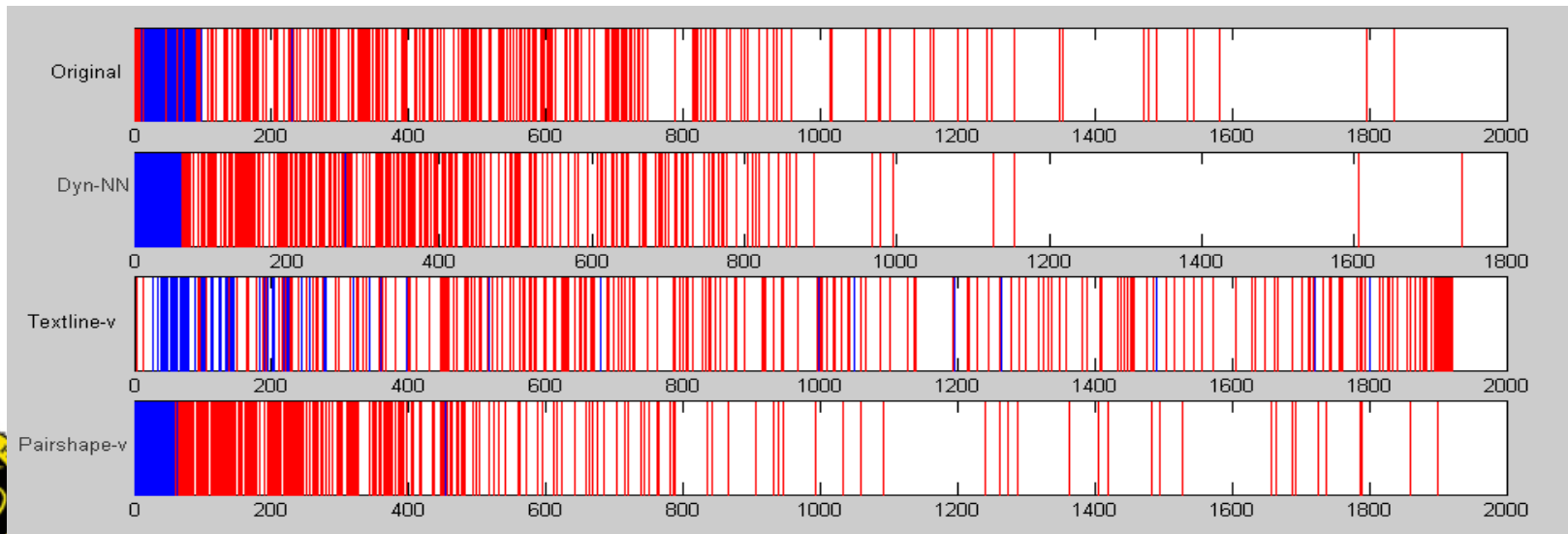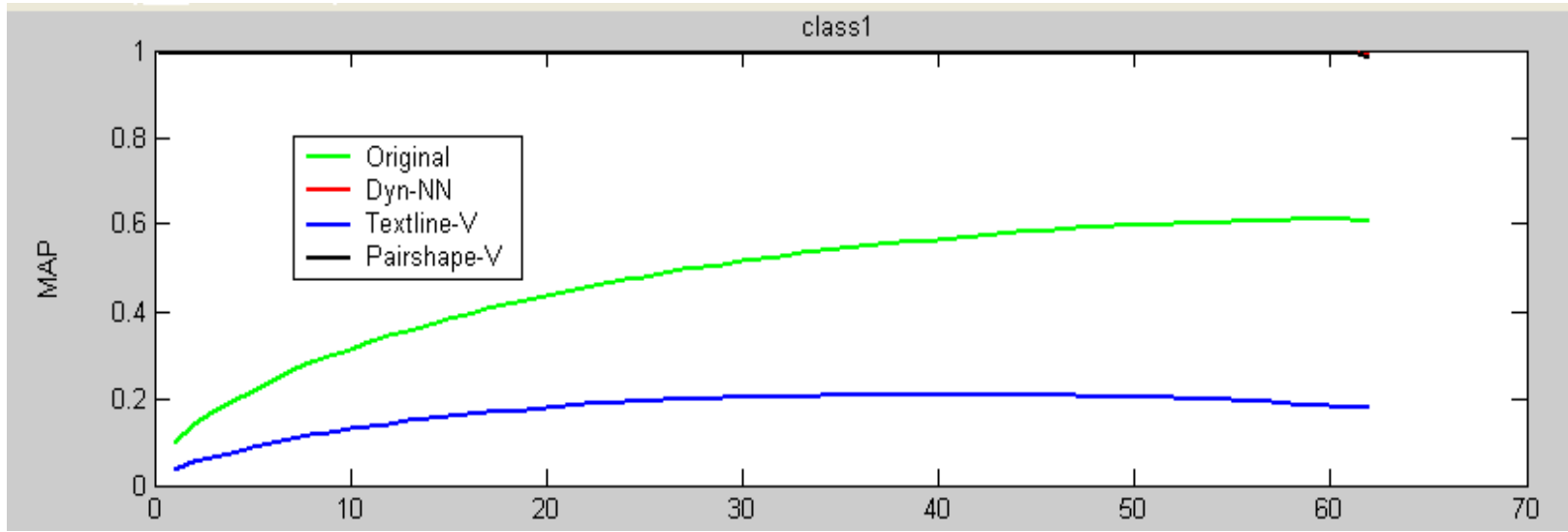# Experimental Results
## --Confusion Matrix

| | 1c | 2c | 1r2c | 3c | 2c_a sym | 2c2c_ asym | class 1 | class 2 | clas s3 | clas s4 |
|---|---|---|---|---|---|---|---|---|---|---|
| **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** (95) | | | | 1 | | 1 | 17 | 27 | 7 | 52 |

# Experiments – ARR Results

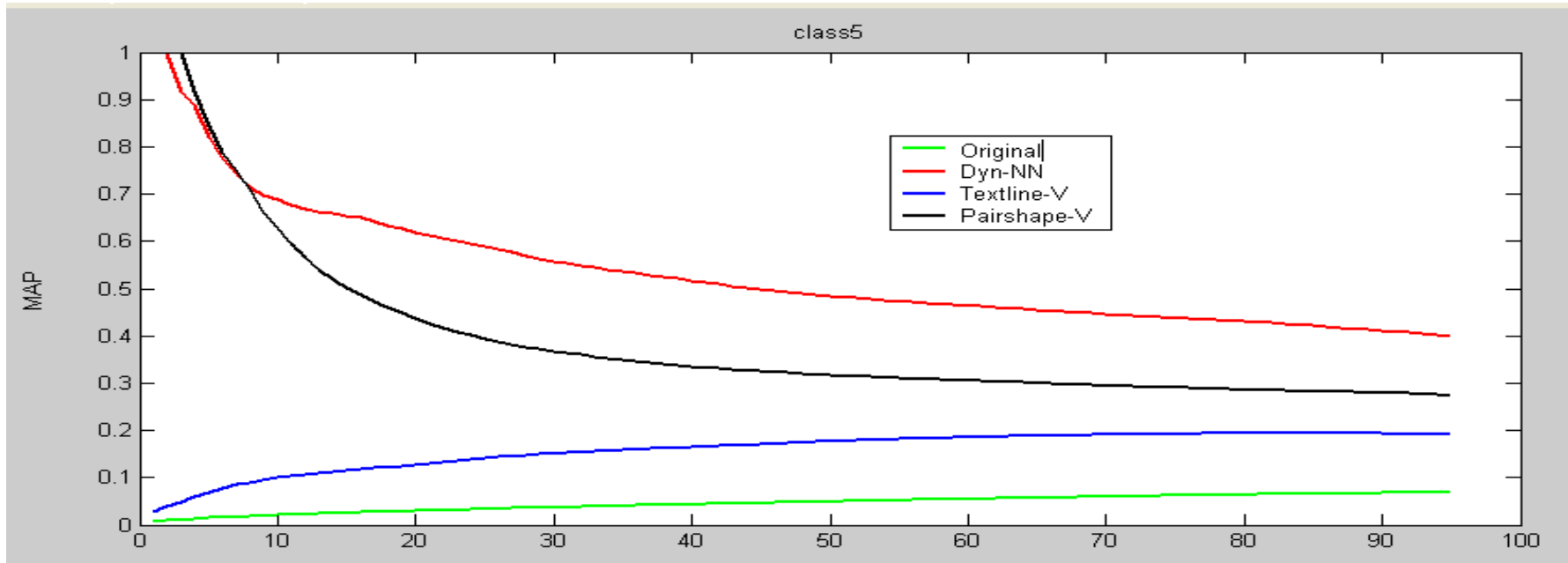| | Original | Dyn-NN | Text-V | Pair_V |
|---|---|---|---|---|
| 1c | 0.450 | 0.011 | 0.038 | 0.043 |
| 2c | 0.062 | 0.010 | 0.324 | 0.087 |
| 3c | 0.028 | 0.0002 | 0.504 | 0.013 |
| 1r2c | 0.148 | 0.063 | 0.245 | 0.105 |
| 1r1r2c | 0.159 | 0.010 | 0.103 | 0.045 |
| 1r2c2c | 0.121 | 0.067 | 0.186 | 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

# Summary and Future Work

- ## Conclusions:
  - Time efficiency
  - Space efficiency: only need to store classifier parameters and class signatures
  - Easy to combine new layout classes
  - Generalizability : is able to tell, to some degree, whether a new pair of instances unseen in the training set are of similar layout

- ## Future Work:
  - Find out the optimal number of chops for a given number of classes
  - Guarantee non-overlapping of classes
  - Try different classifiers, like NB, SVM

# Software

- Currently Implemented as DocLib
- Line Detection Modules Improved

# Technical Presentations

- Page Segmentation (and rule line separation)
- Page Layout Similarity
- Document ID/Script ID

## This afternoon

- Logo Detection and Recognition
- Signature Detection
- Font OCR

# 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.

# ScriptID

- Motivation
- Challenges
- Observation
- Descriptor
- Implementation
- Results

# The Motivation

- Speedup the recognition process
  - Turn on the OCR engine, when necessary;
- Improve the accuracy
  - Select different OCR engines for different scripts;
- Understand the human perception
  - Can we recognize different scripts before recognizing individual characters?

# The Challenge

- Handwritten documents
  - Template matching cannot be used in general.
- The method needs to be fast
  - Naïve trial-and-error methodology doesn't work
- The method needs to be invariant to
  - Scale
  - Rotation

# 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

1　들어　있는　것

어릴　것　같은데

입고　있었다.

स्वायत्त विनियोग

जो राष्ट्रीय आय

वर्तमान आर्थिक

न धारित 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

# Classifier: Support Vector Machine

- Suppose we have *N* classes

- For each class, we train 1 SVM using images from this class vs other classes.

- Result: *N* SVM classifiers (linear classifier in high dimensional space)

# Dataset for Classification

- Arabic

# Dataset for Classification

- Chinese

# Dataset for Classification

- Hindi

# Dataset for Classification

- Korean

# 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)

# ImageID

- Motivation
- Challenge
- The Approach
- Results

# The Motivation

- 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...

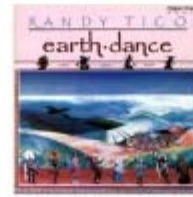google_cd_cov...

google_cd_cov...

google_cd_c...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...
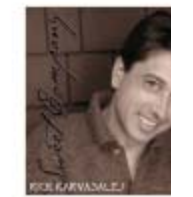
google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_c...

google_cd_cov...
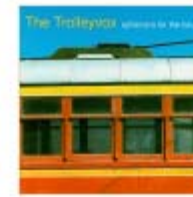
google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_cov...

google_cd_c...

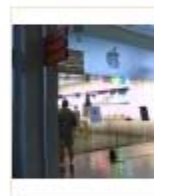google_apple_0.tif  google_apple_1.tif  google_apple_2.tif  google_apple_5.tif  google_apple_9.tif  google_apple_1...  google_apple_1...  google_apple...

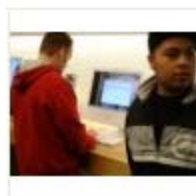google_apple_1...  google_apple_1...  google_apple_1...  google_apple_2...  google_apple_2...  google_apple_2...  google_apple_2...  google_apple...

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...

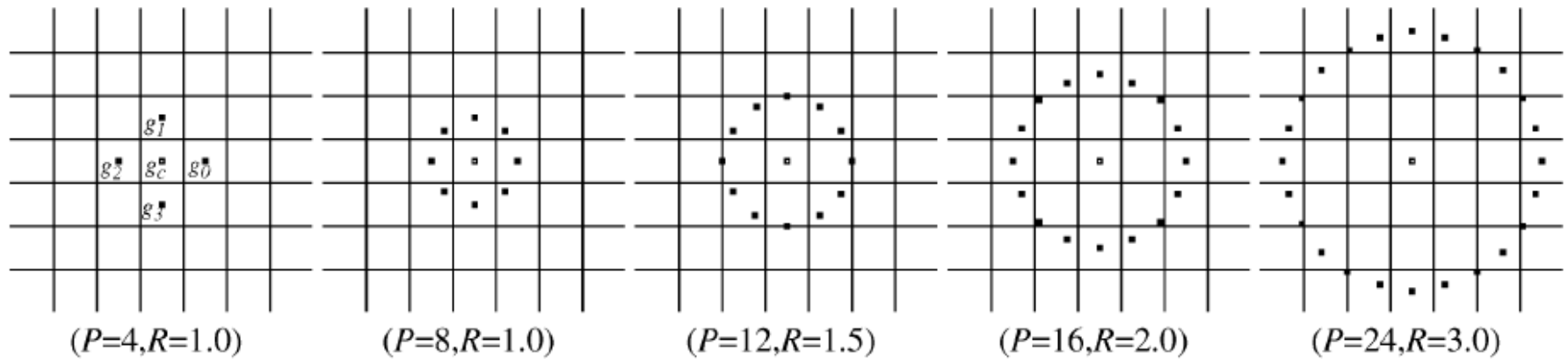google_apple_5...  google_apple_5...  google_apple_5...  google_apple_5...  google_apple_5...  google_apple_5...  google_apple_5...  google_apple...

# The Spatial Descriptor



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

- P: number of neighbor pixels
- R: neighbor size

# LBP: Local Binary Pattern

- Define

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

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

- Example

| $g_3$ | $g_2$ | $g_1$ |
|---|---|---|
| $g_4$ | $g_c$ | $g_0$ |
| $g_5$ | $g_6$ | $g_7$ |

# The LBP representation

- Given an image.
- Transform the distribution vector into an P-bit pattern code ("Binary pattern")

$$LBP_{P,R} = \sum_{p=0}^{P-1} s\,(g_p - g_c)\, 2^p$$

  – s: scale factor

# Other variations of LBP

- Rotation invariant
- Different neighbor points and area
- "uniform" pattern

# The performance

- Confusion matrix

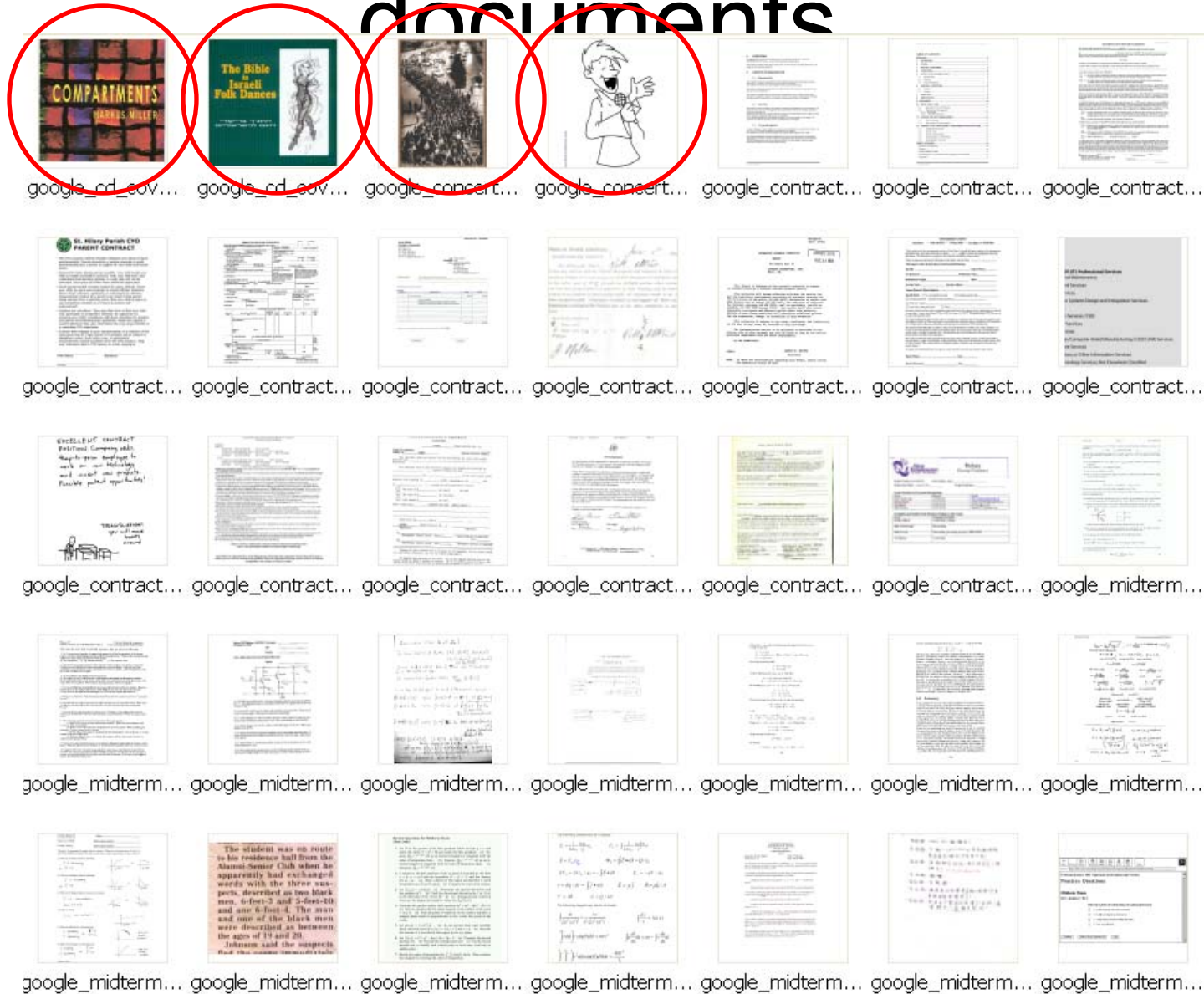|  | Doc | Image w/ | Non doc |
| --- | --- | --- | --- |
| Doc | 0.8557 | 0.1340 | 0.0103 |
| Image w/ | 0.1725 | 0.6011 | 0.2264 |
| Non doc | 0.0444 | 0.1422 | 0.8133 |

# 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

# Improvement

- Incorperate the distribution of grayscale:
  - an important clue for classification
- Try larger neighbor area for LBP
- Combine with other descriptors
  - Appearance model

# Future work

- ## ScriptID
  - Test more scripts. 10-15 would be a reasonable goal

- ## ImageID
  - Improve the performance of the classification of the image w/ text vs images .

# Technical Presentations

- Page Segmentation (and rule line separation)
- Page Layout Similarity
- Document ID/Script ID

## This afternoon

- Logo Detection and Recognition
- Signature Detection
- Font OCR