

CLEAT:

**A Classification, Enhancement and
Analysis Toolkit for Heterogeneous
Document Image Collections**



LAMP History

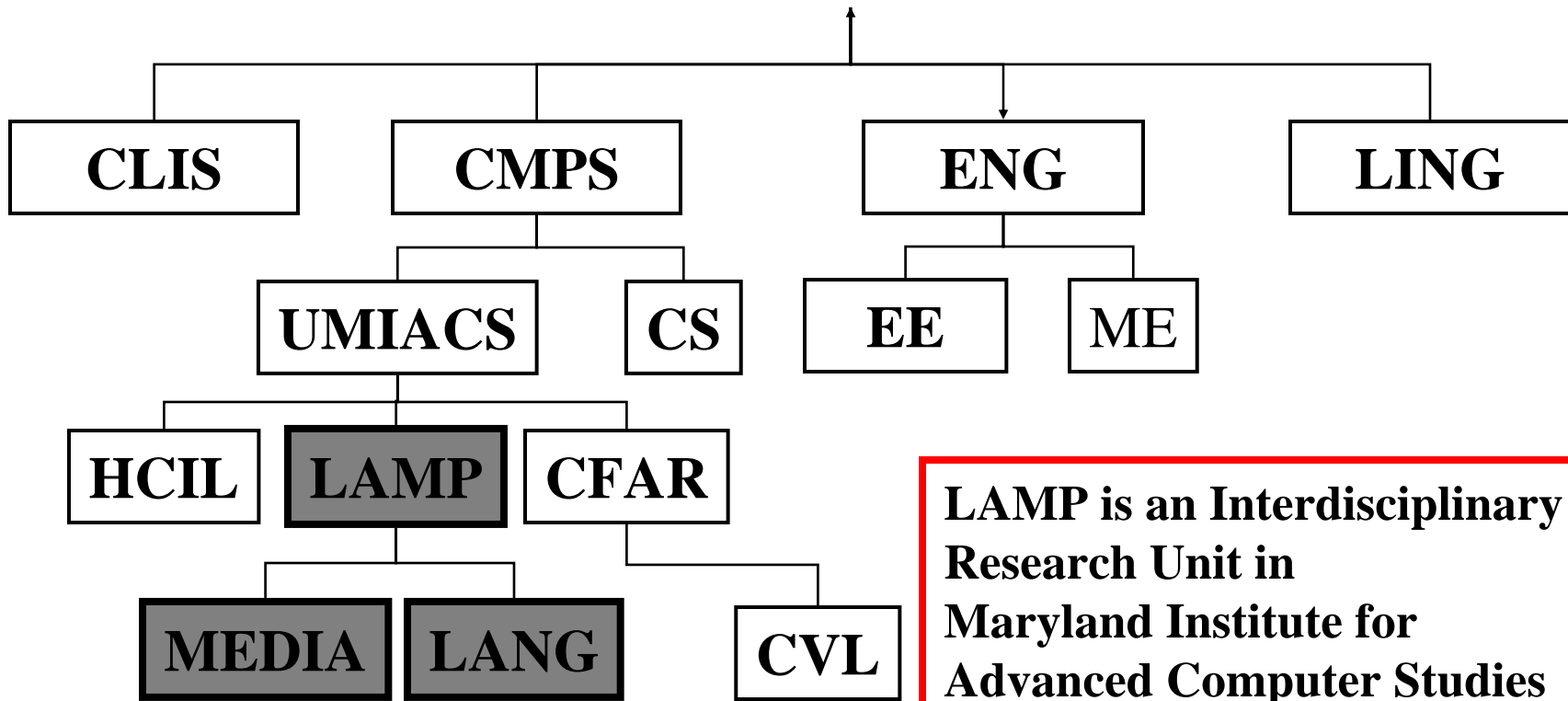
- Began in 1996 with a focus on documents
- Produced 9 PhD (2 more expected in 2007)
- Over 200 scientific publications
- Almost 50 Students (Undergrad-Graduate)
- Numerous Technology Transfer Opportunities



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





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



Media Personnel

DOCUMENTS

Stefan Jaeger

RESEARCH FACULTY

David Doermann

VIDEO

*Daniel DeMenthon***

Graduate Students

Guangyu Zhu
Yi Li
Rajat Ahuja
Nagia Ghanem*
Xu Liu

Burcu Karagol-Ayan
Sameer Kibey*
May Huang
Zhe Lin

David Mihalcik*
Yang Yu
Xiaodong Yu
Mudit Agrawal

Undergraduate Researchers

Mike Roth*

Faculty Research Associate

Tandeep Sidhu



Affiliates

- Research/Sponsor Affiliations
 - Ricoh Japan/CRC (Scanning Equipment)
 - Xerox (OCR software), Panasonic, KLA-Tencor, Hitachi
 - University of Oulu (Multimedia Documents in Telecommunications)
 - Army Research Laboratory, Other Intelligence Community Participants



*National
Security
Agency*



HITACHI
Inspire the Next

RICOH®

THE
DOCUMENT
COMPANY
XEROX



Panasonic
ideas for life



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



Agenda

Project Overview

- Introduction
- Goals and Objectives

Data Collection and Ground Truth

GEDI and Evaluation Framework

Evaluation and Research Components

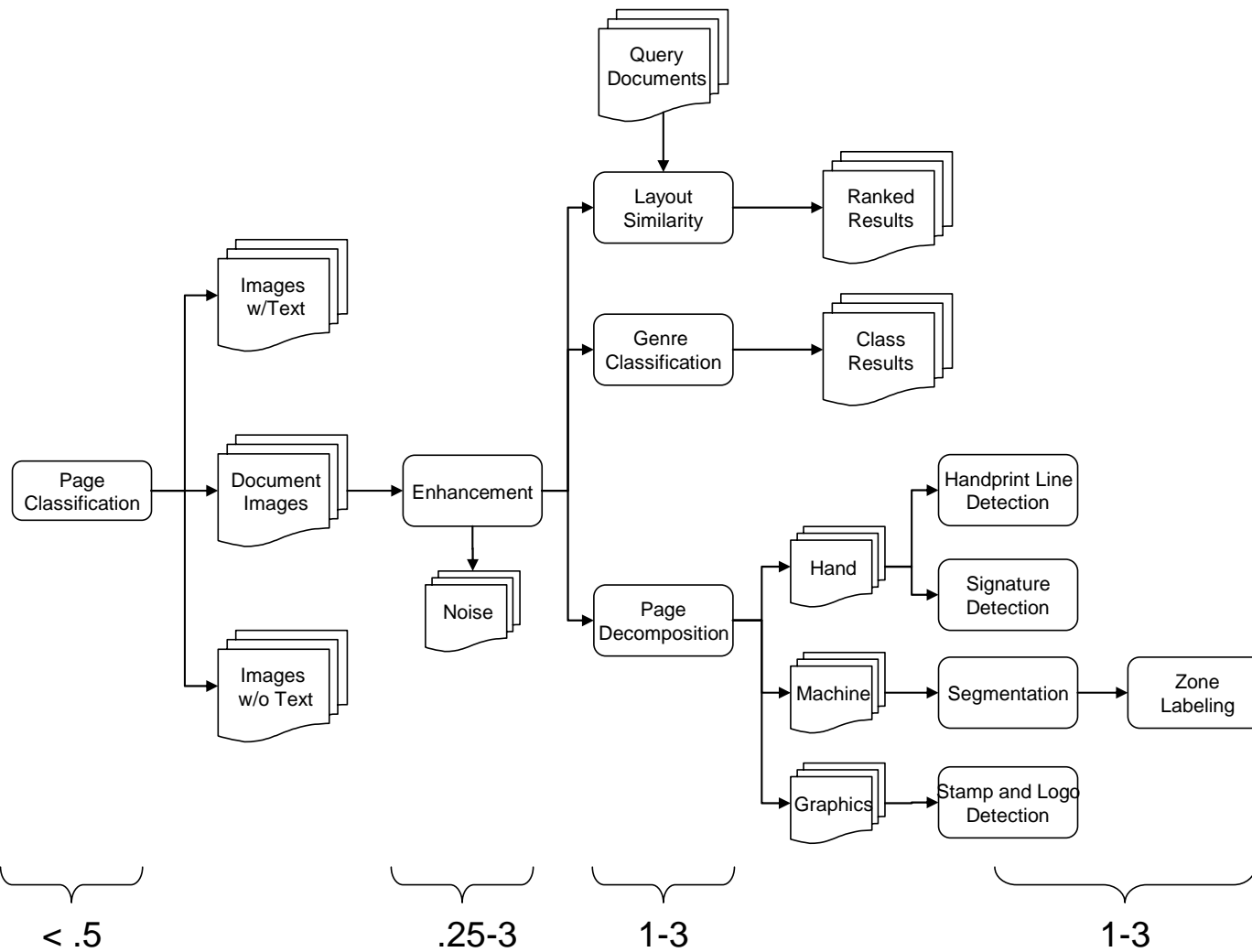
- Unconstrained Signature and Logo Detection and Matching for Off-line Document Image Retrieval
- Document/Non-Document Discrimination

Technical Presentations

- Adaptive OCR
- Document Classification by Layout
- Learning from Web for Shape-based Object Recognition



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 1 - March 31, 2007**
 - Deliver completed CLEAT data collection.
 - Provide ground truth for subset of data including signatures, stamps, logos, handwritten, and machine printed text.
 - Provide document describing evaluation framework.



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.		Newsletters and Flyers	
Forms	650	Google images	2400
Drawing	80	Structured Documents	
Tables	100	Phonebook	229
Chemistry formulae	25	Dictionaries (Chinese English, English Chinese)	1150
Math equations	165	Yellowpage	80
Figures	40	Total	1459
Total	1060	Structured Documents	
Business documents and Memo letters		Phonebook	229
Business documents	50	Dictionaries (Chinese English, English Chinese)	1150
Business documents degraded	2700	Yellowpage	80
Business documents with annotations	160	Total	1459
Memo letters	900	Handwritten	
Total	3810	Chinese	146
Journal and Conference Papers, Articles		Cyrillic	410
English	2785	Japanese	47
German	360	Korean	80
Japanese	480	Thai	319
Total	3625	Hindi	281
			1283



Handwritten notes and a diagram on a piece of paper. The notes include a list of items and a circular diagram with various annotations and arrows. The text is written in German and includes phrases like "Wissen" and "Pflanzen".

MARKTSTÄTTE	MARKT	ANZEIGENKATEGORIE	BEFRAG
DALLS SEE (SONDRIC PRICE)	AKKAPPAJ	NEWSPAPER	15/10/99
BEAGE (SONDRIC PRICE)	AKKAPPAJ	SONDAG SUPPLEMENT	11/10/99
EALE (SONDRIC PRICE)	DIOSKIA	SONDAG SUPPLEMENT	11/10/99
EALE (SONDRIC PRICE)	DIOSKIA	NEWSPAPER	15/10/99
EALE (SONDRIC PRICE)	DIOSKIA	MARKETS	09H 12/10-11/99
RED FAUSE (SONDRIC PRICE)	SEKASAKA MARKETS	09H	1/99
SEAGE	HOLAND	SONDAG SUPPLEMENT	11/10/99
SEAGE	HOLAND	NEWSPAPER	15/10/99
SEAGE	INDONESIA	SONDAG SUPPLEMENT	11/10/99
SIGL QTY READ-09	AKKAPPAJ MARKETS	09H	1/99-11/99
SIGL QTY READ-09	OSKON MARKETS	SONDAG SUPPLEMENT	11/10-11/99
SIGL QTY READ-09	PORTLAND, OREGON	SONDAG SUPPLEMENT	11/10/99
"CORTVETE"	ECHERE, OREGON	SONDAG SUPPLEMENT	11/10/99
"CORTVETE"	AKKAPPAJ MARKETS	SONDAG SUPPLEMENT	11/10/99
"TEAM - LOW PRICE"	AKKAPPAJ	NEWSPAPER	11/10
"BEAGE"	OSKON MARKETS	NEWSPAPER	11/10

Handwritten notes and a diagram on a piece of paper. The notes include a list of items and a circular diagram with various annotations and arrows. The text is written in German and includes phrases like "Wissen" and "Pflanzen".

Handwritten notes and a diagram on a piece of paper. The notes include a list of items and a circular diagram with various annotations and arrows. The text is written in German and includes phrases like "Wissen" and "Pflanzen".

Handwritten notes and a diagram on a piece of paper. The notes include a list of items and a circular diagram with various annotations and arrows. The text is written in German and includes phrases like "Wissen" and "Pflanzen".

The Washington Post

After Zargawi, No Clear Path In Weary Iraq

The Young Apprentice

How U.S. Forces Found Iraq's Most-Wanted Man

DEUTSCHES BILANZKONTO (DEBT)				
DreiM/J.A.	0,270	1,30		
1				
1				

THE WORLD'S RICHEST PEOPLE

SPECIAL ISSUE

Forbes

BILLIONAIRES

946 BILLIONAIRES

946 BILLIONAIRES

946 BILLIONAIRES

Journal of Neural Probes

FINANCIAL TIMES

IRAQ braced as Saddam rejects Bush ultimatum

Fed holds rates but signals uncertain future

US probes Abld collision claim



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



以片治国的遗憾

李周基

赵紫阳不幸逝世，噩讯传来使人感到无比的悲痛，他为中国的改革开放和社会、经济的腾飞作出了不可磨灭的巨大贡献。同时人们称赞他，赞扬他，爱戴他；然而由于“六·四”事件使他丧失了人身自由，十五年漫长的囚禁生活，但他还没有作一个贯穿其一生的评价，使他抱恨终身，交待不遗憾了！

现在不评赵紫阳的“错误”性质如何；就在对他的处理方式而言，也是极端错误的，违背了以片治国的思想。使人产生遗憾的思，古语有“欲求大治，必先严法”（当然也很重要），而另一方面却又背离以片治国，甚至践踏以片治国。本来，赵紫阳是党的国家的重要领导人，党通过总结经验，高而立的总结化，然而，未经任何法律程序而在审判，造成国人的愤慨，就把赵紫阳长期软禁起来，剥夺他的人身自由长达十五年之久。这难道是对以片治国的误解和讽刺，也是对“治

1) आपल्या विरोधा वाड !
आपला विरोध वाढ !
आपला विरोध वाढ !
आपला विरोध वाढ !
आपला विरोध वाढ !
आपला विरोध वाढ !
आपला विरोध वाढ !
आपला विरोध वाढ !
आपला विरोध वाढ !
आपला विरोध वाढ !

2) विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !

3) विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !

4) विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !

5) विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !
विनाशकारी वाड !



Πολύγωνο
5 πλευρών 2521

Πολύγωνο
 Είναι ένα κλειστό σχήμα που αποτελείται από τμήματα ευθείας γραμμής που ονομάζονται πλευρές. Τα σημεία όπου οι πλευρές συναντούνται ονομάζονται κορυφές. Τα πολύγωνα ταξινομούνται με βάση τον αριθμό των πλευρών τους. Ένα πολύγωνο με 5 πλευρές ονομάζεται πεντάγωνο. Η περίμετρος ενός πολύγωνα είναι το άθροισμα των длиνών των πλευρών του. Η επιφάνεια ενός πολύγωνα είναι το εμβαδόν που καταλαμβάνει το σχήμα.

Εάν P_1 είναι η περίμετρος του P_2

$P_1 = 5 \times P_2$

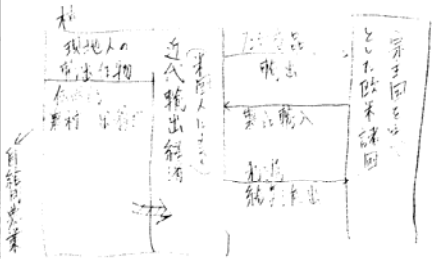
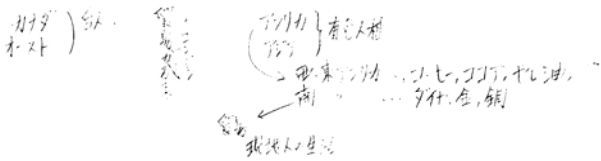
Η περίμετρος του P_1 είναι 5 φορές μεγαλύτερη από την περίμετρο του P_2 .
 Η επιφάνεια του P_1 είναι 25 φορές μεγαλύτερη από την επιφάνεια του P_2 .
 Η απόσταση από την αρχή των αξόνων μέχρι την ευθεία P_1 είναι 5 φορές μεγαλύτερη από την απόσταση από την αρχή των αξόνων μέχρι την ευθεία P_2 .



II 遂上国の社会経済的特徴

1. 基本的構造 ~ 依存型二重経済

この(西)日本経済構造は、作られた、小口村落地型二重経済と呼ぶ。これは、他の(東)日本(本邦)に比べて、発展した急速な作られた、小口村落地型経済を余議破る。とある。



村落地二重経済は、本邦の一次産業を特化したものである。これは、他の(東)日本に比べて、発展した急速な作られた、小口村落地型経済を余議破る。とある。

6. The Three Kingdoms and the Six Dynasties

- 885 ferment n. 불안, 동요, 정치적 동요
 usurp v. 강탈하다
 incessant a. 끊임없는
 shrink v. 줄어드다, 축소되다, 기가 죽다 (p.p. shrunken)
 rashly ad. 무모하게, 빈중하게, 성급하게
 massacre v. 학살하다
 flee v. 관아나다, 도피하다 (p. fled)
 tribe n. 부족, 종족
 abandon v. 버리다, 포기하다
 nomadic a. 유목민의, 유목 생활을 하는
 + Sinicize v. 중국화하다, 중국식으로 만들다
 cavalry n. 기병대, 기마부대
 refugee n. 망명자, 피난자, 난민
 perpetual a. 영구한, 부단한
 turmoil n. 혼란, 동요, 불안
 887 undermine v. 약화시키다, 서서히 무너시키다
 + monastery n. 수도원, 사원, 선원(선)
 vast a. 광대한, 거대한, 엄청난
 + proportion n. 크기, 비례, 비율
 realm n. 왕국, 영토
 | realm |
 bureaucracy n. 관료제, 관료사회, 관료주의
 + exert v. (힘, 권력 등을) 행사하다, (영양을) 미치다
 < Taoism > 도교
 Taoism n. 도교
 + calligraphy n. 서예, 서도, 장필, 필법
 conglomeration n. 융합, 결합

И. Бродский.

Год вперь он видит...

- 1. -

Год вперь он видит, заставши в дверях,
два владника скажут в окрестных полях,
как дует по кругу, сквозь роуц и гар,
и рвало не могут друг друга позарь.
То троеив поворья, пошлится, угаб,
то ехва в игле воздржимо привстав,
и дыстро по светлону склоу жасина,
то в роуц опер, еде егуцаеся тьма.

Два владника скажут в вечерней тьме,
не только от дома, от сердца близки,
друг друга они окликают, зовут,
нейсине рарь за роуц пикбут.
И так шкоча им на свете брвоим
сквозь роуц и гар, сквозь пургой вояем,
не ехарь в виду стационарных поствов,
как будто млет ишим не согна куеров!

Стрхи под эпиграфами И. Бродский.

То, что дозволено Юпитеру,
не дозволено дысу...

Каждый пред Богом нас.

Малок, нас и убо.

В каждой музыке Бах,

В каждом из нас Бог.

Что ведомо — Богом.

Брешось — удем дыков...

Богово етает нам

Сумерсаним догов.

И кадо недом чрисклуга,

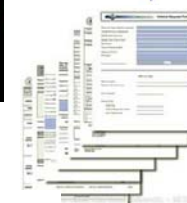
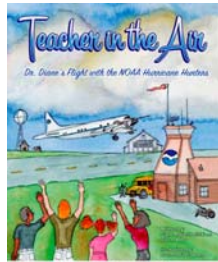
И, может дыгь, невпаад.

Еще нас не раз расплуд

И ехасу потом: расплуд.



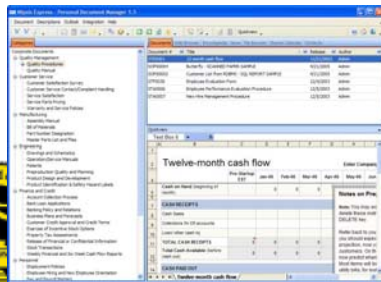
Other "Documents"



Parliamentary Assembly
Assemblée parlementaire



COUNCIL OF EUROPE
CONSEIL DE L'EUROPE



Honey, I think we are beyond the point of me being just your "boyfriend." It's about time you started calling me what I really am.

And that is...?
Your manfriend.



Order tables now at www.poker-wear.com
Prices shown here are subject to change without notice.

Poker Tables from Poker-Wear

You can play poker without a poker table and Poker-Wear.com has a great selection of folding chairs that can be set up for the big game, after and then store with the rest of your gear. Or, you can have those great looking poker tables set up for the time.

Poker-Wear.com is offering **FREE SHIPPING ON poker tables** everywhere in the US states. This is a limited time offer. Order your favorite poker table today.

Poker Table
This table is made of high quality materials, built to last and is easy to move. It's a great choice for your home or office. Price: \$199.99.

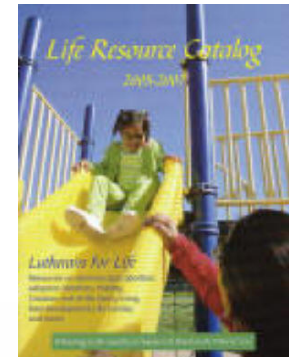
Poker Table
This table is made of high quality materials, built to last and is easy to move. It's a great choice for your home or office. Price: \$199.99.

Poker Table
This table is made of high quality materials, built to last and is easy to move. It's a great choice for your home or office. Price: \$199.99.

Poker Table
This table is made of high quality materials, built to last and is easy to move. It's a great choice for your home or office. Price: \$199.99.

Poker Table
This table is made of high quality materials, built to last and is easy to move. It's a great choice for your home or office. Price: \$199.99.

Poker Table
This table is made of high quality materials, built to last and is easy to move. It's a great choice for your home or office. Price: \$199.99.



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



FFX 115r

MELLO DECIDES NOT 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 41-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 business and that being mayor would mean at least 10 hours per week.

German - Basic

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

Woh! über Tal und Hügel, Er

FFX121-D
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 for help.

The situation is stark.

"People are going to starve to death," warned Fremont Mayor Bob Hall.

The Fremont effort is focused on the area around Lindsay, is of California's largest olive processing plant and, according to city letterhead, the "heart of the California Orange Area."

The severe winter wiped out the winter crop throughout the Lindsay, halfway between Fresno and Bakersfield, is running at 25 percent. Work will not be available in the area, including there until November.

On May 5, a Farmworker Relief Convoy from Fresno, San Francisco, Alameda, and Richmond delivered 12 tons of flour, two tons of rice, two tons of beans and hundreds of pounds of sugar, cereal, cooking oil and used clothing.

GEDI – Java Interface

Image Window

00010005.TIF - DocLib GoundTruthing - Editor and Document Interface (DL-GEDI)

File Edit View Modify Preferences Window Help

Merge Split Save Open Scale: Fit To Window Drag (743, 886)

Current File: 00010005.xml*

Name	Image	Xml
00010003	✓	✓
00010004	✓	✓
00010005	✓	✓
00010006	✓	✓
00010007	✓	✓
00010008	✓	✓
00010009	✓	✓
00010010	✓	✓
00010011	✓	✓

NAME	COLOR	KEY	VISIB...	COU...
DLStamp	Red	None	✓	1
DLSignature	Black	None	✓	1

Attribute	Value
gedi_type	DLStamp
(row,col)(width,hei...	(379,7)(452,300)
Overlap	partial
Quality	poor
Shape	elliptic

Selected Zone Info
(379,7)(831,307)
DLStamp

Browser Window

Type Window

Attribute Window

GEDI

- 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)



Software Architecture

- **Efficient Technology Transfer**
 - **software compatibility**
 - **balance of academia, government, 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



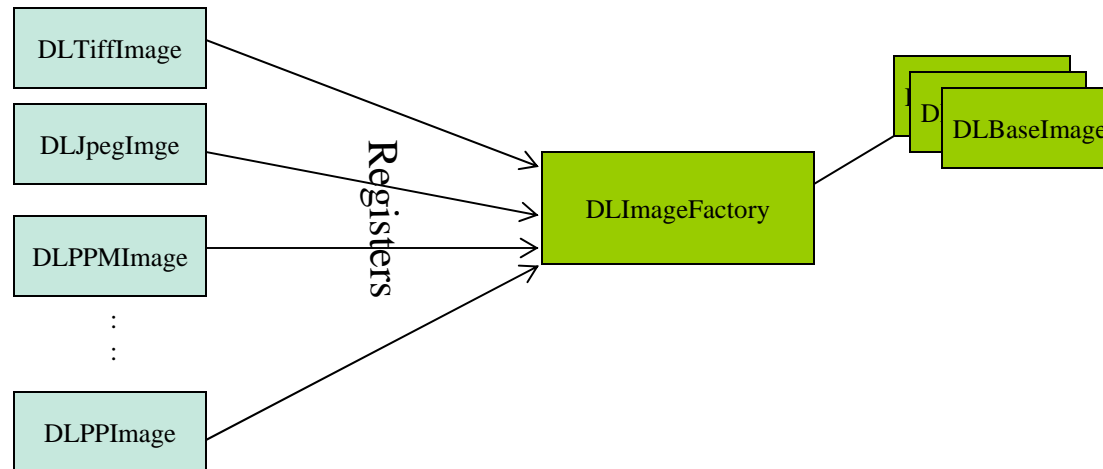
Core vs Add-ons

- 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.

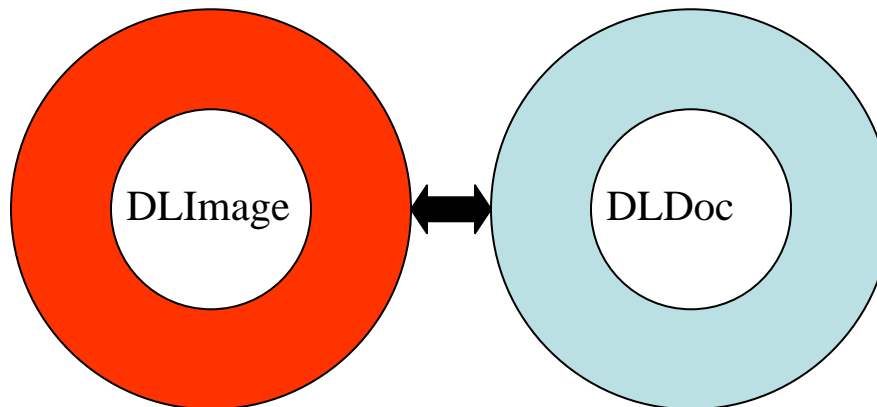


DocLib Architecture

DocLib's architecture rests on two pillars:

DLImage:

➤ **Image Processing**



DLDocument:

➤ **Document Processing**

e.g.

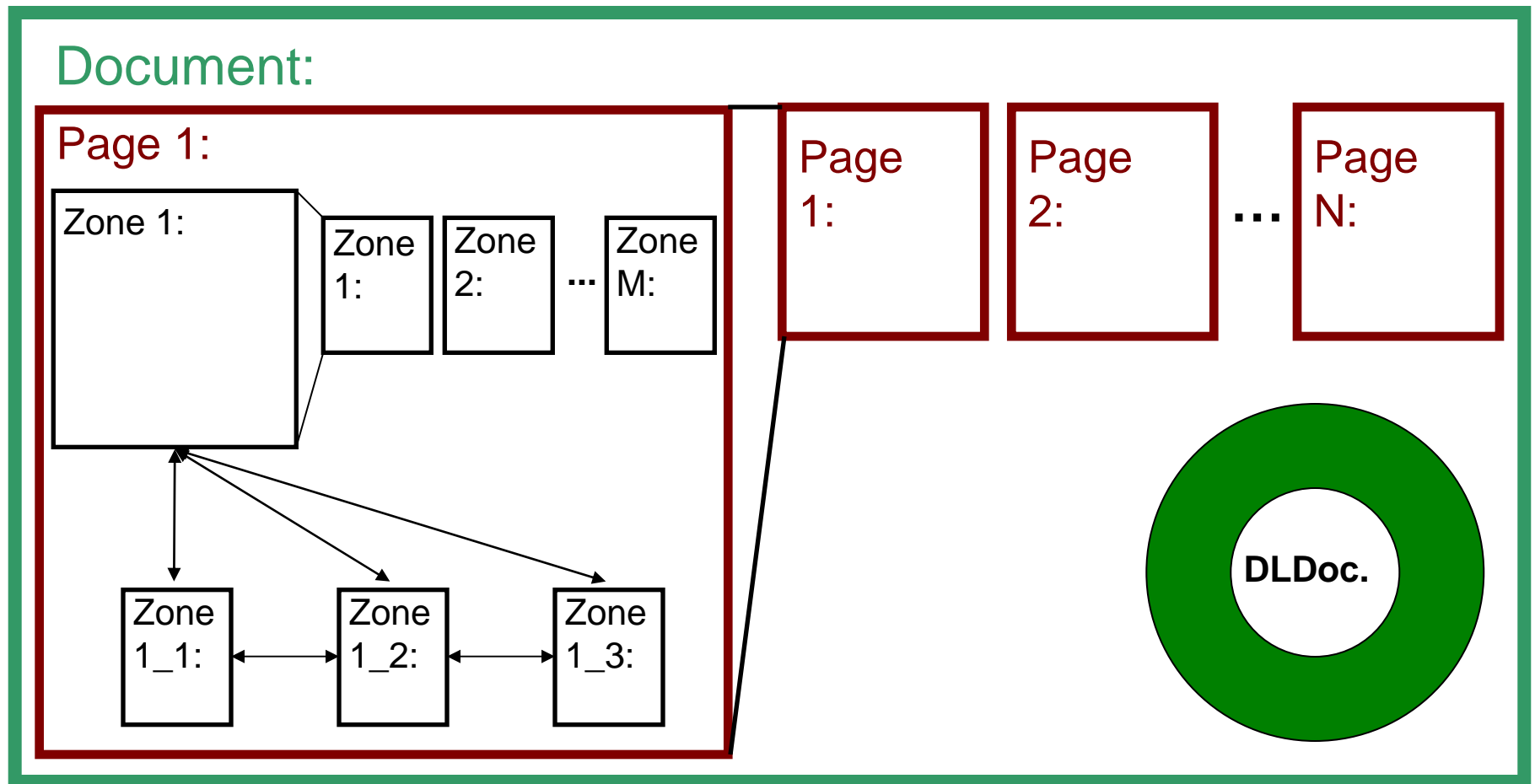
- **image rotation**
- **image deskewing**
- **image conversions**
- **cc calculation**
- **shape drawing**

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



Logo and Stamp Detection



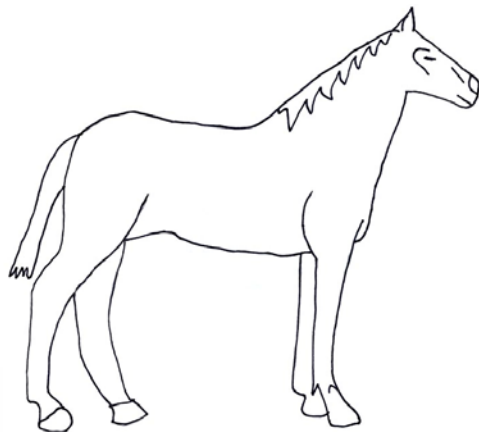
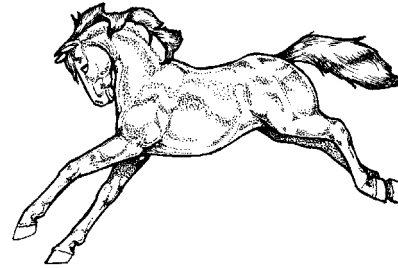
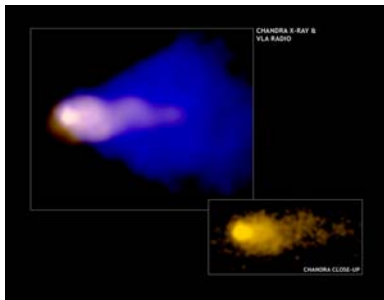
Document/Non-Document Discrimination

- Dataset
 - documents: 422
 - non-documents: 570

	Doc	Non-Doc
Doc	0.9031	0.0969
Non-Doc	0.0941	0.0941



Non-Documents



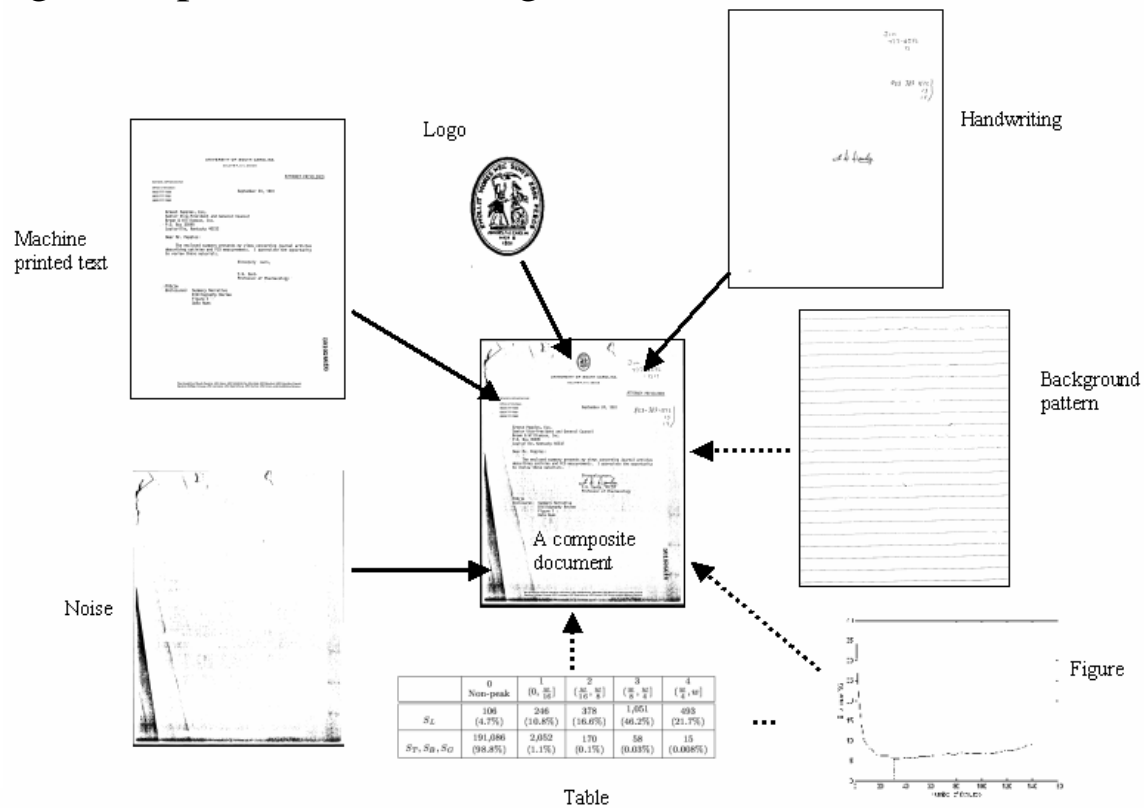
Original Challenges

- Optical Character Recognition (OCR) does not work well in handwritten characters
- Word segmentation in handwritten is challenging
 - Words often touch each other
- Large variation exists even for the handwriting from the same person
- Handwriting often mixes with noise
 - Underground lines
 - Noise introduced by binarization

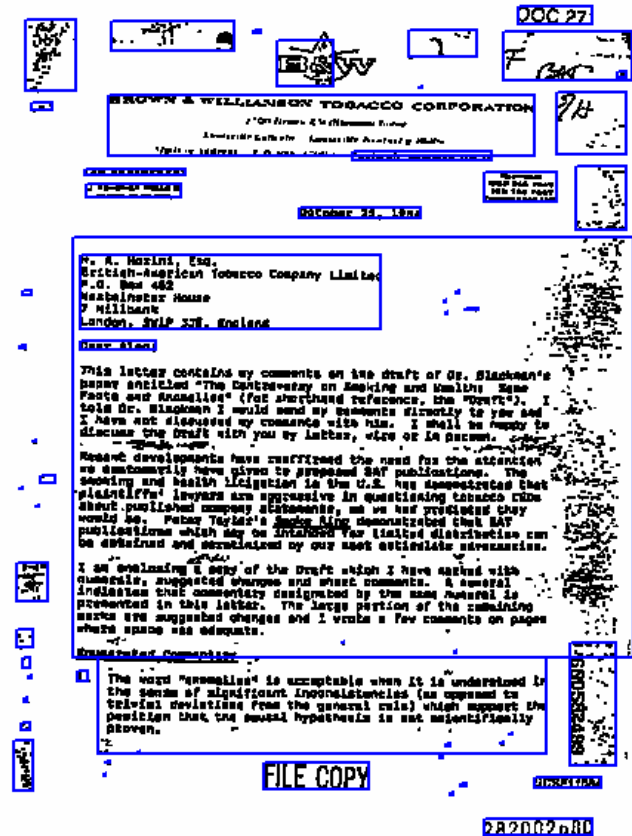
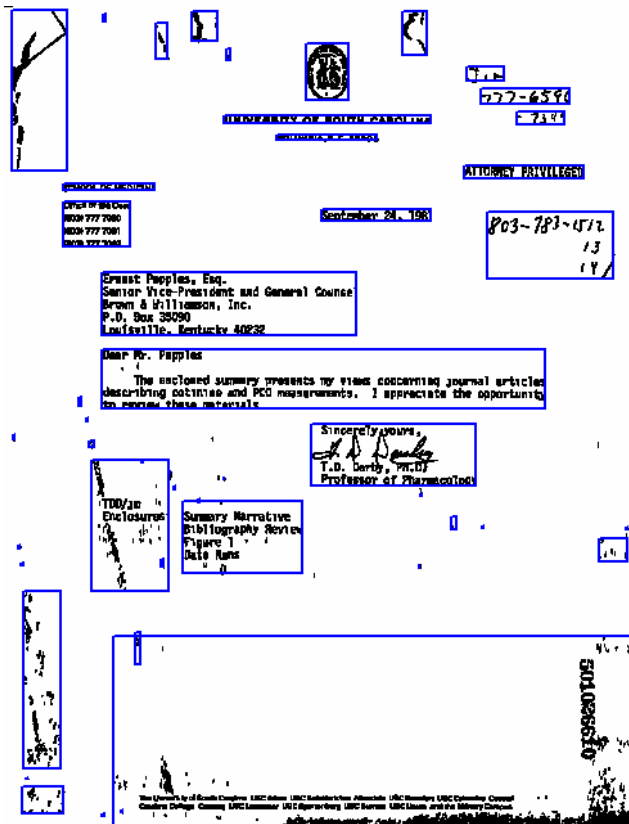


Introduction

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



Page Segmentation for Noisy Documents



* Docstrum page segmentation technique is used



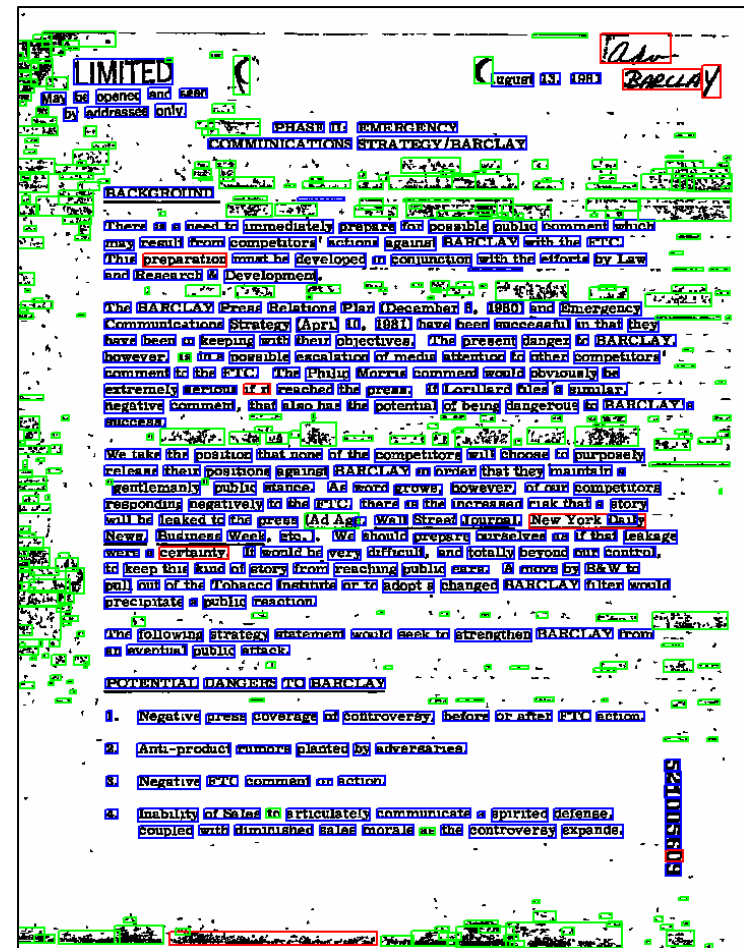
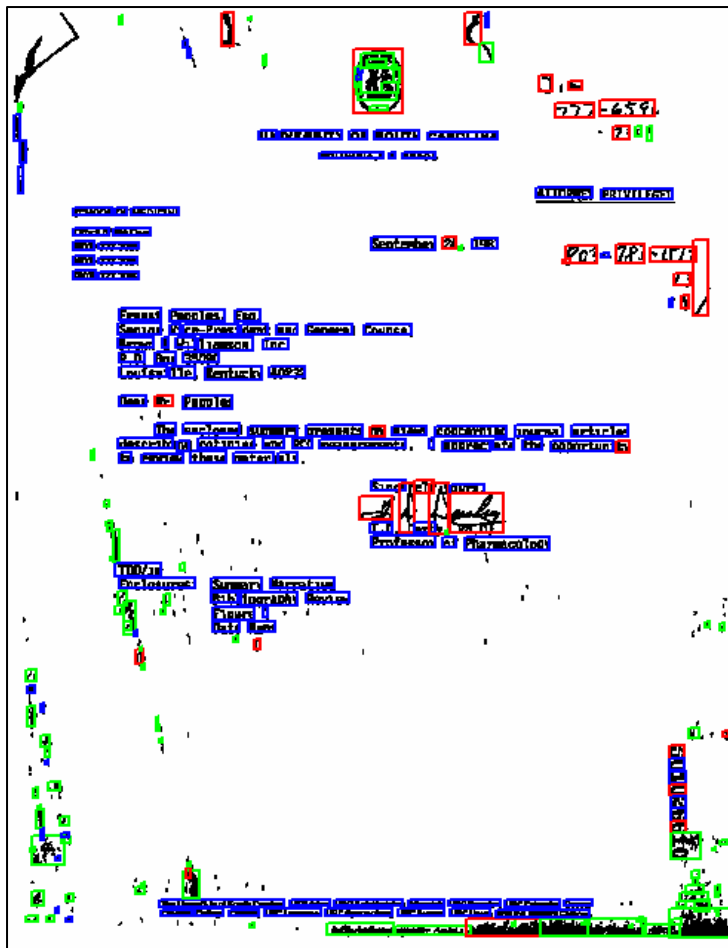
Overview of Our Approach

- 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.



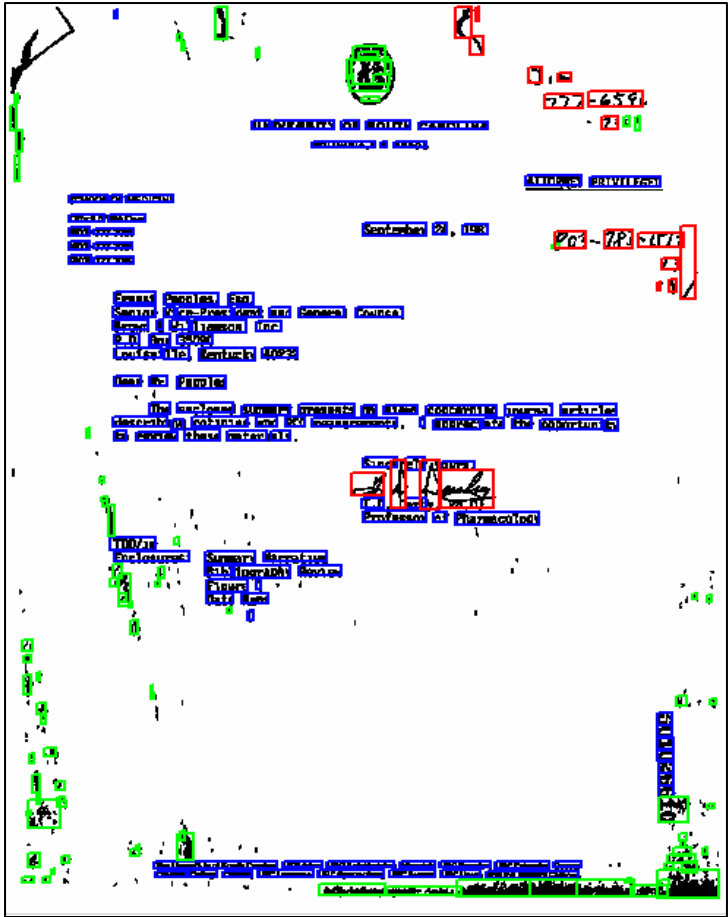
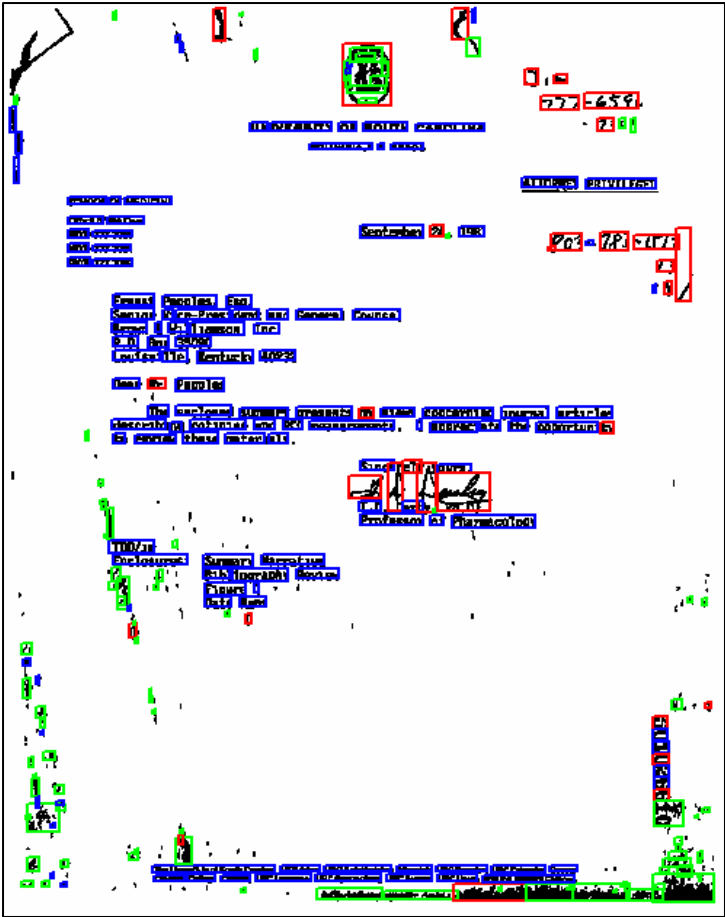
Classification Results with Fisher Classifier

Printed text
Handwriting
Noise



MRF Postprocessing Example

Printed text
Handwriting
Noise



Before MRF-based postprocessing

After MRF-based postprocessing



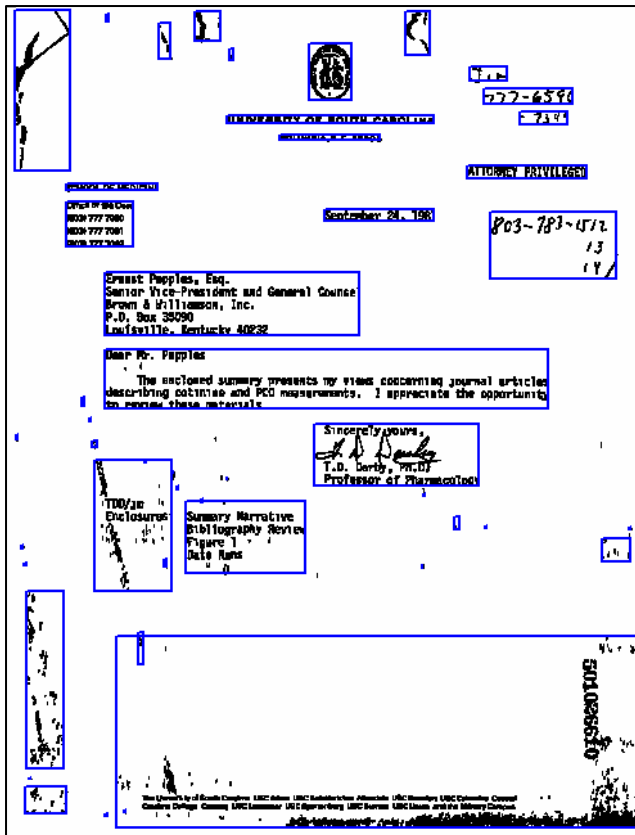
Evaluation

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

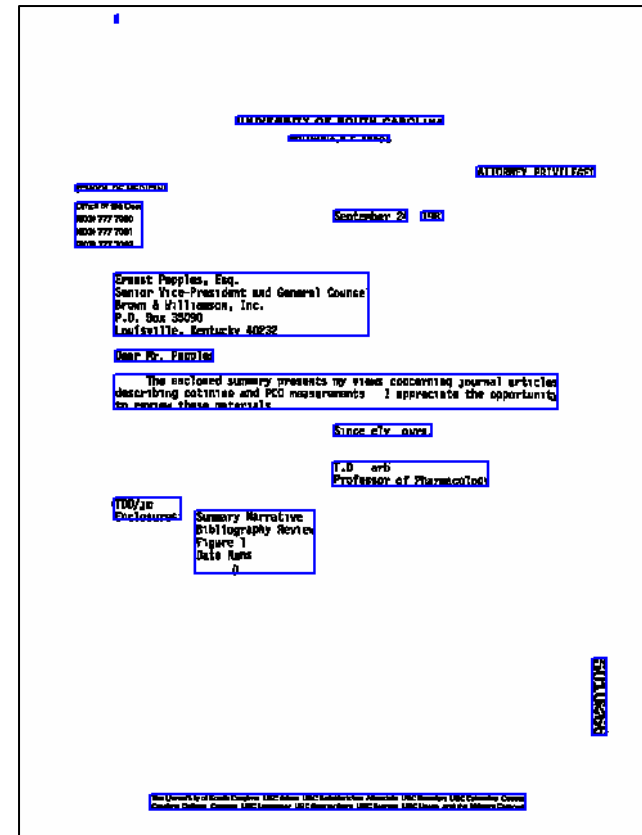
	#Total	Percentage	Before Post-processing		After Post-processing	
			Accuracy	Precision	Accuracy	Precision
Printed Words	19,227	66.9%	95.9%	99.5%	98.0%	99.7%
Handwritten Words	701	2.4%	93.2%	62.9%	93.0%	83.3%
Noise Blocks	8,802	30.7%	96.8%	93.0%	98.6%	96.0%
Total	28,730	100%	96.1%	N/A	98.1%	N/A



Application to Page Segmentation



Before enhancement



After enhancement



Handwritten Line Detection

Handwritten Arabic text, likely a medical or scientific report, with some lines highlighted in red. The text is dense and appears to be a detailed account or analysis.

ATTORNEY CERTIFICATE RE: THE CONTENTS OF BLOOD VALUES AND ALVEOLAR CARBON DIOXIDE (PACO₂) VALUES TO DETERMINE CIGARETTE CONTENT DELIVERED PER CIGARETTE SMOKED

Summary

Basic to a review of data related to concentration measurements of chemical substances in body fluids is an understanding of the kinetics of of uptake, absorption, and elimination. Inhalation of gases and particulate matter can cause absorption, at the alveoli, into the blood. Blood solubility is an important factor in determining the rate of equilibration between inhaled air concentration and total body concentration as represented by blood values. Absorption into organs, tissues, and fluids of the body will affect blood values. Elimination by excretion and metabolism also affects blood values of chemical substances.

The following table lists the major factors which impact on the accuracy of determinations of nicotine content delivered per cigarette based upon cigarette blood values and alveolar carbon dioxide values and used to determine the nicotine content.

Table 1
FACTORS AFFECTING THE REPRESENTATION

1. Sensitivity and accuracy of the method used for determination of cigarette blood values;
 2. Cigarette biological half-life;
 3. Number of cigarettes smoked per day;
 4. Fluctuation in daily cigarette smoking patterns;
 5. Fluctuation pattern of the individual smokers used in the study;
 6. Fluctuation in number of cigarettes smoked per day;
 7. Inter-subject/inter-day variation among study group subjects;
 8. Sex of the subject;
 9. Consistency of alveolar carbon dioxide metabolism to nicotine.
- Factors 1, 2, 4, and 5, listed above, can be controlled or a correction can be made based upon alveolar carbon dioxide (PACO₂).

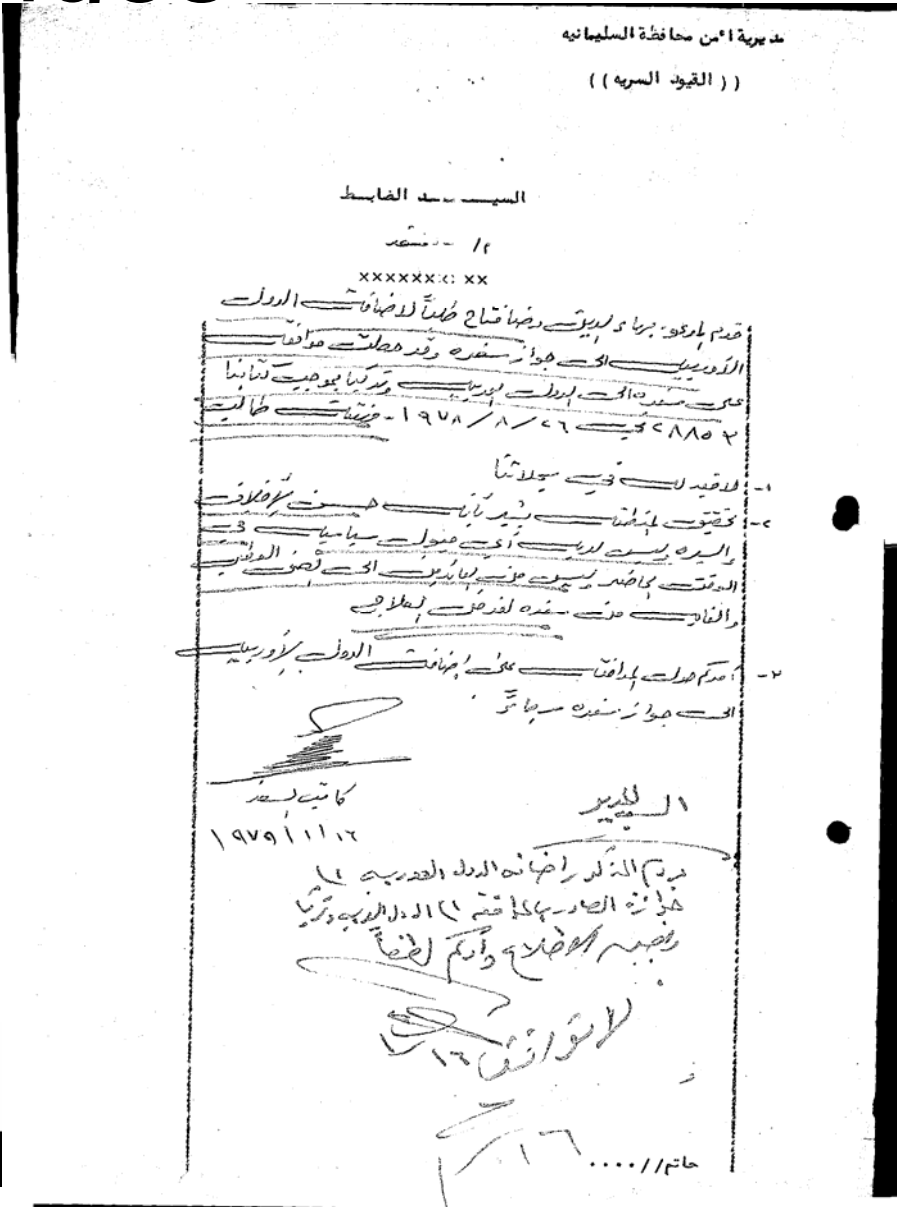


UNIVERSITY



Challenges

- Handwriting
 - Curvilinear
 - Touching lines, words, characters
 - Non-Manhattan layout
 - Multiple orientation
- Machine Printed
 - Straight
 - Significant Gap
 - Regular layout



Our Assumptions

- Text Lines are generally horizontal
- Text line orientation may vary slightly
- There are gaps between neighboring lines
 - Strong for machine printed
 - Weak for handwriting



Intuitive Approach

- Preprocessing
 - Gaussian Blurring using rectangular window
- Segmentation: The Level Set Methods
 - Enforce merging in horizontal direction
- Localization
 - Group nearby isolated connected components



The level set method

$$\frac{\partial f}{\partial t} + \underbrace{\vec{S} \cdot \nabla f}_{\text{Vector Field Based}} + \underbrace{V_N |\nabla f|}_{\text{In Normal Direction}} = \underbrace{b\kappa |\nabla f|}_{\text{Curvature Based}}$$

http://vision.ece.ucsb.edu/~sumengen/level_set_methods

- The moving speed df/dt is changed by forces based on curvature, normal direction or external vector field



Evaluation

- Ground truth
 - VIPER Ground Truth Editor
 - Pixel level: use polygon to represent text lines
- Evaluation method
 - The Hungarian algorithm
 - One-to-one correspondence
- Statistics
 - Four scripts, 100 documents per script
 - 5000 text lines



Comparison

QUANTITATIVE COMPARISON OF OUR APPROACH AND A CONNECTED COMPONENT BASED APPROACH ON 100
HANDWRITTEN HINDI DOCUMENTS.

	Hit Rate	STD of Hit Rate	Detected Text Lines (Total: 1,365)
Our Method	97%	0.02	1,295 (95%)
Projection Profile based method	43%	0.37	673 (49%)
Improved connected component	78%	0.19	1,103 (80%)
Docstrum based method	74%	0.22	996 (72%)

QUANTITATIVE COMPARISON OF OUR APPROACH AND A CONNECTED COMPONENT BASED APPROACH ON 100
HANDWRITTEN CHINESE DOCUMENTS.

	Hit Rate	STD of Hit Rate	Detected Text Lines (Total: 1,672)
Our Method	98%	0.01	1,532 (92%)
Projection Profile based method	58%	0.24	965 (57%)
Improved connected component	64%	0.12	1,131 (67%)
Docstrum based method	94%	0.04	1,389 (83%)



Comparison (cont')

QUANTITATIVE COMPARISON OF OUR APPROACH AND A CONNECTED COMPONENT BASED APPROACH ON 100

HANDWRITTEN WITH BACKGROUND PARALLEL LINES.

	Hit Rate	STD of Hit Rate	Detected Text Lines (Total: 1,589)
Our Method	98%	0.01	1,466 (93%)
Projection Profile	79%	0.10	1,178 (74%)
Improved connected component	54%	0.34	849 (53%)
Docstrum based method	62%	0.29	1,024 (64%)

QUANTITATIVE COMPARISON OF OUR APPROACH AND A CONNECTED COMPONENT BASED APPROACH ON 100 FREESTYLE

HANDWRITTEN DOCUMENTS.

	Hit Rate	STD of Hit Rate	Detected Text Lines (Total: 2,178)
Our Method	93%	0.03	1,863 (85%)
Projection Profile	72%	0.16	1,537 (70%)
Improved connected component	81%	0.13	1,576 (73%)
Docstrum based method	82%	0.15	1,751 (80%)



Robustness Test

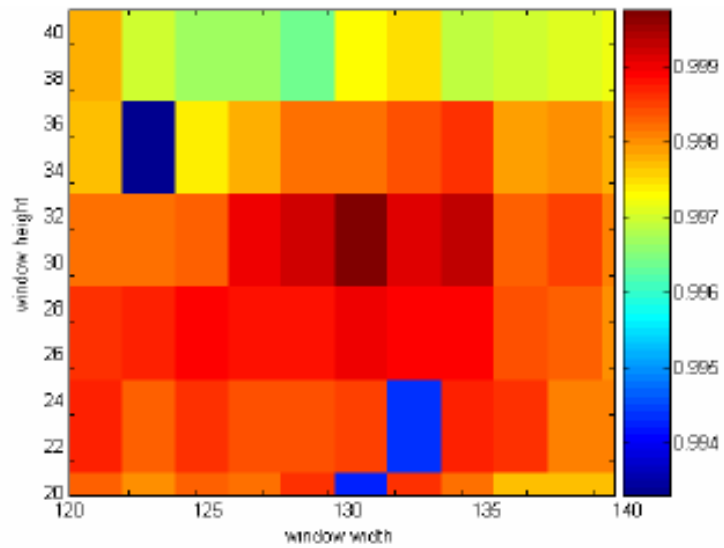
- Change of internal parameter
- Change of input scale
- Rotate input image
- Corrupt input with Noise

Handwritten notes in Chinese, likely a transcription of the slide content, written in a cursive style.

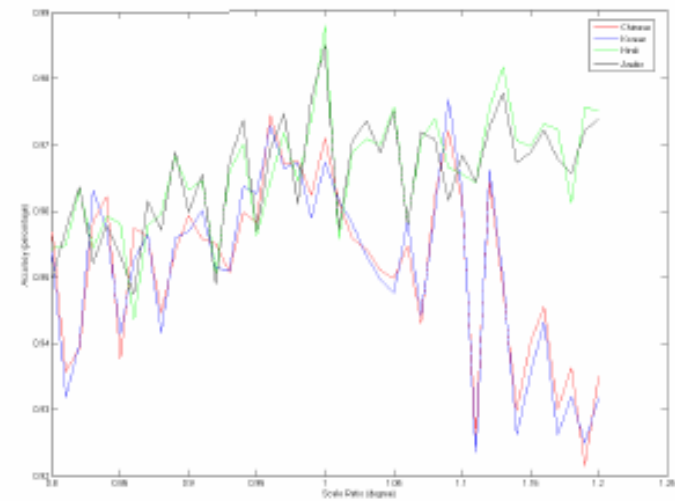
Handwritten notes in Chinese, similar to the first image, with some text highlighted in red.

Handwritten notes in Chinese, similar to the first image, with some text highlighted in red.

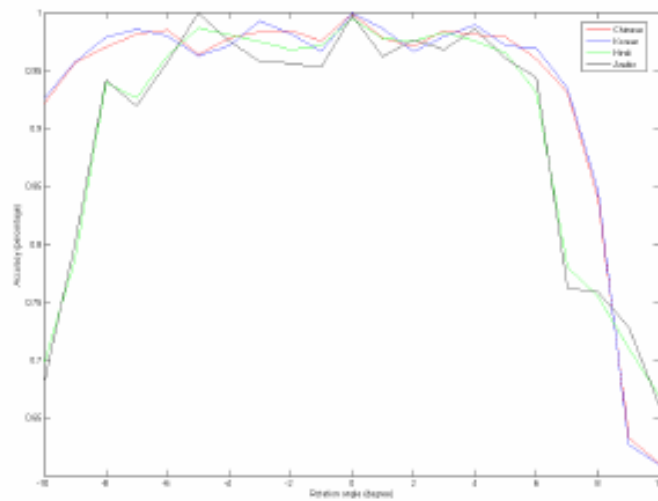
Handwritten notes in Chinese, similar to the first image, with some text highlighted in red.



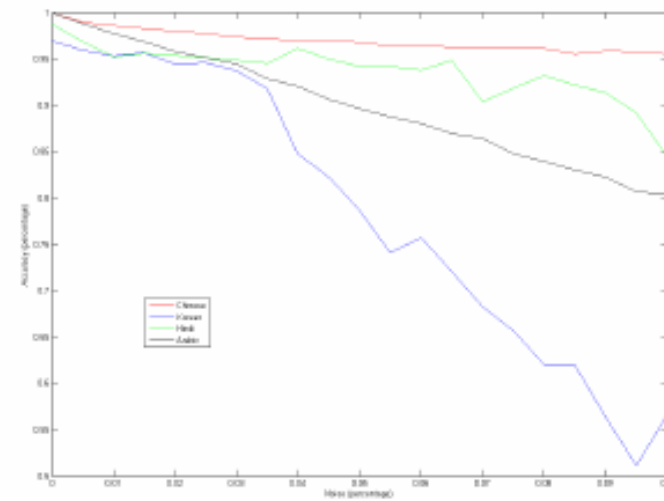
(a)



(b)



(c)



(d)



Fig. 6. Robustness test. (a) Accuracy of different parameters. (b) Accuracy of different scales. (c) Accuracy of different rotation angles. (d) Accuracy of different noises

Publications

- Yi Li, Yefeng Zheng and David Doermann. Detecting Text Line in Handwritten Documents. ICPR'06,
- Yi Li, David Doermann, Stefan Jaeger, and Yefeng Zheng . A new algorithm and its implementation in Detecting Text Line in Handwritten Documents. IWFHR'06, 2006.
- Yi Li, Yefeng Zheng, Stefan Jaeger, and David Doermann. A new algorithm and its evaluation in Detecting Text Line in Handwritten Documents. UMD Technical Report. (in preparation)



DVD

- Dataset without Ground Truth
- Slides from Today



Agenda

Project Overview

- Introduction
- Goals and Objectives

Data Collection and Ground Truth

GEDI and Evaluation Framework

Evaluation and Research Components

- Unconstrained Signature and Logo Detection and Matching for Off-line Document Image Retrieval
- Document/Non-Document Discrimination

Technical Presentations

- **Adaptive OCR**
- **Document Classification by Layout**
- **Learning from Web for Shape-based Object Recognition**

