# Reconocimiento de Escritura
## Lecture 5/5
## More Layout Analysis and Applications

Daniel Keysers

Jan/Feb-2008

UNIVERSIDAD POLITECNICA DE VALENCIA

# Outline

UNIVERSIDAD POLITECNICA DE VALENCIA

# Outline

UNIVERSIDAD
POLITECNICA
DE VALENCIA

# Page Frame Detection

- motivation
  - textual and non-textual noise
  - textual noise results in OCR errors
- idea: detect page contents area
- RAST for page frame detection
  - solves the problem in a general framework
  - robust against the amount of noise
  - robust against overlapping noise
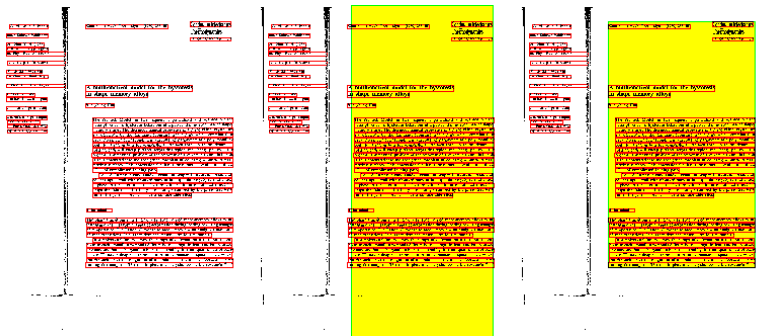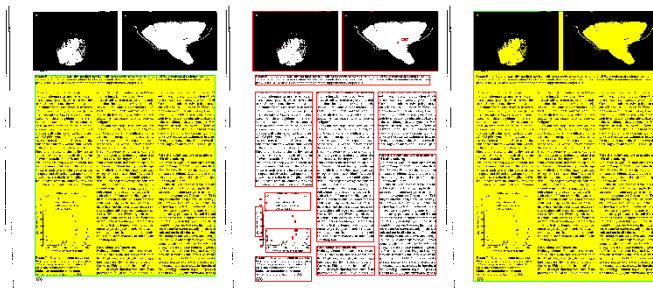
UNIVERSIDAD POLITÉCNICA DE VALENCIA

# Page Frame Detection Method

- two-phase approach
- (1) determine left and right side ($\rightarrow$ RAST)
  - use text-lines
  - the quality function has two parts:
    - the left and right border should have many text-line ends on the inside of the page frame
    - but they should not have many text-line ends on the outside of the page frame
  - use soft term (bounded error) of the form $\max(0, 1 - d^2/\epsilon^2)$
- (2) determine upper and lower side
  - include all character bounding boxes in the range
  - adjust for page numbers and images

# Illustration of the Two Steps

# Inclusion of Images and Page Numbers



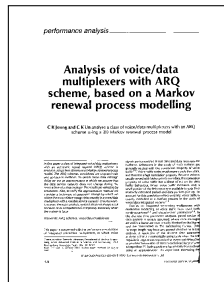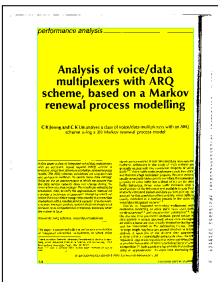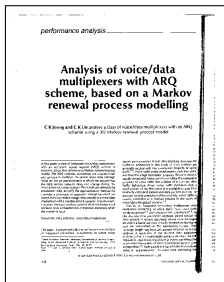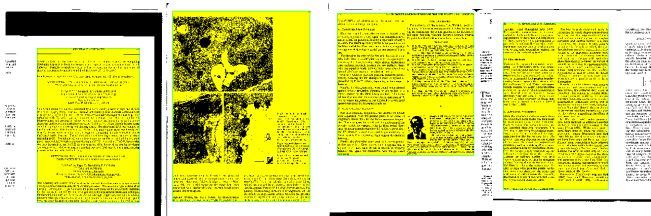center image: zones detected by Voronoi algorithm

# Page Frame Detection Results

| Performance measure | Error rate (%) |
|---|---|
| Area overlap | 4.0 |
| Connected components classification | 1.6 |
| Ground-truth zone detection | 2.8 |

| Application | Error rate (%) | |
|---|---|---|
| | No PFD | With PFD |
| OCR | 4.3 | 1.7 |
| Layout based document image retrieval | 7.0 | 5.4 |

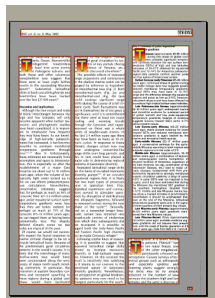evaluated on the UW-III dataset with every tenth document used as training set

UNIVERSIDAD POLITECNICA DE VALENCIA

# Examples of Page Frame Detection

# Robustness against Noise

Noise ratio = outside pixels / total pixels

# Example of OCR Difference



Screen-shot of Omnipage 14 showing the recognized text of the original document (left) and the document cleaned using page frame detection (right). Note that the reading order of the text has changed, probably due to the slightly changed geometry.

UNIVERSIDAD
POLITÉCNICA
DE VALENCIA

# Outline

UNIVERSIDAD POLITECNICA DE VALENCIA

# Adaptation to Urdu/Arabic Script

"We are often reminded that English is blessed with one of the simplest scripts in the world." (Nagy 2000)

UNIVERSIDAD
POLITECNICA
DE VALENCIA

# Motivation for Urdu Document Analysis

- Urdu has more than 150 million speakers
- written in Arabic script with above 20,000 ligatures

قوم کے لیے اپنے بے ہُنر ہاتھوں سے ایک آئینہ خانہ بنایا ہے جس میں آکر وہ اپنے خط و خال
دیکھ سکتے ہیں کہ ہم کون تھے اور کیا ہوگئے۔ اگرچہ اس جانگاہ نظم میں جس کی دُشواریاں لکھنے والے
کا دل اور دماغ ہی خوب جانتا ہے، بیان کا حق نہ مجھ سے ادا ہوا ہے اور نہ ہو سکتا ہے مگر شکر ہے

- no Urdu OCR and layout analysis system
- large potential market

# Diagram of Approach

UNIVERSIDAD POLITÉCNICA DE VALENCIA

# Approach for Urdu Documents

1. Find empty whitespace rectangles that completely cover the page background.
2. The whitespace rectangles are evaluated as candidates for column separators or gutters based on their aspect ratio, width, and proximity to text-sized connected components.
3. Find text-lines that respect the columnar structure of the document.
4. Determine the reading order of the text-lines using constraints on the geometric arrangement of text-line segments on the page. (Change left-to-right model to right-to-left model.)

UNIVERSIDAD POLITECNICA DE VALENCIA

# Preprocessing

# Column Separators

The whitespace rectangles are evaluated as candidates for column separators based on the following constraints:

1. Column-separating rectangles must have an aspect ratio of at least 1:3

2. Column-separating rectangles must have a width of at least 1.5 times of the mode of the distribution of widths of inter-word spaces.

3. Column-separating rectangles must be adjacent to at least four character-sized connected components on their left or their right side.

# Constrained Text-Line Finding

Use Column Separators as Obstacles
adapt Roman script text-line model for Urdu/Arabic:
use two descender lines



Base-line

First-descender-line

Second-descender-line

# Urdu Document Layout Analysis Examples

Note that images and graphics were not removed here, so they result in some spurious text-lines.

# Text-Line Detection Accuracy Results

25 images of Urdu text from different sources
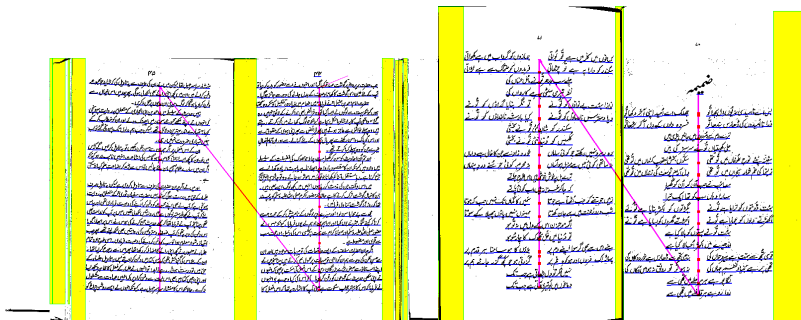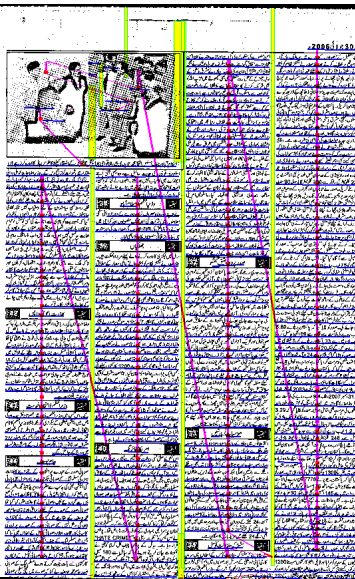five classes: *book*, *poetry*, *digest*, *magazine*, and *newspaper* with 5
images each

| Layout ($n =$) | Correct | Split | Merged | Missed |
|---|---|---|---|---|
| Book (234) | **91.45** | 4.27 | 0.00 | 4.27 |
| Poetry (286) | **92.31** | 4.55 | 0.00 | 3.15 |
| Digest (702) | **80.63** | 11.54 | 0.00 | 7.84 |
| Magazine (1158) | **90.07** | 4.14 | 0.86 | 4.75 |
| Newspaper (819) | **72.16** | 7.81 | 4.15 | 15.87 |

UNIVERSIDAD POLITECNICA DE VALENCIA

# Outline

UNIVERSIDAD POLITECNICA DE VALENCIA

# Block Types

**math**

$$\omega_1 = \Psi_\leftrightarrow \Rightarrow \Psi_\rightarrow \vee \Psi_\leftarrow$$
$$\omega_2 = \Psi_\rightarrow \wedge \Psi_\leftarrow \Rightarrow \Psi_\leftrightarrow$$
$$\omega_3 = \Psi_\rightarrow \wedge \Psi_\leftrightarrow \Rightarrow \Psi_\leftarrow$$
$$\omega_4 = \Psi_\leftrightarrow \Rightarrow \Psi_\leftarrow \vee \Psi_\leftrightarrow$$
$$\omega_5 = \Psi_\rightarrow \wedge \Psi_\leftarrow \wedge \Psi_\leftrightarrow \Rightarrow false$$

$$M_{\mathbf{t}}^{\hat{\mathbf{t}}} \simeq \langle M_{\infty}^{\hat{\mathbf{0}}} \frac{u_t^{\hat{\mathbf{t}}}}{u_\infty^{\hat{\mathbf{t}}}} \rangle \Big/ \frac{T_{\mathbf{t}}}{T_\infty} , \qquad (2.3)$$

$$_3 Y^{ud} =_4 C^{vdev} {}_2 q^u_{uv} . \qquad (28)$$

**text**

signals are transmitted in real-time and data messages are buffered. Difficulties in the study of such systems are generally related with the correlation property of voice traffic[1,2]. Voice traffic varies much more slowly than data, and thus has a high correlation property. Because voice is usually served with higher priority over data, this correlation property of voice traffic has a direct effect on the data buffer behaviour. When voice traffic increases, only a small portion of the link capacity is available to data for a relatively extended period and data packets pile up. To account for this correlation effect properly, voice traffic is usually modelled as a Markov process in the study of voice/data integrated systems[4-6].

Upon graduation (with whatever degree), the young engineer presumptive could expect to work some years under the direction of an experienced engineer (who would continuously critique the performance and output) and in the presence of other engineers of varying experience and knowledge. Together with the neophyte's own learning from the literature and feedback from the plant itself (I mean *real* hardware), this mentoring process could produce (under ideal conditions) outstanding engineers and identify the best career paths for those engineers to pursue.
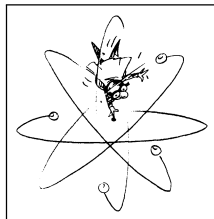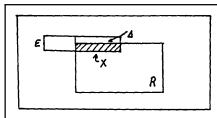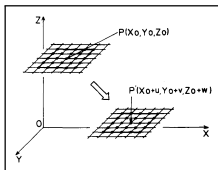
307

table







drawing

# Block Types



logo



halftone

UNIVERSIDAD POLITECNICA DE VALENCIA

# Block Types

ruling

speckles

# Related Work

| reference | # pages | # zones | # types | error [%] |
|---|---|---|---|---|
| Inglis and Witten, 1995 | 1001 | 13831 | 3 | 6.7 |
| Liang et al., 1996 | 979 | 13726 | 8 | 5.4 |
| Sivaramakrishnan et al., 1995 | 979 | 13726 | 9 | 3.3 |
| Wang et al., 2000 | 1600 | 24177 | 9 | 2.5 |
| Wang et al., 2006 | 1600 | 24177 | 9 | 1.5 |
| this work | 713 | 13811 | 8 | 1.5 |

features used in the literature:

connected components, run lengths, cross-correlation between scan-lines, vertical projection profiles, wavelet coefficients, learned masks, black pixel distribution

UNIVERSIDAD POLITECNICA DE VALENCIA

# Data Set Used

- University of Washington III (UW-III) database
- 1600 English document images
- 24177 homogeneous manually labeled page segments/blocks
- degradation types: direct scan, photocopies $\rightarrow$ duplicates
- avoid duplicates here $\rightarrow$ 713 documents
- speckles not annotated but important $\rightarrow$ automatic extraction

UNIVERSIDAD POLITÉCNICA DE VALÈNCIA

# Features

- Tamura texture feature histogram (TTFH)
- relational invariant feature histograms (RIFH)
- down-scaled image $32 \times 32$ (DSI)

- number, mean, and variance of run-lengths (RL{B,W}{X,Y,M,S}V)
- run-length histograms (RL{B,W}{X,Y,M,S}H)
- connected components size histograms (CCXH, CCYH, CCXYH)
- connected components nearest neighbor histogram (CCNNH)
- fill ratio after horizontal smearing (FR)

UNIVERSIDAD POLITECNICA DE VALENCIA

# Classification

- nearest neighbor with leaving-one-out cross-validation
- Jensen-Shannon divergence for histograms,
  Euclidean distance for other features
- weights proportional to inverse of error rate

- for fast and small classifier
  - log-linear classifier using maximum entropy criterion
  - 50/50 split of data for evaluation

## Experimental Results — Single Features

| feature | dim. | extr. [s] | error [%] |
|---|---|---|---|
| TTFH | 512 | 5.51 | **3.4** |
| RIFH | 512 | 12.59 | 7.8 |
| DSI | 1024 | 0.01 | 8.1 |
| FR | 1 | 0.02 | 27.3 |
| CCXH | 8 | 0.04 | 14.5 |
| CCYH | 8 | 0.04 | 14.9 |
| CCXYH | 64 | 0.04 | 6.2 |
| CCNNH | 8 | 0.05 | 19.0 |

| feature | dim. | extr.[s] | error [%] |
|---|---|---|---|
| RLBXH | 8 | 0.01 | 7.9 |
| RLWXH | 8 | 0.01 | 5.1 |
| RLBYH | 8 | 0.01 | 8.2 |
| RLWYH | 8 | 0.01 | 5.6 |
| RLBMH | 8 | 0.01 | 11.8 |
| RLWMH | 8 | 0.01 | 6.6 |
| RLBSH | 8 | 0.01 | 10.5 |
| RLWSH | 8 | 0.01 | 6.2 |
| RLBXV | 3 | 0.01 | 12.9 |
| RLWXV | 3 | 0.01 | 9.7 |
| RLBYV | 3 | 0.01 | 14.6 |
| RLWYV | 3 | 0.01 | 12.1 |
| RLBMV | 3 | 0.01 | 17.2 |
| RLWMV | 3 | 0.01 | 12.6 |
| RLBSV | 3 | 0.01 | 16.7 |
| RLWSV | 3 | 0.01 | 12.2 |

UNIVERSIDAD POLITECNICA DE VALENCIA

# Experimental Results — Combinations

| feature | error [%] |
|---|---|
| RL**V, constant weight | 4.1 |
| RL**H, constant weight | 1.8 |
| RL*, CC*, 1/error weight | **1.5** |
| FR, RL*, CC*, 1/error weight | 1.5 |
| TTFH, FR, RL*, CC*, 1/error weight | 1.5 |
| RL*, CC*, *logistic, 50/50 data split* | 2.1 |

UNIVERSIDAD POLITECNICA DE VALENCIA

# Confusion Matrix

|          | text | speckles | math | drawing | ruling | table | halftone | logo |
|----------|------|----------|------|---------|--------|-------|----------|------|
| text     | **99.8** |      | .1   |         |        |       |          |      |
| speckles | .5   | **99.4** | .1   | .1      |        |       | .1       |      |
| math     | 8.6  |          | **90.8** |     |        |       | .6       |      |
| drawing  | 3.0  | .3       | 1.5  | **86.0** |       | 5.5   | 3.5      | .3   |
| ruling   | 1.3  | 2.2      | .4   | .4      | **96.1** |    |          |      |
| table    | 20.7 |          | .8   | 9.9     |        | **68.6** | .8    |      |
| halftone |      | 1.8      |      | 9.7     | .9     |       | **86.7** | .9   |
| logo     | 36.4 | 9.1      | 9.1  | 9.1     |        |       | 9.1      | **27.3** |
| frequency | 10450 | 2007    | 476  | 401     | 232    | 121   | 113      | 11   |

# Examples of Misclassifications



speckles



text

$$(P_R \sqsubseteq P_{R'}) \wedge (P \sqsubseteq_F P \llbracket P_R \sqsubseteq_F P \llbracket P_{R'})$$

math

1. $P_R \sqsubseteq P_{R'} \Rightarrow \mathscr{R}(P) \sqsubseteq_F P \llbracket P_{R'}$

text



drawing



halftone



drawing



table

UNIVERSIDAD
POLITÉCNICA
DE VALÈNCIA

# Conclusions for Block Classification

- use run-lengths histograms
- background run-lengths more important than foreground
- very competitive error rate of 1.5% using simple features
- simple, fast and accurate classifier at 2.1%
  - run-lengths and connected components distribution
  - maximum entropy log-linear classifier
- probable improvement: use context information

# Outline

UNIVERSIDAD
POLITECNICA
DE VALENCIA

## Document Reflow



T.M. Breuel, W.C. Janssen, K. Popat, H.S. Baird:
"Paper-to-PDA," Procs. ICPR 2002, Quebec City, Quebec, Canada

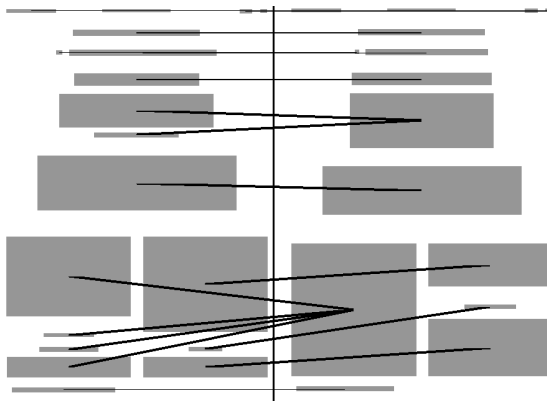UNIVERSIDAD POLITÉCNICA DE VALENCIA

# Outline

UNIVERSIDAD POLITECNICA DE VALENCIA

# Visual Document Search



▶ visual search feature complementing text-based search

# Layout-Based Search



▶ edge cover distance measure for document images

UNIVERSIDAD
POLITÉCNICA
DE VALENCIA

# Example-Based Labeling of Title Page Images

- Labeling of title pages using similar labeled examples
1. Segment document image using layout analysis
2. Search for similar labeled documents in dataset using geometrical and textural features
3. Copy the labels from the best matching document

# Example-based Logical Labeling



- automatic semantic labeling of page segmentations
- example-based approach: match blocks and transfer labels
- extend block distance using texture features
- accuracy on MARG:
  99.6% – unknown document
  98.9% – unknown journal
  94.8% – unknown journal type

UNIVERSIDAD POLITÉCNICA DE VALENCIA

# Outline

UNIVERSIDAD POLITECNICA DE VALENCIA
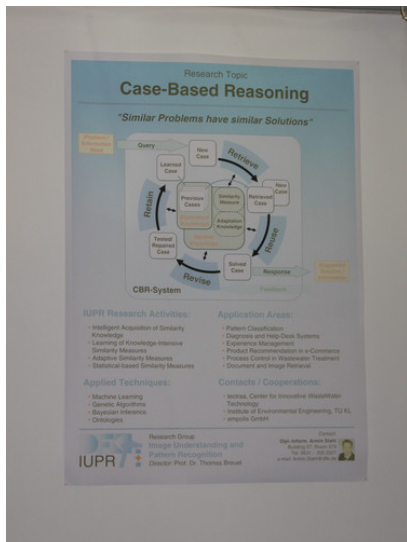
# iDesk

UNIVERSIDAD
POLITECNICA
DE VALENCIA

# iDesk User Interaction
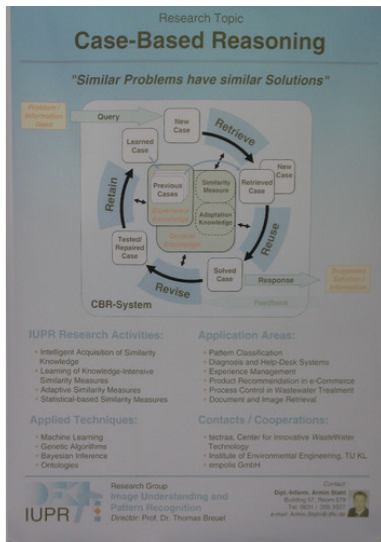


improved user experience:

▶ pointing at region of interest possible

▶ document detection improved

▶ works without separate calibration step

▶ zooming enabled

input

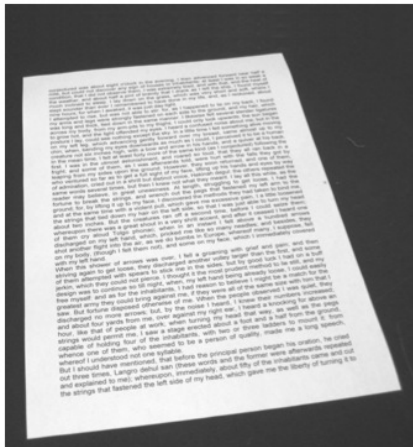output

# OSCAR - One Step Capture and Restoration

UNIVERSIDAD POLITECNICA DE VALENCIA

# Outline

UNIVERSIDAD
POLITECNICA
DE VALENCIA

# Page Surface Dewarping

# Page Surface Dewarping



by my good master, Mr. Pannel, commander; with

UNIVERSIDAD POLITECNICA DE VALENCIA

# Page Surface Dewarping



During the whole of a dull, dark, and soundless day in the autumn of the year, when the clouds hung oppressively low in the heavens, I had been passing alone, on horseback, through a singularly dreary tract of country, and at length found myself, as the shades of evening drew on, within view of the melancholy House of Usher. I know not how it was -- but, with the first glimpse of the building, a sense of insufferable gloom pervaded my spirit. I say insufferable; for the feeling was unrelieved by any of that half-pleasurable, because poetic sentiment, with which the mind usually receives even the sternest natural images of the desolate or terrible. I looked upon the scene before me -- upon the mere house, and the simple landscape features of the domain -- upon the blank walls -- upon the vacant eye-like windows -- upon a few rank sedges -- and upon a few white trunks of decayed trees -- with an utter depression of soul which I can compare to no earthly sensation more properly than to the after-dream of the reveller upon opium -- the bitter lapse into every-day life -- the hideous dropping off of the veil. There was an iciness, a sinking, a sickening of the heart -- an unredeemed dreariness of thought which no goading of the imagination could torture into aught of the sublime. What was it -- I paused to think -- what was it that so unnerved me in the contemplation of the House of Usher? It was a mystery all insoluble; nor could I grapple with the shadowy fancies that crowded conclusion, that while, beyond doubt, there are combinations of very simple natural objects which have the power of thus affecting us, still the analysis of this power lies among considerations beyond our depth. It was possible, I reflected, that a mere different arrangement of the particulars of the scene, of the details of the picture, would be sufficient to modify, or perhaps to annihilate its capacity for sorrowful impression; and, acting upon this idea, I reined my

During the whole of a dull, dark, and soundless day in the autumn of the year, when the clouds hung oppressively low in the heavens, I had been passing alone, on horseback, through a singularly dreary tract of country, and at length found myself, as the shades of evening drew on, within view of the melancholy House of Usher. I know not how it was -- but, with the first glimpse of the building, a sense of insufferable gloom pervaded my spirit. I say insufferable; for the feeling was unrelieved by any of that half-pleasurable, because poetic sentiment, with which the mind usually receives even the sternest natural images of the desolate or terrible. I looked upon the scene before me -- upon the mere house, and the simple landscape features of the domain -- upon the blank walls upon the vacant eye-like windows -- upon a few rank sedges and upon a few white trunks of decayed trees -- with an utter depression of soul which I can compare to no earthly sensation more properly than to the after-dream of the reveller upon opium -- the bitter lapse into every-day life -- the hideous dropping off of the veil. There was an iciness, a sinking, a sickening of the heart -- an unredeemed dreariness of thought which no goading of the imagination could torture into aught of the sublime. What was it -- I paused to think -- what was it that so unnerved me in the contemplation of the House of Usher? It was a mystery all insoluble; nor could I grapple with the shadowy fancies that crowded conclusion, that while, beyond doubt, there are combinations of very simple natural objects which have the power of thus affecting us, still the analysis of this power lies among considerations beyond our depth. It was possible, I reflected, that a mere different arrangement of the particulars of the scene, of the details of the picture, would be sufficient to modify, or perhaps to annihilate its capacity for sorrowful impression; and, acting upon this idea, I reined my

# Page Surface Dewarping



- OCR error rates (commercial): 12.6% → 1.0%

# Outline
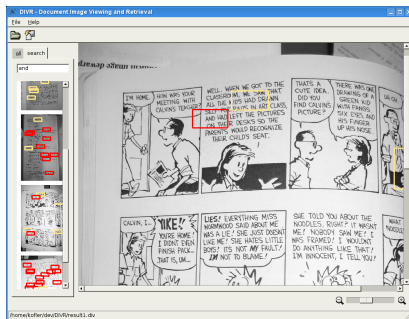
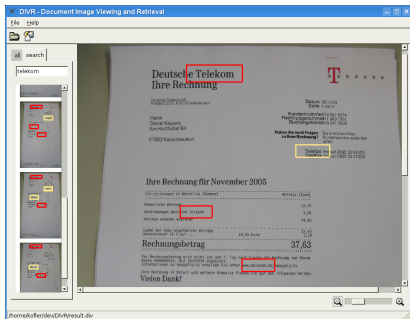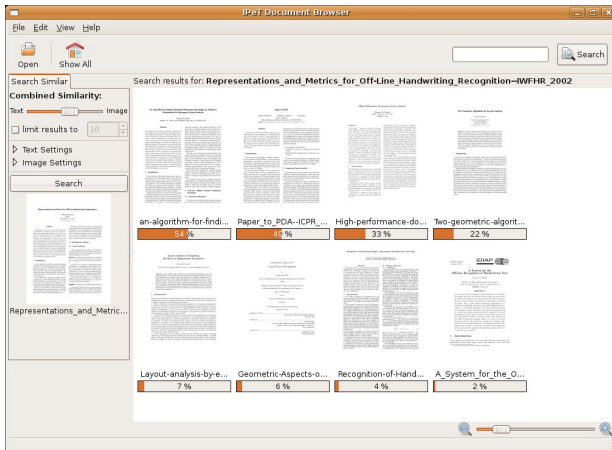# Document Image Viewing and Retrieval



- special OCR for captured documents

# Document Image Viewing and Retrieval



- ▶ special OCR for captured documents

# Document Browser



- stand-alone document browser
- server architecture for visual and textual search
- uses OCR server for scans and PDFs without included text

UNIVERSIDAD POLITÉCNICA DE VALENCIA

# Architectural Overview

# Outline

# Document Image to HTML Conversion

- Put together components of a complete document analysis system
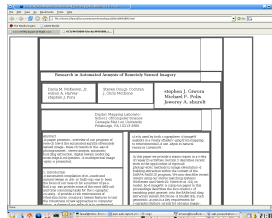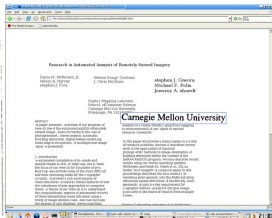- Display the results of OCR and layout analysis for debugging purpose

Figure 6 (a) A clonal, partially mantled ripe fruit with one supplementary separation at the fruit base (position 1) and site of eventual separation present (RHS). No abscission at the position of the supplementary carpe

considerable distances through the cell walls of adjoining tissues), it is evident that substrate specificity exists between the secreted enzymes and the walls of a limited numberofcells which are restricted to the immediate vicinity of the zone [6]. In this way, only certain cells become separated from their neighbours by the enzymes that are induced in the zone. In these dicotyledonous fruits. once the abscission cascade of enzymes is produced (diagnostically, this usually includes a specific9.5pI isozyme of p-1. 4-g1ucanhydrolase) Separation is initiated across all the

within 24 hr of the initial after declines. Ethylene onset of cell separation cate that the Se events a linked. Time-Course eXp shown that only those r riSe in ethylene synthesi tion at the fruit-pedicel j l) figure 7. It is not mesocarp tissue abuts directly onto zone eral layers of cells that s carotene nor storage lip barrier between the meS

classifier confidence information; black: high confidence, blue: medium confidence, red: low confidence; no language model used

UNIVERSIDAD
POLITECNICA
DE VALENCIA

# Outline

# Image-Based HTML Layout Verification

- ▶ rendering of a web page to image
- ▶ layout analysis of the captured image
- ▶ verify the rendered layout of the webpage

# Image-Based HTML Layout Verifier

## Problems in web page rendering

- ▶ Browser incompatibilites

UNIVERSIDAD
POLITÉCNICA
DE VALENCIA

# Image-Based HTML Layout Verifier

Problems in web page rendering

- ▶ Browser incompatibilites
- ▶ Large fonts for visually impaired

# Image-Based HTML Layout Verifier

Problems in web page rendering

- ▶ Browser incompatibilites
- ▶ Large fonts for visually impaired

Solution: Image-based layout verification

- ▶ Rendering of a web page to image
- ▶ Layout analysis and OCR of the captured image
- ▶ Check for
  - ▶ Usable page layouts
  - ▶ Readability of the rendered text
  - ▶ Visibility of all textual contents

# Rendering a Web Page to Image

representative browsers

- ▶ Internet Explorer (Windows)
- ▶ Mozilla Firefox (Linux)
- ▶ Safari (MacOS)

UNIVERSIDAD
POLITÉCNICA
DE VALENCIA

# Web Page Image Binarization

salient features of web page screenshots

- ▶ Extensive use of colors
- ▶ Both normal and inverted text in the same image
- ▶ No noise
- ▶ No skew
- ▶ Perfectly rendered fonts

$\rightarrow$ Seems to be a trivial problem, but

- ▶ Font anti-aliasing, colored non-uniform background, colored text, . . .

UNIVERSIDAD POLITECNICA DE VALENCIA

# Text-Line Extraction

- goal: highlighting of differences in layouts of differently rendered webpages



(a) Internet Explorer      (b) Firefox      (c) Safari

UNIVERSIDAD POLITÉCNICA DE VALENCIA

# Text-Line Matching

# Layout Verification

# Content Verification

- ▶ goal:
  - ▶ identify incorrectly rendered and missing text
- ▶ problem:
  - ▶ low OCR accuracy on screenshots of webpages
- ▶ approach:
  - ▶ OCR on HTML page
  - ▶ highlight incorrectly rendered text

# Image-Based HTML Layout Verifier

# Image-Based HTML Layout Verifier

# Outline

UNIVERSIDAD
POLITECNICA
DE VALENCIA

# Bibliographic Meta-Data Extraction

- ▶ task: extract structured meta-data from references
- ▶ problem: strong variations across different reference styles
  - ▶ subfield ordering
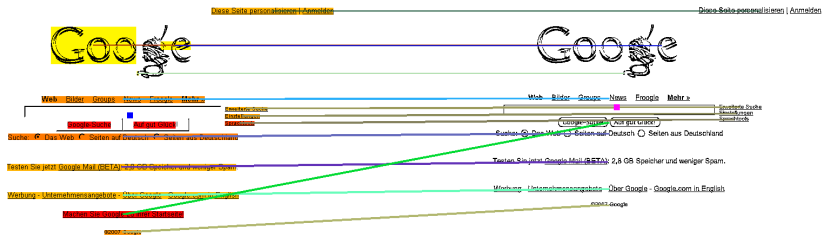  - ▶ partitioning symbols
  - ▶ spacing differences
  - ▶ content representation
- ▶ goal: retrieve labeling of reference according to semantics

- ▶ example:
  - ▶ input: (plain text)
    Davenport, T., D. DeLong and M. Beers, "Successful knowledge management projects," Sloan management review, 39, 2, (1998), 43–57.
  - ▶ output: (bibtex)
    author = "Davenport, T. and DeLong, D. and Beers, M."
    title = "Successful knowledge management projects"
    journal = "Sloan management review"
    volume = "39", number = "2", year = "1998", pages = "43–57"

# probabilistic finite state transducers

- ▶ motivation:
  - ▶ modular and flexible
  - ▶ composition of complex models via abstract operations
  - ▶ training-based derivation of weights
  - ▶ intuitive illustration as directed graph
- ▶ model:
  - ▶ one PFST for each occuring subfield



  - ▶ language model is built as a closure of subfield bigrams
  - ▶ training on dataset yields probabilistically weighted PFST

UNIVERSIDAD
POLITECNICA
DE VALENCIA

# performance evaluation of the system

- ► Cora dataset:
  - ► used for training and evaluation purposes
  - ► publicly available and most commonly applied
  - ► consists of 500 research paper citations
- ► Cora evaluation:

|  | word | field | instance |
|---|---|---|---|
| **CRF** (Peng et al.) | 95.4 |  | 77.3 |
| **PFST** | 88.5 | 82.6 | 42.7 |
| **HMM** (Seymore et al.) | 85.1 |  | 10.0 |
| **INFOMAP** (Day et al.) |  | 73.3 |  |

UNIVERSIDAD POLITECNICA DE VALENCIA

# Outline

UNIVERSIDAD POLITECNICA DE VALENCIA

# Arc and Line Detection - Motivation

- analysis of scanned technical drawings
- reconstruction of CAD data
- use advantages of CAD storage for archives of drawings

# Arc and Line Detection - Method

- ▶ use runs of black pixels
- ▶ explicitly include line thickness
- ▶ no preprocessing like thinning, line adjacency graphs, . . .
- ▶ global optimization, no heuristics
- ▶ use quality function, branch-and-bound, interval arithmetic
  - ▶ keep priority queue of parameter regions
  - ▶ use upper bound of quality estimate for region
  - ▶ on best region:
    - ▶ stop?
    - ▶ output?
    - ▶ split and re-insert

$$q(\vartheta, (x_0, x_1, y)) \ = \ \max\big(0, \ d^{-\frac{1}{2}} \sum_{x=x_0}^{x_1} \operatorname{sgn}(\frac{d}{2} - d_\vartheta(x, y))\big)$$

$$q(\vartheta, (x_0, x_1, y)) \ = \ \max\big(0, \ 1 - \frac{|\frac{d}{2} - d_\vartheta(x_0, y)|}{\sigma^2}\big) + \max\big(0, \ 1 - \frac{|\frac{d}{2} - d_\vartheta(x_1, y)|}{\sigma^2}\big)$$

# Arc and Line Detection - Examples

# Arc and Line Detection - Summary

- globally optimal detection possible
- very exact results
- results (VRI-scores) on GREC 2003 contest images:

    0.757 our method (2005)
    0.609 S. JiQiang (2003)
    0.487 D. Elliman (2003)

- 2nd place in GREC2005 contest
- current implementation memory intensive (500M),
  takes some time ($\sim$5min)

# Outline

UNIVERSIDAD POLITECNICA DE VALENCIA

# Historical Document Revision Detection



- robust document image matching for historical documents
- question:
  Which changes were made between different printings?
- uses geometric matching and Fourier contour descriptors

# RAST: Application in Inspection

# RAST for Object Matching



translation: (-14., 23.1)
rotation: 0.025
scale: 0.84

UNIVERSIDAD
POLITECNICA
DE VALENCIA

# Object-Based Image Retrieval



- detect whether an image contains an object of a given class
- efficient matching for fully-connected patch-based model
- uses our RAST algorithm
- optimal, statistically well-founded

# Object-Based Image Retrieval – Method

- ▶ match weakly annotated reference images, no segmentation
- ▶ patch-based approach (interest points, cluster descriptors)
- ▶ factor dependencies: $x/y$-translation, rotation, scale
- ▶ find optimal match using branch-and-bound approach
- ▶ set of reference patches $R$, test patches $S$

$$\hat{\vartheta}(R, S) := \arg \max_{\vartheta \in T} Q(\vartheta, R, S)$$

$$Q(\vartheta, R, S) := \sum_{p \in R} q(\vartheta, p, S)$$

$$q(\vartheta, p, S) := \begin{cases} 1 & \text{if } \exists p' \in S : l_p = l_{p'} \wedge d(\vartheta, p, p') \leq d_0 \\ 0 & \text{otherwise} \end{cases}$$

UNIVERSIDAD POLITECNICA DE VALENCIA

# Object-Based Image Retrieval – Results

| method | airp. | faces | mot. |
|--------|------:|------:|-----:|
| constellation model A | 32.0 | 6.0 | 16.0 |
| automatic segmentation | 2.2 | 0.1 | 10.4 |
| texture feature combination | 0.8 | 1.6 | 8.5 |
| constellation model B | 9.8 | 3.6 | 7.5 |
| PCA SIFT features | 2.1 | 0.3 | 5.0 |
| discriminative salient patches, SVM | 7.0 | 2.8 | 3.8 |
| spatial part-based model | 6.7 | 1.8 | 3.0 |
| constellation model C | 6.3 | 9.7 | 2.7 |
| patch histograms A | 3.8 | 7.1 | 2.5 |
| features inspired by visual cortex | 3.3 | 1.8 | 2.0 |
| patch histograms B | 1.4 | 3.7 | 1.1 |
| IPeT approach | 4.8 | 2.8 | 1.3 |

- ▶ error rates [%] on Caltech data
- ▶ possible use in DIA → logo recognition

UNIVERSIDAD POLITECNICA DE VALENCIA

# DFKI-IUPR Demos

http://demo.iupr.org/