

Introduction

Document image analysis (DIA) is the subfield of digital image processing that aims at converting document images to a symbolic form for modification, storage, retrieval, reuse, and transmission^[1]. However, more information is conveyed in addition to such a transcription, e.g., writer idiosyncrasies, data arrangement (tables, forms, etc.).

Metadata in Context^[2]:

- Descriptive Metadata: describes a resource for discovery and identification, e.g., author, and writer idiosyncrasies.
- Structural Metadata: indicates how document components are organized, e.g., tables, preprinted ruling lines, etc.

Thesis Hypothesis:

By exploiting various metadata in off-line handwritten documents, we are able to restore the original structure between documents, build new relationships from them, and facilitate problem solving in information retrieval tasks.



A high-level view of modules of

document analysis system for various applications.

Table Detection in HW Documents

Table Detection via Dynamic Programming ^[3] Each page is decomposed into: 1. one table 2. multiple tables Row correlation is based on inside-space stream.

Don Fernando, il projento de su segundo, matrimornio, importa analizarla detendamente por ser annito sobre el que Lievitores están muy discordes. Soman "que an se espresa inumora intre las causas que dividiron al Fry à contracto Exploiting Metada la probudar divisioner que existian en su familia, y la in constatibilidad de su caracter of sus miras con las de su yorno_

Ph.D Candidacy: Jin Chen Thesis Advisor: Daniel Lopresti Computer Science & Engineering, Lehigh University, Bethlehem, PA, United States

Table Detection Performance

- 62 Arabic handwritten documents are scanned at 600 dpi. The ground-truth for text regions (HW vs. MP) is available as bounding polygons.
- 20 pages for testing and the other 42 for training.
- Use area ratio-based measures proposed in Shafait and Smith ^[5], where bounding boxes are used to describe table regions.







Conclusions and Future Work

- Existing methods in MP documents do not completely solve the problem in HW documents.
- Part of the errors are caused by complicated layouts, such as letter forms, signatures, etc.
- Future work includes detection methods requiring no row segmentation, and better layout analysis.

Writer ID w/ Severe Data Constraints Scenarios and Methodologies

- Although usually assumed sufficient for research, data can become a problem for real-world writer ID, e.g., authors are unavailable, uncooperative.
- Model-Generated Handwriting and Model-Perturbed Handwriting can be used.
- We adopt Varga and Bunke's perturbation model ^[6] and employ user studies to calibrate ^[7].

متسامح حدير لا يتعارم موالذهداف العامة
متسامح جدير لايتارم م الذهدان العلمة
سنسادح جبرالا بيعارمام الذهراق العامة
ستسادح حديد لا يتعارمن موالذهداف العادة
سنسامح جبر لا يتعارف موالذهراف العلمة

ICDAR Doctoral Consortium 2011', Beijing, China



			I chen carronisme tracker de constrte ter a	the summer by	trucks is so at the s	A this surraining breaks & or	uño tau a c		
ngba dan agamat ke menan walitati shikada menan angang ingin di mita di sudi di sila dari pang ingi di kemat dangan di sila okada sepak di kemat pang ingi di kemat da kalila okada sebaga tangan mana at at a		يولك ركيب محمد معرفين المسيسين بيليك محمد محمد محمد محمد محمد محمد محمد محمد	More opened to <u>instance of the perform Managers</u> and the second	anner familier in star	gener y utbila folden til atten fra den for som er som atten er som er som att den er som er so att den er som er som att den er som er att den er som	mfall, show appear to summarize the second s	endondi eindoze wei kon maritur fangi ha		
the problem is a function of the problem in the problem in the problem is the pro		<u>An en anno an anna an a</u>	Sheriya i alaya sa mala mbanan - Kariya mar tani da ta alaya da anka na jina kata ka alaya na kariya da anka na jina kata ka alaya na kariya da anka na ta alaya da alaya da alaya bila na manda angana da angana da kariya da alaya da alaya	hi se Penin y Andrew Li Charlos transmi de de grafie andre any de grafie inst transmit de grafie de grafie de grafie de any de grafie de grafie de grafie de grafie de any de grafie de grafie de grafie de grafie de any de grafie de grafie de grafie de grafie de grafie de any de grafie de	e s. regular untrivente entre teure reviewe state la imposfulle, that was the la la imposfulle, that was the la la imposfulle, that was the laterature of the state of the	hi), na Sherina at anhar an anna an Mariana taibida tao an Anna an an a isin an Anna Anna an In gai chaisaga at agus la ionnachdia Ray, bita na annacha agus an Anna Anna Ray, bita na	in a second seco		
The strain of the second secon		The second secon	and the speech of the set of the	I de sense la constant de separa partir et terret à constant de separa partir et terret à constant de separa, a que de sense de sense sense sense sense sense sense que de sense sense sense sense sense sense sense sense sense que de sense sen	ta Pine California a Classical program a Classica	2 Di norme da ma grandera de sela Petade strincte, de la Cara da marte instruction de totas de la con- nectoria de la consecuencia da la con- ción de la consecuencia de marte a del consecuencia de marte			
$\mathcal{X}_{\mathrm{reg}}^{\mathrm{chrom}}(\mathcal{X}) = \mathcal{X}_{\mathrm{reg}}^{\mathrm{chrom}}(\mathcal{X}) $		<u>aller en finite la nata</u> de cas Sen Senale à pr <u>ainduite</u> Senal sight antic <u>e bindenata</u> ingla la indian de nacharde de la cast parase de la case qualitée quara	an a	n de la cata al constante de la constante de	artin spir a britheoppi a' de su sapada - ardiniceria , de gar in andra spir d'ya - Erennan Van a a napan Interna d'Arte à antenda	Biographic Andre Ian anter age in h Des Strande ei preiste de mangele ei ingele sodionie dit her hande einer sodie her inder als and her inder Strand De ingende als her inder anter anderlien di De	telege d deserves , son a depe		
ha yabadar Siriirea ya mgablikidi di su swater .	orstion in in faille a faire. 4 ins wins: to he le is gove.	hereafields thinks by testion to be said and that here and gives a sayah	alar kanisan operantar av se forsker yte inner som	divisione ana contra de se northe grane segnificade de se roater y	atin wa se fashe y ke in. so was sa ke k s gana.	<u>la salada da dinina</u> pa witan na sa fa sayat <u>itid da namata</u> y sa wine sa bi	hu y he inn 		
	-4.16	0.55	0.9			-0.01	0		
σ	5.60	2.18	5.11	0	0.01	0.05	1.41		
Germana Dataset									
μ	-12.30	3.80	12.85	0.25	-0.49	-0.09	0.25		
σ	39.01	20.70	34.03	1.25	2.24	0.26	0.22		
Germana Dataset (an existing method [4])									
μ	-49.40	57.60	70.40	-1.20	-2.40	0.02	0		
σ	62.64	113.30	62.12	5.68	8.82	1.09	0		