

Document Recognition and Retrieval XXII

Eric K. Ringger
Bart Lamiroy
Editors

11–12 February 2015
San Francisco, California, United States

Sponsored by
IS&T—The Society for Imaging Science and Technology
SPIE

Published by
SPIE

Volume 9402

The papers included in this volume were part of the technical conference cited on the cover and title page. Papers were selected and subject to review by the editors and conference program committee. Some conference presentations may not be available for publication. The papers published in these proceedings reflect the work and thoughts of the authors and are published herein as submitted. The publishers are not responsible for the validity of the information or for any outcomes resulting from reliance thereon.

Please use the following format to cite material from this book:

Author(s), "Title of Paper," in *Document Recognition and Retrieval XXII*, edited by Eric K. Ringger, Bart Lamiroy, Proceedings of SPIE-IS&T Electronic Imaging, SPIE-IS&T Vol. 9402, Article CID Number (2015)

ISSN: 0277-786X

ISBN: 9781628414929

Copublished by

SPIE

P.O. Box 10, Bellingham, Washington 98227-0010 USA

Telephone +1 360 676 3290 (Pacific Time) · Fax +1 360 647 1445

SPIE.org

and

IS&T—The Society for Imaging Science and Technology

7003 Kilworth Lane, Springfield, Virginia, 22151 USA

Telephone +1 703 642 9090 (Eastern Time) · Fax +1 703 642 9094

imaging.org

Copyright © 2015, Society of Photo-Optical Instrumentation Engineers and The Society for Imaging Science and Technology.

Copying of material in this book for internal or personal use, or for the internal or personal use of specific clients, beyond the fair use provisions granted by the U.S. Copyright Law is authorized by the publishers subject to payment of copying fees. The Transactional Reporting Service base fee for this volume is \$18.00 per article (or portion thereof), which should be paid directly to the Copyright Clearance Center (CCC), 222 Rosewood Drive, Danvers, MA 01923. Payment may also be made electronically through CCC Online at copyright.com. Other copying for republication, resale, advertising or promotion, or any form of systematic or multiple reproduction of any material in this book is prohibited except with permission in writing from the publisher. The CCC fee code is 0277-786X/15/\$18.00.

Printed in the United States of America.

Paper Numbering: Proceedings of SPIE follow an e-First publication model, with papers published first online and then in print. Papers are published as they are submitted and meet publication criteria. A unique citation identifier (CID) number is assigned to each article at the time of the first publication. Utilization of CIDs allows articles to be fully citable as soon as they are published online, and connects the same identifier to all online, print, and electronic versions of the publication. SPIE uses a six-digit CID article numbering system in which:

- The first four digits correspond to the SPIE volume number.
- The last two digits indicate publication order within the volume using a Base 36 numbering system employing both numerals and letters. These two-number sets start with 00, 01, 02, 03, 04, 05, 06, 07, 08, 09, 0A, 0B ... 0Z, followed by 10-1Z, 20-2Z, etc.

The CID Number appears on each page of the manuscript. The complete citation is used on the first page, and an abbreviated version on subsequent pages.

Contents

- v *Authors*
- vii *Conference Committee*
- ix *Introduction*

SESSION 1 DOCUMENT LAYOUT ANALYSIS AND UNDERSTANDING

- 9402 04 **Ground truth model, tool, and dataset for layout analysis of historical documents** [9402-2]
- 9402 05 **Use of SLIC superpixels for ancient document image enhancement and segmentation** [9402-3]
- 9402 06 **Software workflow for the automatic tagging of medieval manuscript images (SWATI)** [9402-4]

SESSION 2 DOCUMENT STRUCTURE SEMANTICS, FORMS, AND TABLES

- 9402 07 **Math expression retrieval using an inverted index over symbol pairs** [9402-5]
- 9402 08 **Min-cut segmentation of cursive handwriting in tabular documents** [9402-6]
- 9402 09 **Cross-reference identification within a PDF document** [9402-7]
- 9402 0A **Intelligent indexing: a semi-automated, trainable system for field labeling** [9402-8]

SESSION 3 TEXT ANALYSIS

- 9402 0B **Re-typograph phase I: a proof-of-concept for typeface parameter extraction from historical documents** [9402-9]
- 9402 0C **Clustering of Farsi sub-word images for whole-book recognition** [9402-10]
- 9402 0D **Gaussian process style transfer mapping for historical Chinese character recognition** [9402-11]
- 9402 0E **Boost OCR accuracy using iVector based system combination approach** [9402-12]

SESSION 4 HANDWRITING I

- 9402 0F **Exploring multiple feature combination strategies with a recurrent neural network architecture for off-line handwriting recognition** [9402-13]
- 9402 0G **Spotting handwritten words and REGEX using a two stage BLSTM-HMM architecture** [9402-14]
- 9402 0H **A comparison of 1D and 2D LSTM architectures for the recognition of handwritten Arabic** [9402-15]
- 9402 0I **Aligning transcript of historical documents using dynamic programming** [9402-16]
- 9402 0J **Offline handwritten word recognition using MQDF-HMMs** [9402-17]

SESSION 5 QUALITY AND COMPRESSION

- 9402 0K **Separation of text and background regions for high performance document image compression** [9402-19]
- 9402 0L **Metric-based no-reference quality assessment of heterogeneous document images** [9402-20]

SESSION 6 GRAPHICS AND STRUCTURE

- 9402 0M **Clustering header categories extracted from web tables** [9402-21]
- 9402 0N **A diagram retrieval method with multi-label learning** [9402-22]
- 9402 0O **Detection of electrical circuit elements from documents images** [9402-24]

SESSION 7 HANDWRITING II

- 9402 0P **Missing value imputation: with application to handwriting data** [9402-25]

Authors

Numbers in the index correspond to the last two digits of the six-digit citation identifier (CID) article numbering system used in Proceedings of SPIE. The first four digits reflect the volume number. Base 36 numbering is employed for the last two digits and indicates the order of articles within the volume. Numbers start with 00, 01, 02, 03, 04, 05, 06, 07, 08, 09, 0A, 0B...0Z, followed by 10-1Z, 20-2Z, etc.

Barrett, William A., 08, 0A
Bideault, Gautier, 0F, 0G
Blégean, Julien, 0B
Bouville, Thomas, 0B
Breuel, Thomas M., 0H
Brunessaux, S., 0F
Busch, Hannah, 06
Cao, Hongliu, 0B
Cao, Huaigu, 0E
Chanda, Bhabatosh, 0O
Chandna, Swati, 06
Chatelain, Clément, 0F, 0G
Chen, Kai, 04
Clawson, Robert, 0A
Cohen, Rafi, 0I
Das, Amit, 0O
Davis, Brian L., 08
De, Paramita, 0O
El-Sana, Jihad, 0I
Embley, David W., 0M
Essoukri Ben Amara, Najoua, 05
Fan, Wei, 0K
Feng, Jixiong, 0D
Fu, Songping, 0N
Gao, Liangcai, 09
Ghamizi, Salah, 0B
Gomez-Krämer, Petra, 05
Hambarde, Mangesh, 0J
Hennebert, Jean, 04
Héroux, Pierre, 05
Houpin, Romain, 0B
Ingold, Rolf, 04
Jejkal, Thomas, 06
Kabir, Ehsanollah, 0C
Kedem, Klara, 0I
Kochar, Shaivi, 0J
Krause, Celia, 06
Krishnamoorthy, Mukkai, 0M
Lamiroy, Bart, 0B
Lebourgeois, Franck, 0D
Li, Sida, 09
Liu, Lu, 0N
Liwicki, Marcus, 04
Lloyd, Matthias, 0B
Lu, Xiaoqing, 0N
Mandal, Sekhar, 0O
Mehri, Maroua, 05
Mioulet, Luc, 0F, 0G
Mullot, Rémy, 05
Nagy, George, 0M
Naai, Satoshi, 0K
Natarajan, Prem, 0E
Nayef, Nibal, 0L
Ogier, Jean-Marc, 0L
Paquet, Thierry, 0F, 0G
Patil, Ajay, 0J
Peng, Liangrui, 0D
Peng, Xujun, 0E
Prabhune, Ajinkya, 06
Qu, Jingwei, 0N
Rabaev, Irina, 0I
Ramachandrula, Sitaram, 0J
Sahoo, Dushyant, 0J
Seth, Sharad, 0M
Seuret, Mathias, 04
Sliiti, Nabil, 05
Soheili, Mohammad Reza, 0C, 0H
Srihari, Sargur N., 0P
Stalnaker, David, 07
Stotzka, Rainer, 06
Stricker, Didier, 0C, 0H
Sun, Jun, 0K
Swingle, Scott D., 08
Tang, Zhi, 09, 0N
Tonne, Danah, 06
Vanscheidt, Philipp, 06
Wei, Hao, 04
Xu, Zhen, 0P
Yousefi, Mohammad Reza, 0H
Yu, Yinyan, 09
Zanibbi, Richard, 07

Conference Committee

Symposium Chair

Sheila S. Hemami, Northeastern University (United States)

Symposium Co-chair

Choon-Woo Kim, Inha University (Korea, Republic of)

Conference Chairs

Eric K. Ringger, Brigham Young University (United States)

Bart Lamiroy, Université de Lorraine (France)

Conference Program Committee

Gady Agam, Illinois Institute of Technology (United States)

Sameer K. Antani, National Library of Medicine (United States)

Elisa H. Barney Smith, Boise State University (United States)

William A. Barrett, Brigham Young University (United States)

Kathrin Berkner, Ricoh Innovations, Inc. (United States)

Bertrand Coüasnon, Institut National des Sciences Appliquées de
Rennes (France)

Hervé Déjean, Xerox Research Centre Europe Grenoble (France)

Xiaoqing Ding, Tsinghua University (China)

Jianying Hu, IBM Thomas J. Watson Research Center (United States)

Ergina Kavallieratou, University of the Aegean (Greece)

Christopher Kermorvant, A2iA SA (France)

Laurence Likforman-Sulem, Télécom ParisTech (France)

Xiaofan Lin, A9.com, Inc. (United States)

Marcus Liwicki, Deutsches Forschungszentrum für Künstliche
Intelligenz GmbH (Germany)

Daniel P. Lopresti, Lehigh University (United States)

Umapada Pal, Indian Statistical Institute (India)

Sargur N. Srihari, University at Buffalo (United States)

Venkata Subramaniam, IBM Research - India (India)

Kazem Taghva, University of Nevada, Las Vegas (United States)

George R. Thoma, National Library of Medicine (United States)

Christian Viard-Gaudin, Université de Nantes (France)

Pingping Xiu, Microsoft Corporation (United States)

Berrin Yanikoglu, Sabanci University (Turkey)

Richard Zanibbi, Rochester Institute of Technology (United States)

Jie Zou, National Library of Medicine (United States)

Conference Review Committee

Jin Chen, Lehigh University (United States)

Session Chairs

Keynote Session I

Bart Lamiroy, Université de Lorraine (France)

Eric K. Ringger, Brigham Young University (United States)

Keynote Session II

Bart Lamiroy, Université de Lorraine (France)

Eric K. Ringger, Brigham Young University (United States)

- 1 Document Layout Analysis and Understanding
Richard Zanibbi, Rochester Institute of Technology (United States)
- 2 Document Structure Semantics, Forms, and Tables
Xiaofan Lin, A9.com, Inc. (United States)
- 3 Text Analysis
Eric K. Ringger, Brigham Young University (United States)
- 4 Handwriting I
Daniel P. Lopresti, Lehigh University (United States)
- 5 Quality and Compression
William A. Barrett, Brigham Young University (United States)
- 6 Graphics and Structure
Bart Lamiroy, Université de Lorraine (France)
- 7 Handwriting II
Daniel P. Lopresti, Lehigh University (United States)

Introduction

On behalf of the program committee, we welcome you to the twenty-second Document Recognition and Retrieval conference (DRR 2015) in San Francisco, California, USA. DRR is held annually as part of the IS&T/SPIE Symposium on Electronic Imaging. It is one of the leading international conferences on document recognition, including related research on information retrieval and text mining, especially as they pertain to document images.

This year we received 44 paper submissions. Of those, 23 papers were accepted, for an overall acceptance rate of 52%. In this year's edition of the conference, all papers will be presented as oral presentations. We want to sincerely thank the program committee members and additional reviewers for helping us create a strong technical program. The additional reviewers: William Lund (BYU, USA), Douglas Kennard (BYU, USA), Oliver Nina (University of Central Florida, USA), Cérés Carton (IRISA, France), Aurélie Lemaitre (IRISA, France), Yann Ricquebourg (IRISA, Rennes, France), Santosh K.C. (NIH-NLM, USA), Jin Chen (Lehigh University, USA).

The program encompasses the following sessions: Document Layout Analysis and Understanding; Document Structure Semantics, Forms, and Tables; Text Analysis; Handwriting I & II; Quality and Compression; and Graphics and Structure. Topics that are current elsewhere in the pattern recognition and machine learning communities, including probabilistic graphical models and deep neural networks, also make important appearances in the technical program.

This year's invited Keynotes will be given by Brewster Kahle, founder of the Internet Archive, and by Dan Klein, professor of computer science at the University of California, Berkeley. Brewster will address the challenges of large document repositories such as the Internet Archive, while Dan's talk will address unsupervised transcription for documents and music. Our invited speakers are helping to build bridges between DRR and neighboring disciplines such as computer vision and NLP.

Once again we recognize the strong contributions of student primary authors. For the best student paper award, we are grateful to Elisa H. Barney Smith (chair) and the award committee for carrying out the difficult task of choosing winning papers. The winners will be announced at the award ceremony of the conference. This year, thanks to the generous sponsorship by A9 and Google, we are able to complement these awards with a generous gift. We are very grateful for this continued support by our sponsors.

This year promises to hold change for DRR. After many years, SPIE and IS&T have announced that they are parting ways in their sponsorship of the annual Electronic

Imaging conference, DRR's umbrella venue. Therefore DRR will feature a panel discussion in which we address the options for future editions.

We hope that attendees have an excellent experience at DRR XXII and that readers of the proceedings find work that helps them to advance the state of the art in document recognition and retrieval!

Eric K. Ringger
Bart Lamiroy