Next

# Digital Document Processing

Bidyut B. Chaudhuri (Ed.)

Major Directions and
Recent Advances

Springer

# Digital Document Processing

Bidyut B. Chaudhuri
(Ed.)

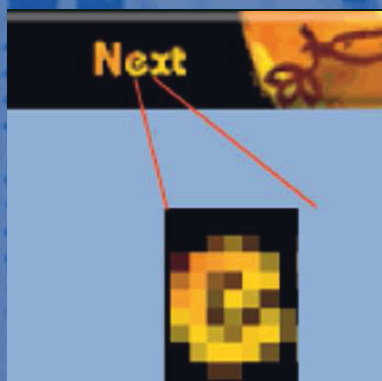## Major Directions and Recent Advances

Springer

Advances in Pattern Recognition

**Advance in Pattern Recognition** is a series of books which brings together current developments in all areas of this multi-disciplinary topic. It covers both theoretical and applied aspects of pattern recognition, and provides texts for students and senior researchers.

Springer also publishes a related journal, **Pattern Analysis and Applications**. For more details see: http://link.springer.de

The book series and journal are both edited by Professor Sameer Singh of Exeter University, UK.

*Also in this series:*

Bidyut B. Chaudhuri (Ed.)

# Digital Document Processing

**Major Directions and Recent Advances**

Springer

Bidyut B. Chaudhuri, PhD
Indian Statistical Institute, Kolkata, India

*Series editor*
Professor Sameer Singh, PhD
Department of Computer Science, University of Exeter, Exeter, EX4 4PT, UK

9 8 7 6 5 4 3 2 1

springer.com

# Preface

The field of automatic document processing, more than a century old, is quite a mature one. It started with an attempt to automatically read printed text in the era well before the birth of digital computers and has been since continuing on various topics like document image enhancement, document structure and layout analysis, handwritten character recognition, document data compression, graphics recognition, document information retrieval and meta-data extraction. In addition to OCR, applications like tabular form and bank cheque processing or postal mail reading are of great interest to the software industry.

This edited book is a compilation of twenty chapters written by leading experts on the above and several other important topics. These chapters describe the state of the art in these areas, enumerate the research challenges and propose one or more possible solutions in a systematic way. Usually, edited books are compiled on some special area of a general discipline. But this one attempts to cover wider aspects of digital document processing and hence has the flavour of a handbook. Since there is no standard textbook with a wide coverage of the subject, this book will immensely help students taking undergraduate and graduate courses in digital document processing. Also, it is hoped that the researchers in the field will benefit from this book. About 9 years ago, Bunke and Wang edited a useful book of similar nature. However, the activities on document analysis have since advanced much, with newer techniques being invented by the community and younger disciplines like super-resolution text processing, handwriting individuality identification or web document mining gaining importance. Even on older topics like OCR, good work is in progress on challenging problems like the reading of Indian and Tibetan scripts. All these advancements contributed to the need for another edited book.

This book starts with an excellent introduction to the general discipline of document processing, followed by studies on document structure analysis. The next few chapters describe OCR systems for some difficult printed scripts that are followed by advances in on-line and off-line handwritten

text recognition, with applications to postal automation and bank cheque processing. Then come the special topics of mathematical expression and graphics recognition, as well as super-resolution text analysis. Other problems like image degradation modelling, meta-data extraction, document information retrieval are addressed in the next few chapters. Some emphasis has been given on web document analysis and data mining problems. The last three chapters of this book are dedicated to these topics.

This book is the outcome of my interaction with the authors over a reasonably long period. In this endeavour, my sincerest thanks go to Professor Horst Bunke with whom I initially made plans to co-edit this book. At the preparatory stage, I received generous help from his vast experience on various issues like framing of book structure, the choice of chapter topics, the choice of potential authors as well as fixing other subtle points. However, because of heavy workload and subsequent illness, he had to discontinue his involvement in the editorial process. Nevertheless, I am grateful for his continuous advice, encouragement and best wishes for the completion of this project.

I express my sincere thanks to all authors for their hard work in preparing their manuscripts. Special mention should be made of Prof. L. Schomaker who not only wrote an introductory chapter but also patiently read all other chapters in order to write critical summaries of them. The readers will get a quick idea of all the topics by reading this chapter alone.

During my work on this book I enjoyed a sabbatical leave as well as a *Jawaharlal Nehru Fellowship*. The support of my institute and of the Jawaharlal Nehru Memorial Fund is gratefully acknowledged. I am thankful to my colleagues, Dr Utpal Garain, Mr Chittaranjan Das and Ms Shamita Ghosh, who helped me in various stages of the editorial work. The understanding and support of my family is also highly appreciated. Finally, I thankfully acknowledge the patience of Ms Catherine Brett of Springer Verlag for replying to my numerous e-mails and other support provided to me during the editing of this book.

Bidyut B. Chaudhuri
Kolkata, India

# Contents