The automatic recognition of text in images and its transformation into machine-readable formats is a long-standing promise of research in computer science to the text-based humanities and natural sciences but also to libraries. But while users can expect virtually error-free results for modern sources, the situation is different for more challenging materials: For example, newspapers from the first half of the 20th century already show many difficulties in recognition. The prospects for documents from the 15th-18th century are even worse.

Libraries expect no more and no less than reliable and efficient methods for the text digitization of their own collections. This is associated with the hope of greater outreach, especially via the Internet. With the help of virtual collections, libraries may reach substantially larger user groups which would hardly find their way into the local reading rooms. Unfortunately, text digitization is currently all too often implemented as a by-product of image digitization. A thorough analysis of the results is only carried out to a very limited extent.

On the part of the text-based humanities and natural sciences, there is a desire for large amounts of text data, which do not have to be collected by means of time-consuming, detailed transcription work but can be conveniently retrieved, searched and analyzed. Scientific investigations could thus concentrate on the actual research instead of spending a great deal of time and effort on data collection and processing. However, this requires easy-accessible research data of high quality.

For computer sciences, automatic text recognition is primarily a scientific problem. The aim is to develop algorithmic solutions the quality of which can be verified in a comparable manner using standardized, more or less representative data sets. The results which are presented, practically domain-independent error rates of below 1 %, always read promisingly, but often have only a very limited validity for their application in mass digitization scenarios and consequently for the creation of a reliable basis for text-based research. It took a critical examination of the results of industrial projects such as Google Books or research projects such as IMPACT (Improving Access to Text) to bring the task of text recognition back into the focus of scientific efforts.

However, there are legitimate reasons to hope that error rates which can be observed for modern sources may also be in range for historical documents: The use of statistical learning methods based on (deep) neural networks has led to an enormous leap in quality also in the field of text and layout recognition. This is reflected in a significant reduction of error rates due to a substantially higher tolerance to variances in the material which is to be processed. In addition, the task of handwriting recognition and the recognition of printings are regarded as instances of the same scientific problem.

A prerequisite for the successful use of machine learning methods, however, are great amounts of texts and structural annotations which are sufficiently accurate and representative of the subject (and therefore called ground truth). They serve as a reference point for training and evaluation. The creation of such materials is therefore at least as important as the development and adaptation of algorithms which make use of ground truth data.

This volume of the *Journal for Language Technology and Computational Linguistics* aims to meet both basic requirements for successful automatic text recognition. It contains methodological contributions concerning the actual recognition process as well as contributions dealing with the creation and optimization of necessary training and evaluation data. It is completed by the presentation of concrete resources and tools intended to contribute to the quality of automatic text recognition results with the perspective to qualify them as genuine research data in the not too distant future.

We would like to thank all authors for their contributions. Next, we would also like to thank the reviewers, who have contributed to the quality of the published articles in an excellent way. Last but not least we want to express our gratitude to the German Society for Computational Linguistics and Language Technology and especially Adrien Barbaresi and Lothar Lemnitzer for providing a forum for the topic of automatic text recognition.

The guest editors, Kay-Michael Würzner, Alexander Geyken, Günter Mülberger.