# The Information Retrieval Series    Volume 27

Donald Metzler

# A Feature-Centric
# View of Information
# Retrieval

Springer

Donald Metzler
Natural Language Group
Information Sciences Institute
University of Southern California
4676 Admiralty Way, Suite 1001
Marina del Rey, CA 90292
USA
metzler@isi.edu

*Cover design*: VTeX UAB, Lithuania

Printed on acid-free paper

*To* 晓黎

# Acknowledgements

# Contents