# Contents

## PART III  BIG DATA SECURITY AND PRIVACY

# PART IV  BIG DATA APPLICATIONS

List o

**T. Achalake**
King Mong

**P. Ameri**
Karlsruhe I

**A. Berry**
Deontik, Br

**N. Bojja**
Machine Ze

**R. Buyya**
The Univer
Australia

**W. Chen**
University c

**C. Deeroseja**
King Mongk

**A. Diaz-Pere**
Cinvestav-Ta

**H. Ding**
Xi'an Jiaoto

**X. Dong**
Huazhong U

**H. Duan**
The Univers

**S. Dutta**
Max Planck

**A. Garcia-Ro**
Cinvestav-Ta

**V. Gramoli**
University of

**X. Gu**
Huazhong U

**J. Han**
Xi'an Jiaotong

**B. He**
Nanyang Tec