# Handbook of Statistical Systems Biology

Edited by

**MICHAEL P. H. STUMPF**
*Imperial College London, UK*

**DAVID J. BALDING**
*Institute of Genetics, University College London, UK*

**MARK GIROLAMI**
*Department of Statistical Science, University College London, UK*

# Contents