Editors

# Gil Alterovitz and Marco Ramoni

# Knowledge-Based
# Bioinformatics

## From analysis to interpretation

# Contents