# Table of Contents

# 3

# Next-Generation Sequencing    49

# 4

## Advanced NGS Data Processing                                                93

# 5

## Working with Genomes                                                       121

# 8

# Using the Protein Data Bank                                                217

# 9

# Bioinformatics Pipelines                                                   249

# 10

## Machine Learning for Bioinformatics          273

# 11

## Parallel Processing with Dask and Zarr          291

# 12

## Functional Programming for Bioinformatics    313