# *Contents*

## Section I   Introduction to Cellular and Molecular Biology

**Section II    Introduction to Next-Generation
                Sequencing (NGS) and NGS Data Analysis**

# Section III   Application-Specific NGS Data Analysis

**Section IV    The Changing Landscape
of Next-Generation Sequencing
Technologies and Data Analysis**