



# 新一代测序数据分析

NEXT-GENERATION SEQUENCING DATA ANALYSIS

[美] 王忻琨 著  
陈浩峰 主译



科学出版社



新生物学丛书

# 新一代测序数据分析

Next-Generation Sequencing Data Analysis

〔美〕王忻琨 著

陈浩峰 主译

科学出版社

北京

## 内 容 简 介

本书是一部介绍有关新一代测序（NGS）数据分析方法的著作。书中全面系统地介绍了新一代测序技术的生物学意义、测序原理、分析过程和应用领域等；详细介绍了新一代测序数据的分析方法，包括其在基因组从头测序和重测序、转录组测序、小RNA测序、ChIP测序、表观基因组测序及宏基因组测序等应用中的具体分析方法，对读者学习新一代测序技术、促进该技术在生命科学各个领域中的应用有着重要意义。

本书可以作为具有生物学、遗传学、基因组学与生物信息学、农学等生命科学背景的高等学校师生和科研院所研究人员学习新一代测序技术原理及进行测序数据分析的参考书，也可供各大医院和医学院对于精准医学感兴趣的专业人士和其他生命科学从业者参考。

Next-Generation Sequencing Data Analysis, by Xinkun Wang.

Copyright © 2017 by Taylor & Francis Group LLC.

All Rights Reserved. Authorized translation from English language

edition published by CRC Press, an imprint of Taylor & Francis Group LLC.

本书封面贴有 Taylor & Francis 集团防伪标签，未贴防伪标签属未获授权的非法行为

### 图书在版编目( CIP )数据

新一代测序数据分析 / (美) 王忻琨著；陈浩峰主译. —北京：科学出版社，2018.2

(新生物学丛书)

书名原文：Next-Generation Sequencing Data Analysis

ISBN 978-7-03-056469-6

I. ①新… II. ①王… ②陈… III. ①生物信息论 IV. ① Q811.4

中国版本图书馆 CIP 数据核字 (2018) 第 019094 号

责任编辑：罗 静 / 责任校对：彭 涛

责任印制：张 伟 / 封面设计：刘新颖

科 学 出 版 社 出 版

北京东黄城根北街 16 号

邮政编码：100717

<http://www.sciencecp.com>

北京京华彩印刷有限公司 印刷

科学出版社发行 各地新华书店经销

\*

2018 年 2 月第 一 版 开本：720 × 1000 1/16

2018 年 2 月第一次印刷 印张：13 1/4

字数：262 000

定价：98.00 元

(如有印装质量问题，我社负责调换)

## 译者名单

---

主 译 陈浩峰

翻译人员 陈浩峰 王 静 韩 瑶 齐 洛

译者单位 中国科学院遗传与发育生物学研究所

# 译者前言

---

随着现代科学技术的发展，生命科学研究已经进入了组学时代。新一代测序（NGS）技术的出现及其广泛应用，极大地推动了生命科学各个领域的发展，促生了包括“精准医学”在内的一大批新兴学科。人们运用 NGS 技术获得了海量的测序数据，如何从这些数据中分析总结出新知识、新发现，并将其应用于生命科学科研和生物产业及医疗产业，是摆在广大 NGS 应用者面前的一大难题。大多数 NGS 用户都仅具有生命科学或医学背景，对 NGS 数据分析过程了解不足，急需学习这方面的知识。目前，读者获取的 NGS 数据分析的相关知识，往往来自于散在的研究或综述论文，或者来自于互联网，不利于构成完整的知识体系。鉴于此，我们组织翻译了由美国西北大学王忻琨博士所著的《新一代测序数据分析》一书，以期填补国内中文书籍这方面的空白。

本书是一部全面介绍新一代测序及基因组学和生物信息学的著作，全面兼顾了生物学基本原理、测序技术、测序数据分析处理等方面的知识，为读者全方位详细介绍了新一代测序数据分析所涉及的方方面面的内容，可以作为有生物学、遗传学背景的教师及学生学习新一代测序技术原理和测序数据分析的参考。我们将其介绍给国内的读者，希望能够对促进新一代测序技术在我国生命科学研究中的应用，乃至精准医学的发展有所裨益。

本书译者曾编写过主要阐述 NGS 实验方法的《新一代基因组测序技术》一书，希望本书的出版可以与之相得益彰，共同构成有关新一代测序的完整知识体系。

感谢中国科学院遗传与发育生物学研究所基因组生物学研究中心和植物基因组学国家重点实验室对本书翻译工作的支持，感谢 Illumina 中国公司、北京英木和生物技术有限公司的赞助。

限于译者水平，译文的不妥之处在所难免，望广大读者予以指正。

陈浩峰

2017 年 8 月 1 日

# 目 录

---

## 第一部分 细胞与分子生物学概论

1	细胞系统与生命密码 .....	3
1.1	细胞面临的挑战 .....	3
1.2	细胞如何面对挑战 .....	3
1.3	细胞内的各种分子 .....	4
1.4	细胞内结构与空间 .....	4
1.4.1	细胞核 .....	5
1.4.2	细胞膜 .....	6
1.4.3	细胞质 .....	6
1.4.4	内体、溶酶体和过氧化物酶体 .....	7
1.4.5	核糖体 .....	7
1.4.6	内质网 .....	8
1.4.7	高尔基体 .....	8
1.4.8	细胞骨架 .....	8
1.4.9	线粒体 .....	9
1.4.10	叶绿体 .....	10
1.5	细胞是一个系统 .....	11
1.5.1	细胞系统 .....	11
1.5.2	细胞的系统生物学 .....	11
1.5.3	如何研究细胞系统 .....	12
2	DNA 序列：基因组基础 .....	13
2.1	DNA 双螺旋和碱基序列 .....	13
2.2	DNA 分子如何复制和保持稳定性 .....	13
2.3	DNA 中保存的遗传信息如何转化为蛋白质 .....	15
2.4	基因组概览 .....	16
2.4.1	最小基因组 .....	16
2.4.2	基因组大小 .....	17

2.4.3 基因组中的蛋白质编码区 .....	17
2.4.4 基因组非编码区 .....	18
2.5 DNA 包装、序列访问和 DNA- 蛋白质互作 .....	20
2.5.1 DNA 包装 .....	20
2.5.2 序列访问 .....	20
2.5.3 DNA- 蛋白质互作 .....	20
2.6 DNA 序列的突变与多样性 .....	21
2.7 基因组演化 .....	23
2.8 表观基因组与 DNA 甲基化 .....	24
2.9 基因组测序与疾病风险 .....	25
2.9.1 孟德尔（单基因）疾病 .....	25
2.9.2 多基因控制的复杂疾病 .....	25
2.9.3 基因组不稳定导致的疾病 .....	26
2.9.4 表观基因组 / 表观遗传疾病 .....	26
3 RNA：转录后的序列 .....	27
3.1 RNA 作为信使 .....	27
3.2 RNA 的分子结构 .....	27
3.3 mRNA 的产生、加工与周转 .....	28
3.3.1 DNA 模板 .....	28
3.3.2 原核生物基因的转录 .....	28
3.3.3 真核生物基因 pre-mRNA 的初始转录 .....	30
3.3.4 从 mRNA 前体到成熟的 mRNA .....	31
3.3.5 运输与定位 .....	33
3.3.6 稳定性与降解 .....	33
3.3.7 mRNA 转录水平上调控的主要步骤 .....	34
3.4 RNA 不仅仅是信使 .....	35
3.4.1 核酶 .....	35
3.4.2 核小 RNA 和核仁小 RNA .....	36
3.4.3 端粒复制中的 RNA .....	36
3.4.4 RNAi 和非编码小 RNA .....	36
3.4.5 长非编码 RNA .....	39
3.4.6 其他非编码 RNA .....	40
3.5 细胞转录组学研究概览 .....	40

## 第二部分 新一代测序技术及数据分析概论

4 新一代测序技术的来龙去脉 .....	43
4.1 怎样做 DNA 测序：从第一代到新一代 .....	43
4.2 典型的 NGS 实验流程 .....	45
4.3 不同 NGS 测序平台的详细介绍 .....	48
4.3.1 Illumina 可逆染色终止子测序 .....	48
4.3.2 Ion Torrent 半导体测序 .....	52
4.3.3 PacBio 单分子实时测序 .....	53
4.4 测序的偏好性及其他影响 NGS 数据准确性的负面因素 .....	54
4.4.1 文库构建中的偏好性 .....	55
4.4.2 测序过程中的偏好性和其他因素 .....	56
4.5 NGS 的主要应用 .....	56
4.5.1 转录组特征和可变剪接检测 .....	56
4.5.2 遗传突变与变异的发现 .....	57
4.5.3 基因组的从头组装 .....	57
4.5.4 蛋白质与 DNA 的互作分析（ChIP-Seq）.....	57
4.5.5 表观基因组学与 DNA 甲基化研究 .....	57
4.5.6 宏基因组学 .....	58
5 新一代测序数据前期分析的常见步骤 .....	59
5.1 碱基识别、FASTQ 文件格式和碱基质量值 .....	60
5.2 NGS 数据的质量控制与处理 .....	61
5.3 读段的定位 .....	63
5.3.1 定位方法与算法 .....	63
5.3.2 定位算法和参考基因组序列的选择 .....	65
5.3.3 标准定位文件格式 SAM/BAM .....	66
5.3.4 定位文件的检验与操作 .....	67
5.4 第三阶段分析 .....	70
6 新一代测序数据管理与分析的计算能力需求 .....	71
6.1 NGS 数据的存储、传输与共享 .....	71
6.2 NGS 数据分析所需的计算能力 .....	72
6.3 NGS 数据分析所需软件 .....	74
6.4 NGS 数据分析所需的生物信息学技能 .....	75

### 第三部分 新一代测序数据分析的具体应用

7 转录组测序 .....	79
7.1 转录组测序的原理 .....	79
7.2 实验设计 .....	79
7.2.1 因子设计 .....	79
7.2.2 重复与随机化 .....	80
7.2.3 样本制备 .....	80
7.2.4 测序策略 .....	81
7.3 转录组测序数据分析 .....	82
7.3.1 数据质控与读段定位 .....	82
7.3.2 转录组测序数据的均一化 .....	84
7.3.3 差异表达基因的鉴定 .....	85
7.3.4 可变剪接分析 .....	87
7.3.5 转录组测序数据的可视化 .....	88
7.3.6 被识别基因的功能分析 .....	88
7.4 利用转录组测序发现新基因 .....	88
8 小 RNA 测序 .....	90
8.1 小 RNA 新一代测序数据生成和上游处理 .....	91
8.1.1 数据生成 .....	91
8.1.2 预处理 .....	92
8.1.3 定位 .....	92
8.1.4 小 RNA 的注释和预测 .....	93
8.1.5 均一化 .....	94
8.2 鉴别差异表达的小 RNA .....	94
8.3 已鉴定小 RNA 的功能分析 .....	94
9 用全基因组重测序方法分析基因型和发现基因组变异 .....	96
9.1 数据预处理、比对、再比对和再校准 .....	96
9.2 单碱基变异和 indel 检测 .....	98
9.2.1 SNV 检测 .....	98
9.2.2 新突变位点的检测 .....	99
9.2.3 Indel 检测 .....	99
9.2.4 转录组测序数据的变异检测 .....	101
9.2.5 变异检测格式文件 .....	101

9.2.6 评估 VCF 结果 .....	102
9.3 结构变异检测 .....	103
9.3.1 基于配对读段的 SV 检测 .....	103
9.3.2 断点的确定 .....	104
9.3.3 基于从头组装的 SV 检测 .....	104
9.3.4 CNV 检测 .....	104
9.3.5 综合 SV 分析 .....	105
9.4 检测变异的注释 .....	105
9.5 变异与疾病或性状关联的检验 .....	105
10 用新一代测序结果进行基因组从头组装 .....	107
10.1 从头组装的基因组因素与测序策略 .....	107
10.1.1 影响从头组装的基因组因素 .....	107
10.1.2 从头组装的测序策略 .....	108
10.2 重叠群的组装 .....	109
10.2.1 测序数据的预处理、错误修正与基因组特征的评估 .....	109
10.2.2 重叠群组装的算法 .....	111
10.3 组装骨架 .....	112
10.4 组装质量评估 .....	113
10.5 补齐缺口 .....	114
10.6 局限性与未来的发展 .....	114
11 用 ChIP-Seq 法对蛋白质 -DNA 互作定位 .....	116
11.1 ChIP-Seq 的原理 .....	116
11.2 实验设计 .....	118
11.2.1 实验对照 .....	118
11.2.2 测序深度 .....	118
11.2.3 重复 .....	118
11.3 读段定位、峰值确定与峰值可视化 .....	119
11.3.1 数据质控与读段定位 .....	119
11.3.2 峰值确定 .....	121
11.3.3 峰值可视化 .....	127
11.4 不同的结合点分析 .....	127
11.5 功能分析 .....	129
11.6 基序分析 .....	129
11.7 整合 ChIP-Seq 数据分析 .....	130
12 用新一代测序进行表观基因组学和 DNA 甲基化分析 .....	132
12.1 DNA 甲基化测序策略 .....	132

12.1.1	全基因组亚硫酸氢盐测序 .....	133
12.1.2	简化的亚硫酸氢盐测序 .....	134
12.1.3	基于甲基化 DNA 富集的甲基化测序.....	134
12.1.4	区分胞嘧啶甲基化与亚硫酸氢盐测序中去甲基化产物 ...	135
12.2	DNA 甲基化测序数据分析 .....	135
12.2.1	数据质量控制和预处理 .....	135
12.2.2	读段定位 .....	135
12.2.3	DNA 甲基化的定量 .....	137
12.2.4	DNA 甲基化数据的可视化 .....	138
12.3	甲基化胞嘧啶位点及差异区域的检测 .....	140
12.4	数据检验、核实和解析 .....	140
13	<b>用新一代测序进行宏基因组学研究 .....</b>	142
13.1	实验设计与样本制备 .....	143
13.1.1	宏基因组样本采集 .....	143
13.1.2	宏基因组样本制备 .....	144
13.2	测序方法 .....	145
13.3	全基因组鸟枪法宏基因组测序数据分析 .....	145
13.4	测序数据的质控和预处理 .....	147
13.5	微生物群落的分类学特征 .....	147
13.5.1	宏基因组的组装 .....	147
13.5.2	序列的分 bin .....	148
13.5.3	在宏基因组序列中识别可读框和其他基因组元素 .....	149
13.5.4	系统遗传学标记分析 .....	150
13.6	微生物群落的功能性特征 .....	150
13.6.1	基因功能注释 .....	150
13.6.2	代谢途径的重建 .....	151
13.7	比较宏基因组分析 .....	151
13.7.1	宏基因组测序数据均一化 .....	152
13.7.2	识别不同丰度的物种或操作分类单位 .....	152
13.8	整合宏基因组数据分析管道 .....	152
13.9	宏基因组数据库 .....	153
<b>第四部分 发展中的新一代测序技术与数据分析</b>		
14	<b>新一代测序将走向何方? .....</b>	157
14.1	发展中的新一代测序 .....	157

14.2 高通量测序数据分析的生物信息学工具的快速涌现与变化 .....	159
14.3 NGS 分析管道的规范化与流程化 .....	160
14.4 并行计算 .....	160
14.5 云计算 .....	161
<b>参考文献 .....</b>	<b>164</b>
<b>附录 A 新一代测序数据分析常用文件格式 .....</b>	<b>188</b>
<b>附录 B 词汇表 .....</b>	<b>190</b>

---

**第一部分**

**细胞与分子生物学概论**

---



# 1 细胞系统与生命密码

---

## 1.1 细胞面临的挑战

虽然细胞体积非常微小，直径只有不到  $50\mu\text{m}$ ，但其运行的精巧程度可以和一切人类设计的系统相媲美，并且它可以通过以 DNA 形式储存的密码来实现自身繁殖。假如我们希望设计一个和细胞同样复杂的系统，将需要面临很多看起来无法逾越的挑战。细胞具有复杂的内部结构，含有多种生物大分子组分，它需要完成一系列任务来维持自身系统的稳定，其中最根本的任务是在复杂多变的环境条件下，维持它的内部秩序，防止其出现故障或者崩溃，并完成自身繁殖，甚至进一步提升自身系统的功能。

要维持细胞内部系统的秩序，就需要输入能量。根据热力学第二定律，一个系统如果没有持续的能量输入，其熵值将逐渐上升，最终将导致系统的解体。同时，因为细胞内部结构的存在状态是动态的，需要响应不断变化的外部环境，所以除了能量输入以外，细胞还需要不断补充“原材料”来更新它的内部部件，或者制造新的部件。因此，细胞要维持其内部平衡及与所处环境之间的互动，就需要不断地输入能量与原材料，并且排出废物。细胞的系统信息储存在它的 DNA 序列中，用于指导细胞获取生存所需的能量和原材料，以及实现细胞繁殖等生命过程。

进化使得大量的生物种类都不再像单细胞生物那样简单。例如，人体就是由数千亿个细胞组成的。在多细胞系统中，细胞发生分化以执行特定的功能，如人体胰腺中的  $\beta$  细胞执行合成与释放胰岛素的功能，大脑中的皮层神经元执行学习和记忆的神经生物学功能。尽管功能多样，但对这些细胞的某个个体来说，它所面对的挑战与单细胞生物是一样的。不同的是，其不直接面对外部环境，而是需要面对其周边微环境的变化。

## 1.2 细胞如何面对挑战

很多细胞，如藻类与植物的细胞，直接从阳光或其他来源汲取能量。另外一些细胞（或者生物体）则作为异养生物从环境中汲取能量。在原材料方面，细胞可以利用简单有机化合物中携带的能量来固定空气中的二氧化碳，或者从环境中直接获取有机分子转化为其所需物质。同时，细胞内现存的物质也可以被降解，作为新的原材料被重新加以利用。这个能量的捕获、利用、合成、转化，以及分

子物质的降解再利用过程，构成了细胞代谢过程。代谢是细胞最基础的特征，包含了大量的生物化学反应。

对环境中各种信号的接收与转导对细胞的生存至关重要。细胞依靠其表面分布的受体接收信号，有些信号则需要依靠细胞内结构接收。信号转导到细胞内以后，通常在胞内引起一系列反应，使原始信号得以放大与调节，与之对应，细胞的代谢过程也发生相应的变化。细胞的信号接收与转导网络是由不同的代谢途径环路组成的。如果这些代谢途径失去功能，将会损害细胞对环境的反应能力，最终导致细胞死亡。

细胞系统依靠 DNA 复制和细胞分裂来完成其繁殖和进化。DNA 的复制（详见第 2 章）是一个高度保真的过程，但并非完全不出错。在维持系统稳定的同时，这个过程也提供了细胞系统分化与进化的机制。细胞分裂在各个阶段都是被精密调控的，以保证复制后的 DNA 可以平均分配到子细胞中。对于大多数有性生殖的生物来说，在其生殖细胞的形成过程中，DNA 只复制一次，而细胞却分裂两次，其结果是每个配子中 DNA 量减半。而父本 DNA 和母本 DNA 发生的重组，又使得子代进一步发生分化。

### 1.3 细胞内的各种分子

细胞中的各种过程需要不同类型的生物分子去实现。在一个典型的细胞中大部分是水，约占细胞总重量的 70%。除了水之外，细胞中还含有大量不同大小的分子。小分子主要包括无机离子（如  $\text{Na}^+$ 、 $\text{K}^+$ 、 $\text{Ca}^{2+}$ 、 $\text{Cl}^-$ 、 $\text{Mg}^{2+}$  等）、单糖、脂肪酸、氨基酸及核苷酸。绝大多数的大分子是多糖、脂类、蛋白质和核酸（DNA 和 RNA）。在上述组分中，无机离子对信号转导过程（如  $\text{Ca}^{2+}$  的波动是重要的胞内信号）、细胞能量的储备（如  $\text{Na}^+/\text{K}^+$  形成跨膜浓度梯度）、蛋白质的结构 / 功能（如  $\text{Mg}^{2+}$  是很多金属蛋白的关键辅助因子）至关重要。碳水化合物（包括单糖与多糖）、脂肪酸及脂类是细胞中提供能量的主要分子。其中脂类是细胞膜的主要成分。蛋白质由 20 种氨基酸以不同的顺序与长度组成，参与几乎所有的细胞活动，包括代谢、信号转导、DNA 复制和细胞分裂。它们还是很多胞内结构，如细胞骨架（详见 1.4 节）的组成部分。核酸以它几乎无穷的核苷酸排列顺序来储存生命的编码信息，它不仅指导细胞内蛋白质的合成组装，而且可以根据环境的变化对这些蛋白质如何组装进行调控。

### 1.4 细胞内结构与空间

细胞具有组织精巧的内部结构（图 1.1）。根据其内部结构的复杂程度，可将细胞分为两大类：原核细胞与真核细胞。它们之间的根本区别在于是否具有细胞

核。原核细胞是较为原始的形态，它没有细胞核，其DNA位于一个没有包被的区域，称为核区；原核细胞也没有细胞器。与之相比，真核细胞具有明显的细胞核，用于DNA的储存、维护与表达。真核细胞还具有各种特化的、可以将细胞各种功能分隔开的细胞器，如内质网（ER）、高尔基体、细胞骨架、线粒体及植物细胞特有的叶绿体等。这里我们对各种细胞内的结构与空间，包括细胞核、各种细胞器及其他亚细胞结构与空间，如细胞膜和细胞质等作一简单描述。

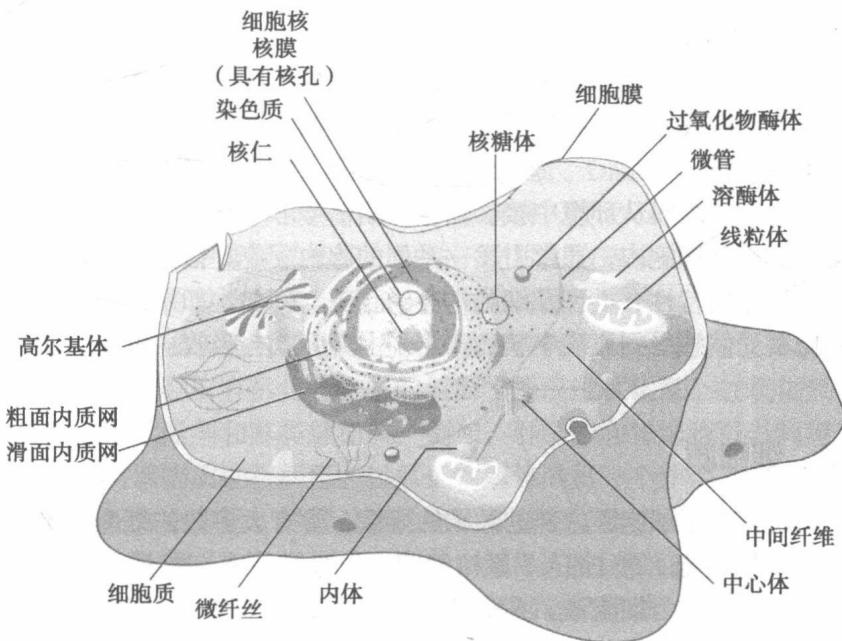


图 1.1 典型真核细胞的结构图，本图所示为一个动物细胞。

### 1.4.1 细胞核

DNA是储存生命密码的载体，因此它需要妥善维护，以避免可能遭受的损伤，并且需要保证它的准确性与稳定性。又因为将储存在DNA中的遗传信息准确地表达出来是细胞正常发挥功能的关键，所以基因表达在任何时候都需要受到严格调控。在绝大多数真核细胞中，细胞核都位于中心位置，它可以为DNA的储存、维护和基因表达提供一个受到良好保护的环境。细胞核区被双层的核膜包被，核膜上分布有跨内膜与外膜的核孔，通过这些核孔可以进行蛋白质和各种RNA的跨膜运输，这对于基因表达过程非常关键。细胞核骨架是由核纤层蛋白组成的网络结构，为细胞核提供机械支持。在细胞核内，细长的DNA分子缠绕在一种称为组蛋白的蛋白质上，被紧密包裹，以便存储在有限的细胞核空间里。在原核细胞中，有一个细胞核形状的不规则区域，称为核区，外部没有核膜包被。核区可以为DNA提供一些类似真核细胞核膜的保护（但其功能不如核膜）。