

引文格式：

林艺生, 武瑞君, 钱珑, 等. 基于纳米孔读取的 DNA 存储发展与展望 [J]. 集成技术, 2024, 13(3): 54-73.

Lin YS, Wu RJ, Qian L, et al. Development and prospects of DNA storage based on nanopore readout [J]. Journal of Integration Technology, 2024, 13(3): 54-73.

基于纳米孔读取的 DNA 存储发展与展望

林艺生¹ 武瑞君² 钱珑^{3*} 张成^{1*}

¹(北京大学计算机学院 北京 100871)

²(中国生物技术发展中心 北京 100039)

³(北京大学定量生物学中心 北京 100871)

摘要 现代社会正处于数据爆炸的时代, 全球对数据存储的需求已经远远超过了已有的存储能力。DNA 是一种天然的遗传信息载体, 可实现稳定、高效、低功耗的数据存储。目前的 DNA 存储过程主要分为 6 个环节: 编码、写入、保存、检索、读取、解码。纳米孔测序技术被广泛应用于读取 DNA 中所存储的信息。该文系统地介绍了纳米孔分子检测技术的原理和发展历史, 及其在 DNA 存储中的应用。此外, 该文总结了机器学习在纳米孔检测技术中的应用, 着重介绍了结合机器学习的新型纳米孔检测技术。该文为纳米孔检测技术的发展提供了新方向, 也为新型实用化 DNA 存储系统的发展奠定了基础。

关键词 纳米孔测序; 纳米孔检测; DNA 存储; 机器学习

中图分类号 Q 819 **文献标志码** A **doi**: 10.12146/j.issn.2095-3135.20231030003

Development and Prospects of DNA Storage Based on Nanopore Readout

LIN Yisheng¹ WU Ruijun² QIAN Long^{3*} ZHANG Cheng^{1*}

¹(School of Computer Science, Peking University, Beijing 100871, China)

²(China National Center for Biotechnology Development, Beijing 100039, China)

³(Center for Quantitative Biology, Peking University, Beijing 100871, China)

*Corresponding Authors: long.qian@pku.edu.cn; zhangcheng369@pku.edu.cn

Abstract In the era of data explosion, the global demand for data storage has far exceeded the existing storage capacity. DNA, as a natural carrier of genetic information, provides a stable, efficient, and sustainable

收稿日期: 2023-10-30 修回日期: 2024-02-25

基金项目: 国家重点研发计划项目 (2021YFF1200100); 国家自然科学基金项目 (62073133, 62273008, 12090054); 之江实验室 (2022RD0AB03); 预研项目 (31511090301)

作者简介: 林艺生, 硕士研究生, 研究方向为 DNA 存储、纳米孔测序分析、生物信息学; 武瑞君, 副研究员, 研究方向为生物技术与生物医药领域战略; 钱珑 (通讯作者), 副研究员, 研究方向为复杂生物系统的进化理论、计算合成生物学、DNA 分子信息存储, E-mail: long.qian@pku.edu.cn; 张成 (通讯作者), 副研究员, 研究方向为生物计算 (DNA 计算和分子电路)、DNA 存储、纳米孔测序、生物信息技术、纳米智能机器, E-mail: zhangcheng369@pku.edu.cn.

data storage solution. The current process of DNA storage is divided into six main parts: encoding, writing, preservation, access, reading, and decoding; and nanopore sequencing technology has been widely used to read information stored in DNA. This review systematically introduces the principle and research history of nanopore-based DNA signal detection, and the applications of nanopore sequencing in DNA storage. Moreover, it summarizes the applications of machine learning in nanopore detection, particularly highlighting the integration of novel nanopore detection techniques with machine learning. This review presents a new direction for the development of nanopore technology and lays the foundation for building a better practical DNA storage system.

Keywords nanopore sequencing; nanopore detection; DNA storage; machine learning

Funding This work is supported by National Key Research and Development Program of China (2021YFF1200100), National Natural Science Foundation of China (62073133, 62273008, 12090054), Zhejiang Lab (2022RD0AB03), Prestudy Project (31511090301)

1 引言

人类文明的发展离不开信息的传递与交互。通过获取和输出不同的信息, 人类认识并改变世界。目前正处于一个数据爆炸的时代, 数字化技术和互联网技术的蓬勃发展使得大量的数据可以被收集、存储和分析。根据国际数据公司(International Data Corporation)预测, 2025 年, 全球数据量将达到 175 ZB^[1], 这些海量数据的保存将需要大量的存储介质。然而, 在可预见的未来, 传统存储介质将不可避免地陷入资源枯竭的困境^[2]。因此, 寻找新的存储介质已经迫在眉睫。DNA 作为遗传信息的载体, 有着许多传统存储介质无法比拟的优点: 超高的信息密度(455 EB/g^[3])、超强的稳定性(半衰期超过 500 年^[4-5])、相对较低的维护能耗、可编程性和可寻址性^[6]。因此, 大量的 DNA 存储研究工作喷涌而出。

与传统的存储技术一样, DNA 存储过程同样分为 6 部分: 编码、写入、保存、检索、读取、解码^[7]。其中, 写入和读取通常对应着合成和测序。目前, 固相芯片合成技术和高通量

测序技术可以迅速、低价地大规模合成和读取 DNA^[8-9]。但高通量测序技术的测序长度十分有限(<200 bp)^[10], 所以, 样品的处理和大量短序列组装造成了较长的数据读取延迟。因此需要一种实时读取且精确可靠的 DNA 测序技术来实现分子存储的高效读取。

受分子识别和跨膜运输的生物过程启发, 科学家开发了纳米孔检测技术, 并将其作为单分子检测的超灵敏分析工具^[11-13]。通常情况下, 单分子在外加电位的影响下进入纳米孔, 在孔中结合或反应, 从而引起电流变化。通过分析电流在幅度、持续时间和频率等特征上的变化, 纳米孔可被用于检测 DNA^[14-18]、RNA^[19]、肽^[20]、蛋白质^[21-22]、代谢物和蛋白质-DNA 复合物^[23]等。其中, 纳米孔检测技术在 DNA 测序上的应用最为成熟, 并在商用中取得了巨大成功。事实上, 牛津纳米孔科技(Oxford Nanopore Technologies, ONT) 已经于 2014 年开发出了商用测序仪 MinION^[24]。目前, 纳米孔检测技术已被广泛用于 DNA 存储系统, 用于实现海量分子数据的精确访问和读取。

本综述将对以下方面进行系统性阐述: (1) 纳

米孔检测技术的发展；(2) 纳米孔检测技术的应用；(3) 基于纳米孔数据读取的 DNA 存储技术；(4) 机器学习在纳米孔检测技术中的应用。本研究团队希望通过本文的综述推进纳米孔检测技术和 DNA 存储技术的发展。

2 纳米孔检测技术的发展

纳米孔检测技术是一项起源于 1976 年的电化学单分子检测技术^[25]，通过实时检测纳米孔的电流，可以分析出孔内分子的功能、结构、动力学等信息。与其他单分子检测技术相比，纳米孔检测技术可以在不破坏原有分子结构的情况下无标记地检测分子^[26]，并且十分便携(ONT 的 MinION 测序仪仅 90 g)^[24]。

2.1 纳米孔检测仪的一般原理

纳米孔检测仪的工作原理类似于 Coulter 计数器，单个带电分子在电场力的作用下通过镶嵌在膜上的纳米孔产生瞬时阻断电流^[27]。通常情况下，纳米孔蛋白被固定在具有电阻性的聚合物膜中，并被浸泡在电解质溶液中。当施加恒定电压时，电解质离子会通过纳米孔，并在膜两侧产生电流。带负电荷的单链 DNA 或 RNA 分子会从负极一侧(*cis*)经过纳米孔向正极一侧(*trans*)移动。转移速度由一种马达蛋白(motor protein)控制，它会以一定的方式使核酸分子通过纳米孔^[28]。在核酸分子的转移过程中，离子电流的变化与纳米孔中存在的核苷酸序列相对应，因此可以使用特定的算法将离子电流解码成核苷酸序列，从而实现单个分子的实时测序。除了控制转移速度外，马达蛋白还具有解旋酶的作用，能够使双链 DNA 或 RNA-DNA 双链解旋成单链分子，以通过纳米孔。

通过多个纳米孔组成阵列，可以实现大规模测序。以 ONT 的 MinION 为例(如图 1 所示)，每一个纳米孔集成于微支架(microscaffold)，多个微

支架组成微支架阵列(array of microscaffolds)，可使多个纳米孔在运输和使用过程中保持稳定。每个微支架与其自身电极相对应，该电极连接至传感器阵列芯片(sensor chip)中的通道上。每个纳米孔通道均由定制的专用集成电路单独控制和测量，支持同时进行多个纳米孔实验。

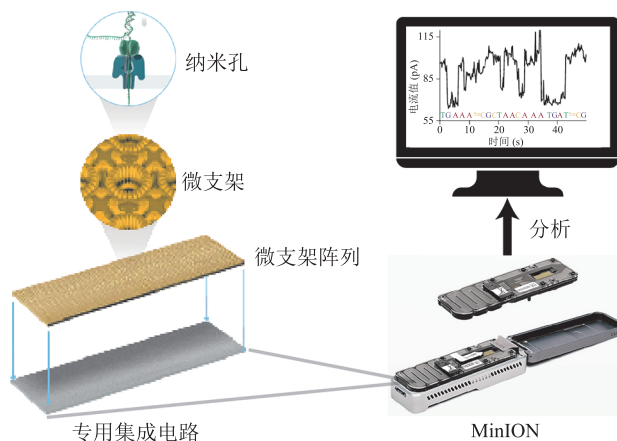


图 1 纳米孔测序原理^[29]

Fig. 1 Principle of nanopore sequencing^[29]

2.2 纳米孔的种类

目前，纳米孔主要有两大类：生物纳米孔和固态纳米孔。生物纳米孔由肽、蛋白质亚基在脂质双分子层或多聚物膜中自组装形成，能够识别直径为 1~10 nm 的分子^[11-12,26]。固态纳米孔是在超薄的人造膜(如 SiN_x)中制造而成的，这使得纳米孔的直径可以扩大至数百纳米，因此可以分析大型生物分子和复合物。用于制造固态纳米孔的技术(如电子/离子铣削^[30-31]和激光蚀刻^[32-33])可以在纳米尺度上操纵纳米孔的大小，但是与生物纳米孔相比，固态纳米孔只能操纵表面结构。然而，生物纳米孔的化学修饰可以进一步增强纳米孔与特定分子之间的相互作用，提高整体灵敏度和选择性，使得纳米孔可以在溶液中可控地捕获、识别、运输各种分子^[12,34-36]。但生物纳米孔也有缺点：缺乏调整蛋白质孔径和几何结构的自由，以及只有在一定的 pH、温度和离子浓度下，才具有有限的稳定性。

一方面, 生物纳米孔主要有 α -HL (α -hemolysin)^[37]、MspA (Mycobacterium smegmatis porin A)^[38]、Ael (aerolysin)^[39]、phi29 motor^[40]、ClyA (cytolysin A)^[41]、OmpG (outer membrane protein G)^[42]等。就孔径(蛋白质孔隙在天然形态下的最窄收缩直径)而言, α -HL (~1.4 nm)、MspA (~1.2 nm)、Ael (~1.0 nm) 和 OmpG (~1.3 nm) 适用于单链 DNA 的研究, 但是不适用于较大的分析物, 如双链 DNA 和折叠蛋白质。而 phi29 motor (~3.6 nm) 和 ClyA (~3.3 nm) 就足够通过双链 DNA 分子。另一方面, 纳米孔的长度对检测灵敏度至关重要。在 DNA 链的检测中, 孔中核苷酸的离子电流在纳米孔中移动时受邻近核苷酸的影响, 如果孔道长度过长, 就会导致电流信号很难解码成核苷酸。例如, 在双链 DNA 中, 两个碱基之间的距离约为 0.34 nm, 在具有 5 nm 长的 β -管的 α -HL 孔中, 将同时存在 15 个左右的碱基, 会有 4^{15} 种碱基组合, 导致难以解码。而 MspA 只有 0.6 nm 的收缩区, 因此是最适用于 DNA 测序的纳米孔之一。

固态纳米孔的制作材料主要有 SiN_x 、 SiO_2 、Si-C、 Al_2O_3 、石墨烯等。与生物纳米孔相比, 固态纳米孔在稳定性、电流噪声等方面有着显著优势, 但是受限于如今的半导体工艺制造水平, 固态纳米孔的制造还较为复杂和昂贵^[30-31]。总之, 每种纳米孔都有着不同的特性, 应该针对不同的任务应用合适的纳米孔。

3 纳米孔检测技术的应用

最初, 纳米孔检测技术是为了实时检测离子和小分子^[12,43], 随后, 纳米孔检测的发展集中于 DNA 测序^[11,14-16]。然而, 纳米孔的应用远不止于测序, 它已经被用于研究许多生物化学系统中分子的特性。第一, 纳米孔的一个关键优势在于它能够以较高速率、有顺序地连续捕获大量单

分子, 这使得纳米孔能够在短时间内检测许多分子, 从而实现实时检测。第二, 纳米孔将分析物的结构和化学性质转化为可测量的离子电流信号, 并可以通过一定的算法实现电流信号到单分子的映射^[44]。综合以上两个特点, 纳米孔检测技术可用于研究多种分子的特征, 并且是无标记识别, 不会影响分子的原有的特征。此外, 纳米孔能够识别缺乏合适信号放大标签或其信息隐藏在噪声中的单分子^[45]。因此, 纳米孔可以很好地用于精准医学所需的分子诊断应用, 实现核酸、蛋白质或代谢物等分析物和其他生物标志物的识别。

3.1 核酸识别

纳米孔测序的概念出现在 20 世纪 80 年代, 并通过纳米孔和相关马达蛋白的一系列技术进步实现^[46]。 α -HL 是第一个可以检测到 RNA 和 DNA 阻塞电流的纳米孔^[11]。野生型 α -HL 可以识别 DNA 的 4 种碱基, 这是迈向单核苷酸分辨率的关键一步。但是 α -HL 的孔道较长, 无法识别复杂序列。在图 2(a) 中, Lieberman 等^[18]通过引入 phi29 DNA 聚合酶, 极大地减缓了 DNA 的过孔速度, 从而实现了复杂序列的识别。与之前缺乏控制的 DNA 过孔实验相比, 马达蛋白的引入不仅减缓了过孔速度, 而且减少了过孔时 DNA 链的动力学波动, 从而提高了数据质量^[29]。

除了通常的 DNA、RNA 测序外, 目前的纳米孔测序技术还能够检测非天然核酸。非天然核酸是一种人工合成的核酸类似物, 具有不同于天然核酸 DNA 和 RNA 的糖骨架。在图 2(b) 中, Yan 等^[35]利用 NIPSS (nanopore-induced phase shift sequencing) 实现了 FANA (2-deoxy-2-fluoroarabinoic acid) 的直接测序。通过将 FANA 与 DNA 驱动链连接, 他们发现, 利用 phi29 DNA 聚合酶的 NIPSS 可以实现 FANA 的直接测序。Yan 等^[35]的研究证实, 纳米孔技术能够清晰地识别并检测 DNA、RNA 及非天然核酸。

3.2 蛋白质识别

与 DNA 测序技术相比, 蛋白质纳米孔测序技术面临更多挑战, 主要原因是: 蛋白质的三维结构复杂性、多样化的电荷分布及更庞大的氨基酸种类。尽管如此, 在克服这些难题上, 科研人员已经取得了一系列进展。Yusko 等^[21]利用双层固态纳米孔确定了单个蛋白质的近似形状、体积、电荷、旋转扩散系数和偶极矩。此外, 配体(如生物素^[47]、蛋白结构域^[48]或抗体^[49])可以直接附着在生物纳米孔上, 即使在复杂溶液(如血清)中, 也能识别特定氨基酸或蛋白质。在图 2(c)中, Fahie 等^[50]利用生物素修饰的 OmpG

生物纳米孔能够特异性地识别两种抗生素抗体。此外, Bell 等^[51]还可以利用经过蛋白质特异性结合物修饰的 DNA 载体识别蛋白质。针对蛋白质线性化及电渗效应驱动的过孔问题, 一种策略是借助高浓度变性缓冲液构建纳米孔平台, 使得蛋白质伸展变直, 并产生足够的电渗动力, 推动长度为 250~750 个氨基酸的蛋白质穿越 α -HL^[52]。另一种策略是在 CytK 纳米孔内部设计带负电荷的环结构, 诱导电渗流, 以促进蛋白质有效穿过纳米孔^[53]。

近年来, 科研人员在精确识别 20 种天然氨基酸方面取得了重大突破。早在 2020 年, 野生型 Ael

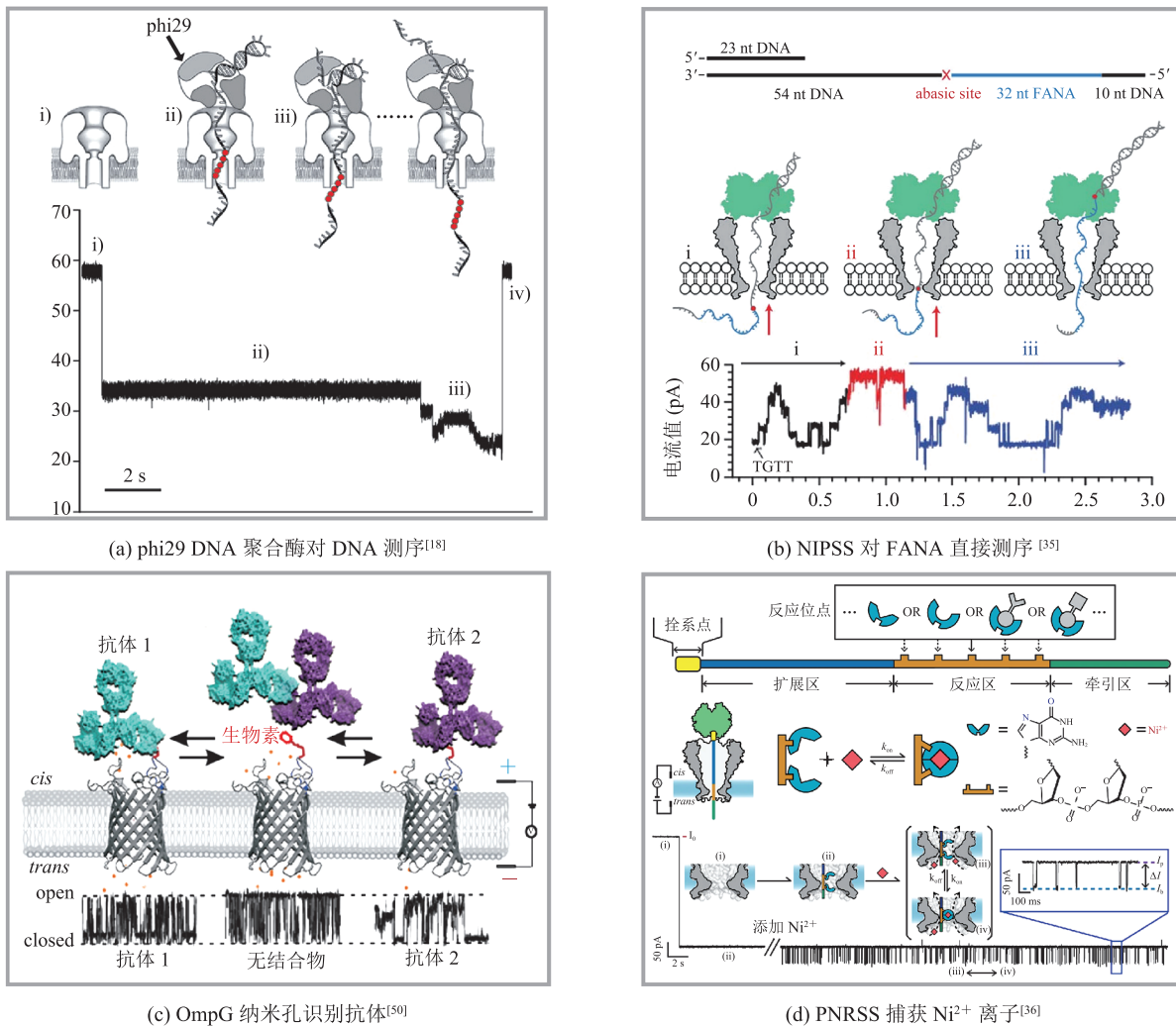


图 2 纳米孔检测技术的应用

Fig. 2 The application of nanopore detection technology

纳米孔已被证实可以识别 13 种氨基酸^[54]。最近的两项研究进一步实现了对 20 种氨基酸的明确鉴定, 并且还进行了短肽的序列测定。改造过的 MspA 异八聚体纳米孔不仅能鉴别 20 种氨基酸, 还能识别 4 种蛋白质翻译后修饰^[55]。另一项研究则展示了含苯丙氨酸的肽探针与葫芦脲复合物如何通过 α -HL 成功区分 20 种氨基酸^[56]。

这些显著的科学进展不仅证明了纳米孔检测技术在解决蛋白质三维结构复杂性、电荷分布多样性及氨基酸种类繁多等难题上的可行性, 而且为未来提升信息存储密度开辟了新的可能。随着对 20 种天然氨基酸及多种蛋白质翻译后修饰的识别能力的不断提升, 可以预见, 蛋白质分子将成为 DNA 存储之外的新型信息载体。理论上, 蛋白质序列的多样性远超 DNA 序列, 这使得每一种氨基酸组合都可能编码更多信息, 从而大幅提高信息存储容量。此外, 蛋白质特定的修饰状态可以进一步扩展编码空间, 实现更复杂和多层次的信息存储。

3.3 无机分子识别

除了检测生物分子外, 纳米孔还可用于检测无机化学分子。Jia 等^[36]于 2021 年提出了一种基于纳米孔测序技术的可编程随机纳米传感器 (programmable nano-reactors for stochastic sensing, PNRSS), 可以在单分子水平上实时监测化学反应。通过 PNRSS 可以直接观察到金属离子、简单有机化合物 (如乳酸) 和核苷类似物等广泛的单分子反应。如图 2(d) 所示, PNRSS 可以识别 Ni^{2+} 。此外, 他们利用机器学习的方法提高了 PNRSS 的传感分辨率, 能够实现多种分子的同时识别。

3.4 生物标记物检测与定量

与实验室制作的样品相比, 医学标本具有更大的复杂性和异质性, 且需要极高的精度、特异性和敏感性。由于纳米孔的灵敏性和可编程性较好, 因此有望作为新型的体液检测仪。利用

α -HL 纳米孔可以用于选择性检测溶液中杂交的 microRNA 分子与寡核苷酸探针原理, Wang 等^[57]从纯化的肺癌患者血浆样品中成功定量检测了 microRNA-155 生物标志物。此外, Rozevsky 等^[58]开发了 RT-qNP (reverse transcription quantitative nanopore sensing), 一种无须纯化与扩增的 RNA 定量方法。RT-qNR 在人类细胞系中准确量化了癌细胞转移相关基因 MACC_1 和 S_{100A_4} , 而且还实现了 SARS-CoV-2 的无聚合酶链式反应检测。

4 基于纳米孔数据读取的 DNA 存储技术

携带数据信息的 DNA 链虽然在保存过程中可能受到外界环境因素 (如紫外线、极端温度变化、细菌和病毒的生物污染) 的干扰和破坏, 但是, DNA 具有超高的信息密度和超长的寿命, 因此依然有望成为下一代新型存储介质。DNA 保护技术和保护设备的发展将大大提高 DNA 存储的存储时间和存储质量^[59-61]。目前, 纳米孔可以识别 DNA/RNA 的序列, 以及序列上的修饰, 而序列及其上的修饰均可被认为是信息位, 因此, DNA 存储主要分为两类: (1) 基于序列的 DNA 存储; (2) 基于结构/修饰的 DNA 存储。

4.1 基于序列的 DNA 存储

如图 3(a) 所示, DNA 存储的流程通常分为 6 部分^[7]。(1) 编码: 将数据通过特定算法编码成 DNA 序列。(2) 写入: 使用 DNA 合成技术合成编码信息的 DNA 链。(3) 保存: 在体内/体外恰当地保存 DNA。(4) 检索: 随机访问特定 DNA 链。(5) 读取: 利用测序技术读取所访问 DNA 链上的信息。(6) 解码: 从测序数据解码信息。目前, 纳米孔检测技术被普遍用于读取步骤。

2018 年, Organick 等^[62]编码并存储了 35 个不同的文件 (总共 200 MB), 共使用了 1 300 万个寡核苷酸, 并且使用纳米孔测序, 无错误地随机访问了每个文件。此外, 他们还设计了新的解码

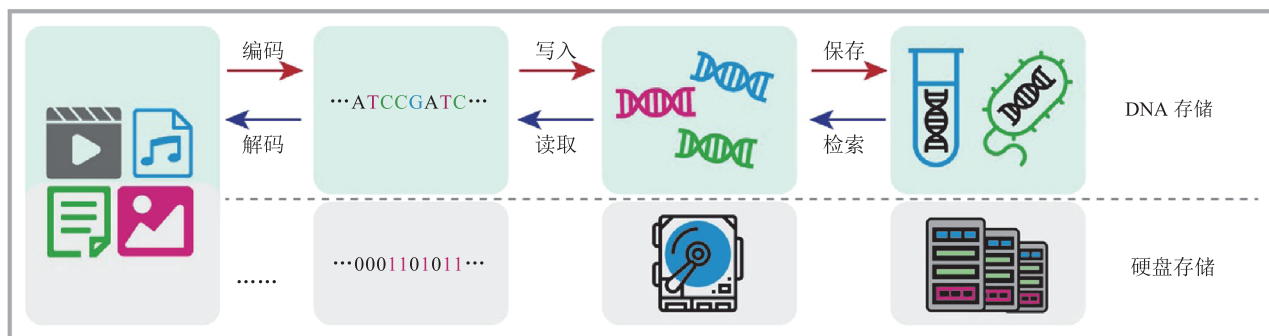
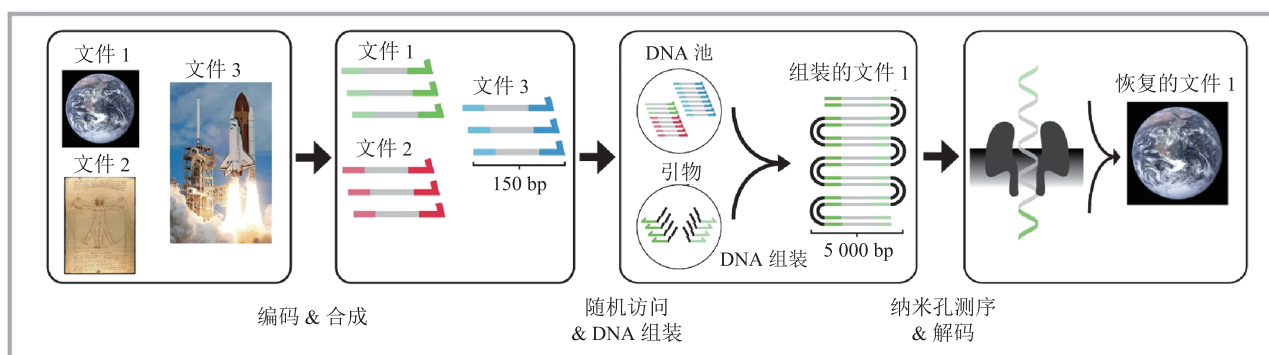
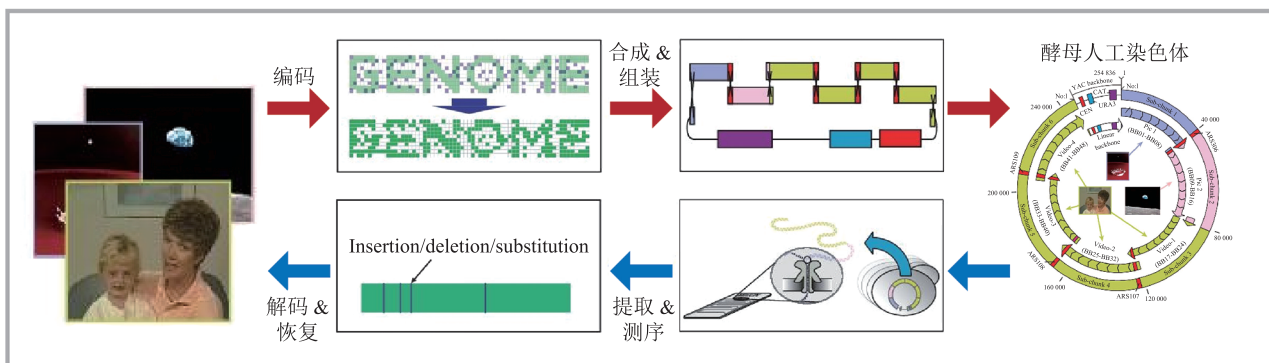
(a) DNA 存储流程与传统硬盘存储类比^[7](b) Lopez 等^[63]的 DNA 存储流程(c) 酵母人工染色体实现 DNA 存储的流程^[64]

图 3 基于合成的 DNA 存储系统

Fig.3 Synthetic-based DNA storage systems

算法，以最大化读出信息，从而降低无错误解码所需的测序覆盖率。2019年，Lopez 等^[63]使用便携式纳米孔测序设备 MinION 成功解码了存储在 DNA 中的 1.6 MB 信息。如图 3 (b) 所示，他们结合 DNA 自组装和纳米孔测序，大大提高了纳米孔测序的吞吐量，并且这种组装技术适用于任何短 DNA 扩增子的纳米孔测序。

不仅仅局限于 DNA 体外存储，纳米孔测序还

能用于 DNA 体内存储。在图 3 (c) 中，Chen 等^[64]设计并合成了一条酵母人工染色体 (25.4 万个碱基)，实现了 DNA 体内存储。他们结合 LDPC (low-density parity-check) 和伪随机序列设计编码算法，可以在一定程度上容忍测序时碱基的插入/删除错误，十分适用于纳米孔测序。通过这一编码算法，在 10.79% 的错误率下，仍然实现了数据的可靠恢复。此外，他们还将该人

工染色体繁衍了 100 代, 证明了其稳定性。2023 年, Sun 等^[65]利用细菌基因组构建便携式 DNA 数据存储系统。结合 Reed-Solomon 码和 RaptorQ 码, 他们设计了 MEPCAL (Mixed Error Processing Coding for Arbitrary Length) 编解码算法, 将 5.56 KB 的压缩数据编码成了 50 540 bp 的 DNA 序列。这些序列在合成以后, 通过他们所开发的 RSGE (recombinase-based site-specific genome engineering) 工具箱, 迭代整合进大肠杆菌和蓝晶盐单胞菌的基因组中。最后利用纳米孔测序实现了无损信息恢复。从这两项工作中可以看出体内存储的可行性。

4.2 基于结构/修饰的 DNA 存储

纳米孔测序技术可以检测 DNA 链上的修饰及 DNA 的二级结构, 为 DNA 存储开拓了一个新方向。可编程的 DNA 纳米结构可以作为一种检索数据的方式^[66-67]。如图 4(a) 所示, 2018 年, Chen 等^[68]研发了一种利用 DNA 发夹结构 (DNA hairpins) 的 DNA 存储方案。在这项工作中, 不同长度的 DNA 发夹被认为是不同的码元, 并通过固态纳米孔识别不同的发夹结构。利用内径约为 5 nm 的石英纳米孔可以识别长度为 8 bp 和 16 bp 的 DNA 发夹。因此, 将 8 bp 对应比特 0, 16 bp 对应比特 1, 将 56 个发夹连接到 7 228 bp 的 DNA 链上, 可以存储 56 bits 的信息。同样是利用 DNA 发夹结构, 2016 年, Bell 等^[51]基于 DNA 折纸 (DNA origami) 原理设计了一个 DNA 纳米结构库。库中的每个成员都包含了一个独一无二的条形码, 其中, 每一位由多个 DNA 哑铃发夹结构的存在或不表示。通过固态纳米孔, 他们实现了 3 bits 长度的条形码识别, 准确率达到 94%。

除了利用 DNA 的纳米结构和传统的碱基外, 还有利用生物聚合物序列的 DNA 存储技术。2020 年, Cao 等^[69]利用专门定制的生物聚合物序列作为比特信息存储载体。如图 4(b) 所示,

生物聚合物序列是由两种不同的单体 (*n*-propyl phosphate 和 [2,2-diynyl]-propyl phosphate) 和天然核苷酸组成的生物杂交大分子, 其中, 两种单体分别对应 0 和 1。他们利用 α -HL 纳米孔成功识别了单比特分辨率的生物聚合物。此外, Cao 等^[69]使用深度学习的方法实现了 4 bits 的高精度编解码。

基于传统的 4 种碱基 DNA 存储系统, 2022 年, Tabatabaei 等^[70]引入了 7 种化学修饰的核苷酸, 将分子字母表的大小从 4 扩大到了 11, 信息密度提高了 1.73 倍。如图 4(c) 所示, 除了 4 种传统碱基 A、C、G、T 外, 还引入了 7 种化学修饰碱基。在扩大分子字母表后, 他们利用 MspA 纳米孔和链霉亲和素减缓过孔速度, 成功识别了 4 bp ($4 \times \log_2 11 \approx 14$ bits) 的信息。更进一步, 他们运用了 ONT 公司研发的 GridION 测序平台, 并结合深度学习技术, 成功实现了对多达 66 种不同信息的同时识别与解析。

4.3 其他 DNA 存储读取技术

除了纳米孔测序外, DNA 存储的读取技术还包括 Sanger 测序^[71]、二代测序^[72]、单分子实时测序^[73]等。

Sanger 测序^[71]是一种基于链终止法的 DNA 测序技术。它通过在聚合酶链式反应中引入不同长度的链终止核苷酸, 使得每个 DNA 片段具有不同的长度, 并在凝胶电泳中分离, 从而确定碱基顺序。Sanger 测序的准确率高, 但速度慢、成本高, 且难以实现大规模并行化操作, 对于大数据量的 DNA 存储应用来说, 效率较低。

二代测序^[72] (如 Illumina 公司的边合成边测序技术) 利用大规模并行化的原理快速生成大量短序列数据, 然后通过算法拼接成完整序列。二代测序技术具有高通量、低成本的特点, 适合基因组研究和部分生物医学应用; 但在 DNA 存储领域, 由于其读长短、需要复杂的序列组装过程, 因此在大数据的读取中存在较高的延迟, 需要在算法

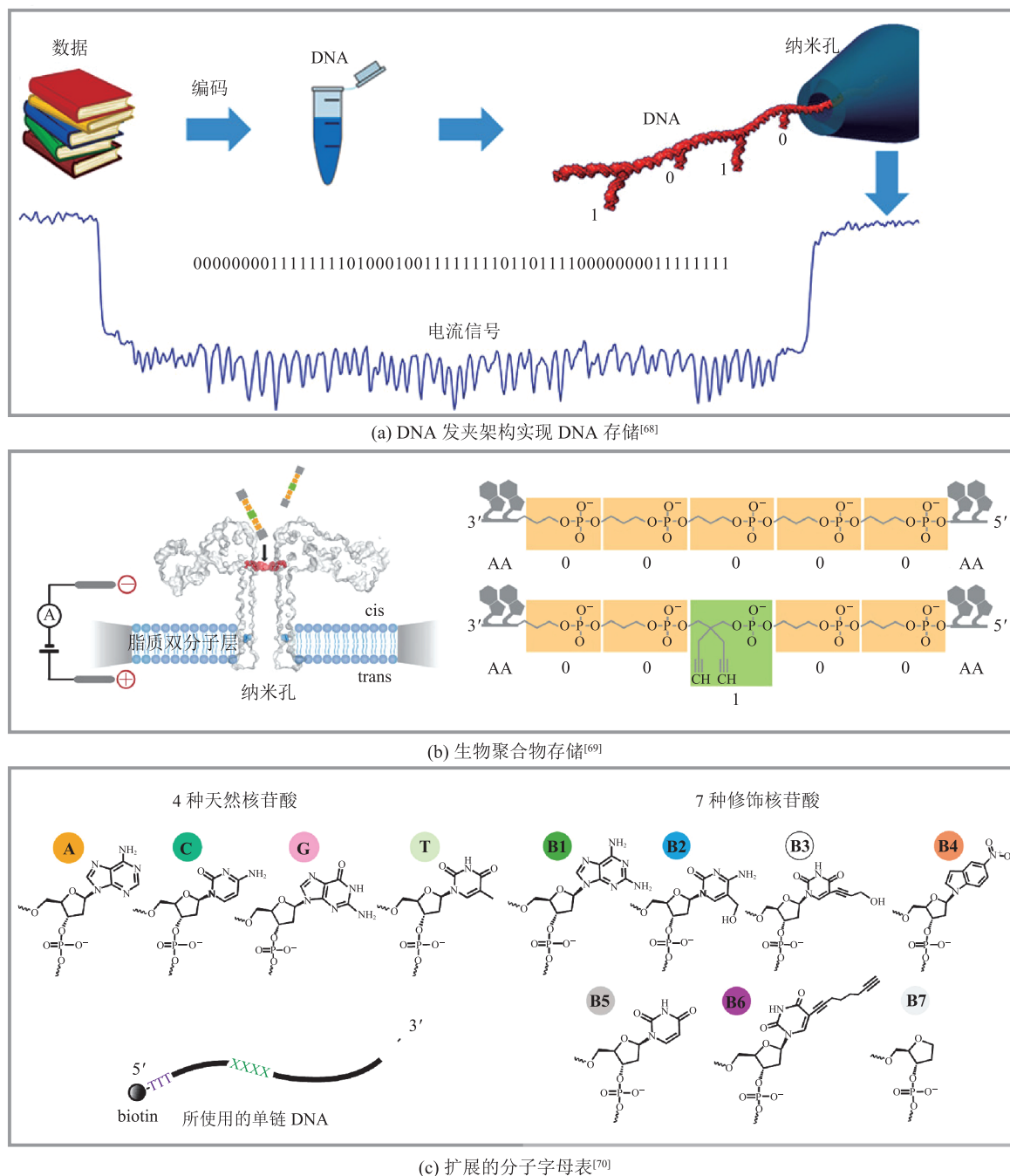


图 4 基于结构/修饰的 DNA 存储系统

Fig. 4 Structure/modification based DNA storage systems

和读取技术上继续探索这一问题的解决方案。

单分子实时测序 (single molecule real-time sequencing, SMRT)^[73]能够对单个 DNA 分子进行长时间连续观测,并记录其合成过程中的荧光信号变化,以确定碱基序列,尤其擅长读取长序

列。SMRT 的显著优势在于能通过 HiFi 模式获得串联的多轮测序结果,从而在 DNA 测序时达到 99.9% 的准确率,且其错误模式没有明显偏好性,易于修复,这对复杂区域的解码或 DNA 存储中连续数据的读取极其有利;然而,该技术的

通量较低, 读长较纳米孔测序更短, 单位成本比其他技术更加高昂。通过提高通量, SMRT 有望继续在需要长序列分析的 DNA 存储场景下发挥作用。

5 机器学习在纳米孔检测技术中的应用

近年来, 机器学习(machine learning)在生物信号处理中发挥着巨大作用^[74]。深度学习(deep learning)是机器学习领域中的一个研究方向, 通过组合低层特征形成更加抽象的高层表示属性类别或特征, 以发现数据的分布特征表示。深度神经网络、卷积神经网络(convolutional neural network, CNN)和递归神经网络(recurrent neural network, RNN)等深度学习模型已经在解决诸如视觉识别(visual recognition)、语音识别(speech recognition)和自然语言处理等很多问题方面都表现出非常好的性能^[75]。而纳米孔测序会产生大量的原始测序数据, 因此, 许多针对纳米孔测序进行数据分析的深度学习模型喷涌而出, 这些模型能够精确地进行测序数据的处理, 如 Basecalling、错误校正、修饰碱基检测等。

5.1 Basecalling

Basecalling 将纳米孔测序的原始信号解码成核苷酸序列, 这对数据精确性和碱基修饰的检测至关重要。ONT 占有了最大多数的纳米孔测序市场, 因此, 以下围绕 ONT 测序数据的 Basecalling 工具进行综述。Basecalling 的发展历史主要分为 4 个阶段^[76-80]: (1) 2016 年, 利用隐马尔可夫模型和递归神经网络对分割后的电流数据进行 Basecalling; (2) 2017 年, 直接对原始电流信号进行 Basecalling; (3) 2018 年, 使用一种 flip-flop 的模型识别单个核苷酸; (4) 2019 年, 根据需求数据训练个性化的 Basecalling 模型。目前, ONT 开发了一系列的 Basecalling 工具作为“技术演示”的软件(如 Nanonet、Scrappie 和

Flappie), 这些工具随后被投入到正式可用的软件包中(如 Albacore 和 Guppy)。

ONT 的测序仪可以进行每秒数千次的电流测量。DNA 或 RNA 分子过孔时会有多个连续的核苷酸(即 k-mer)在孔内, 导致该时刻的电流由多个核苷酸共同决定。原始电流可以根据电流的变化分割成电流信号段, 每个电流信号段对应着一个 k-mer。一个电流信号段包含了多个原始信号值, 这一段信号的平均值、方差和持续时间共同构成事件(event)数据。每个事件对相邻核苷酸的依赖性与马尔可夫链十分相似, 因此, 隐马尔可夫模型自然而然地在早期的 Basecalling 中得到应用, 例如, ONT 的 Nanocall^[81]。ONT 随后的 Nanonet 和 DeepNano^[82]训练 RNN 从事件数据中预测 k-mer, 从而提高 Basecalling 的准确率。其中, Nanonet 还是一个双向递归神经网络(bidirectional RNN), 结合上下游的信息更加精准地进行 Basecalling。

然而, 在原始电流信号到事件的转换中, 会丢失部分信息, 有可能降低 Basecalling 的准确性。ONT 的开源软件 Scrappie(同时用于 Albacore 和 Guppy)和第三方软件 Chiron^[83]率先实现了直接将原始电流信号转换为 DNA 序列。接着, ONT 发布了 Flappie, 该模型使用 flip-flop 模型直接从原始信号中识别单个碱基, 实现了单碱基分辨率的 Basecalling。此外, 如图 5(a)所示, Causalcall 软件通过使用时间卷积网络(temporal convolutional network)和连接时间分类(connectionist temporal classification, CTC)解码器识别远程序列特征, 改进 Basecalling^[84]。与固化的 Basecalling 模型相比, ONT 引入了 Taiyaki(用于 Guppy), 以实现个性化 Basecalling, 例如, 针对特殊序列/物种训练个性化模型。Taiyaki 通过使用自然语言处理的先进技术处理原始电流信号的高度复杂性和长期依赖性。不仅仅止于 Basecalling, Taiyaki 还可以通过训练模型识别修

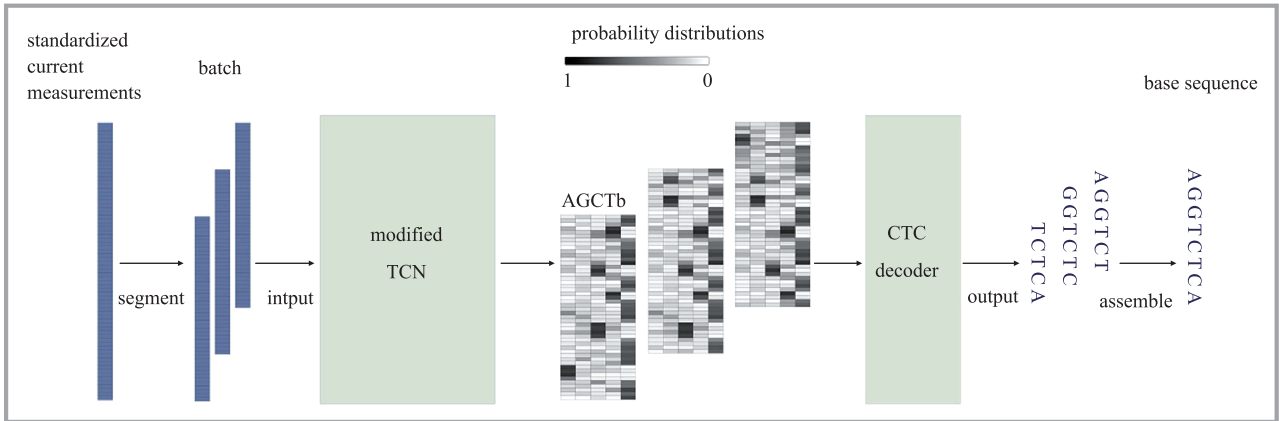
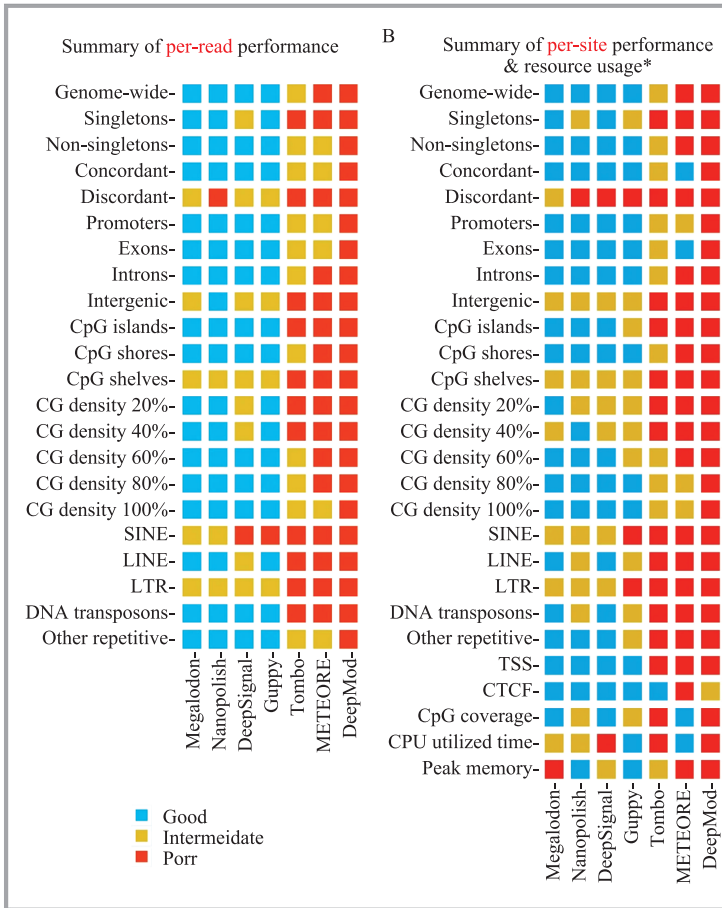
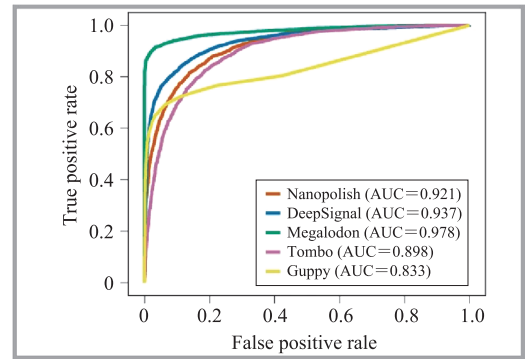
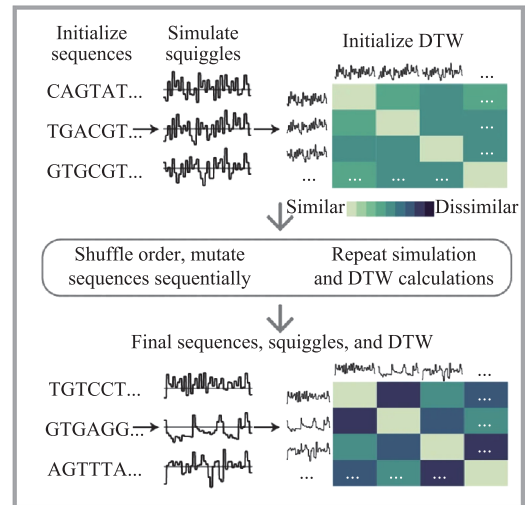
(a) Casualcall 神经网络架构^[84](b) 修饰检测工具 per-read 水平和 per-site 水平基准测试^[103](c) 修饰检测工具 ROC 对比^[104](d) Porcupine 设计序列算法流程^[109]

图 5 机器学习在纳米孔检测技术中的应用

Fig. 5 The applications of machine learning in nanopore detection technology

饰的碱基，如 5-甲基胞嘧啶 (5-methylcytosine, 5mC) 和 N⁶-甲基腺嘌呤 (N⁶-methyladenine, 6mA)。目前，ONT 的 Guppy 是最为广泛使用的 Basecalling 工具，并且可以利用显卡加速运

算，与之前的 Basecalling 工具相比，在准确性和速度上有着无可比拟的优势^[76]。ONT 在原有的 Guppy 基础上开发了新一代架构的 Basecalling 软件 Dorado。与 Guppy 相比，Dorado 在保持计算

资源成本相对可控的基础上, 提升了 3.8 倍的性能^[85]。不仅如此, Dorado 还增强了对修饰碱基的检测能力, 能够在进行 Basecalling 的同时, 以较小的额外计算成本实现对 5mC 及 5-羟甲基胞嘧啶(5hmC)等修饰碱基的高灵敏度检测, 从而为表观遗传学研究提供更为强大的支持。

5.2 序列纠错

尽管纳米孔测序的平均准确率正在提高, 但部分序列的准确率仍然非常低, 如 ONT 的 1D 和 2D/1D² 短序列的错误率远高于二代测序生成的数据。因此, 需要对测序结果进行序列纠错(error correction), 以保证下游分析的正确性。目前主要有两种基于传统算法的序列纠错方案^[86-87]: 自我纠错(self-correction)和混合纠错(hybrid-correction)。自我纠错技术通过运用基于图谱的算法, 对来自同一源头的不同分子构建一致性序列, 并将这些序列与源自同一分子的纳米孔测序结果进行比对, 借此实现错误修正, Canu^[88]和 LoRMA^[89]等工具采用了此类方法。而混合纠错利用高准确率的短序列, 通过基于对齐的算法(如 LSC^[90]和 Nanocorr^[91])、基于图的算法(如 LoRDEC^[92])、基于双端对齐/图的算法(如 HALC^[93])纠正长序列。混合纠错工具结合足够的短读段覆盖度, 能够将长读段的错误率降低到与短读段相似的水平(为 1%~4%)^[86-87], 而自我纠错则能将错误率降至 3%~6%^[87], 这可能是 ONT 数据中存在非随机的系统性误差所致。

除了传统算法外, 还有科研人员利用机器学习算法进行序列纠错, 如 NanoReviser^[94]。Wang 等^[94]通过 CNN 和双向长短期记忆网络(Bi-LSTM), 充分利用原始电信号及碱基识别器提供的碱基识别序列中的信息, 在经过训练后, NanoReviser 能使 *E. coli* 数据集和人类数据集的测序错误率均降低 5% 以上, 并且在加入甲基化注释后, NanoReviser 使 *E. coli* 数据集的错误率降低了 7%。

结合传统算法和机器学习算法, 对 Basecalling 的结果进一步进行序列纠错, 可以明显降低序列信息的错误率, 从而提高从 DNA 存储到基因组学各个领域的的数据解析质量。

5.3 基因组从头组装

纳米孔的长序列测序有着较高错误率, 尽管如此, 随着技术的不断进步和算法优化, 科研人员已经开发出多种策略来克服纳米孔测序错误率高的问题。

Canu^[88]不仅依赖于传统的 overlap-layout-consensus 框架, 还创新性地运用了自适应错误修正算法, 在组装过程中实时识别并纠正碱基识别错误, 从而显著提升了长序列数据在基因组从头组装中的应用效果。miniasm^[95]则以其极简快速的特点著称, 尤其适用于小且简单的基因组从头组装, 通过全对全序列映射实现超快构建初步遗传图谱, 尽管初始组装准确性可能受限, 但后续可通过诸如 Nanopolish^[96]等工具进行 polishing 和序列质量提升。

此外, 基于机器学习方法, 如 ONT 公司的 Medaka, 能够有效利用神经网络模型分析原始电流信号, 以提高碱基调用精度, 并在此基础上改善组装结果。此类方法通过深度学习算法捕捉复杂的信号模式, 实现了对纳米孔测序错误的智能预测和校正, 从而有助于进一步降低组装误差。

5.4 DNA 和 RNA 的修饰检测

纳米孔有着极强的灵敏性, 在通过正常碱基和修饰碱基时, 两者的电流会有细微区别, 因此, 通过机器学习的方法可以区分两者电流, 直接检测一些 DNA 和 RNA 的修饰。近年来, 已经开发了一些 DNA 和 RNA 的修饰检测工具。Nanorow(集成于 ONT 的 Tombo)能够从 ONT 的测序数据中检测 3 种 DNA 修饰: 5mC、6mA 和 N⁴-甲基胞嘧啶(N⁴-methylcytosine, 4mC)^[97]。随后, 出现了许多 DNA 修饰检测工具: Nanopolish(5mC)^[98]、signalAlign(5mC,

5hmC 和 6mA)^[99]、DeepMod(5mC 和 6mA)^[100]、DeepSignal(5mC 和 6mA)^[101]和 NanoMod(5mC 和 6mA)^[102]。如图 5(b)~(c)所示,在两次基准测试中^[103-104],Nanopolish、Megalodon 及 DeepSignal 3 款工具均展现出在单分子层面实现高精度单核苷酸分辨率 5mC 检测的能力。此外,ONT 的 Dorado 也有修饰检测的功能,在性能上远高于 Guppy^[96]。与另一单分子测序技术 Pacbio SMRT 相比,ONT 在检测 5mC 时有优势,在检测 6mA 时,精度较低^[80,98,105]。

目前,ONT 的直接 RNA 修饰测序已经取得了重大进展,能够获得可靠的测序数据。检测 RNA 修饰的工具也随之出现,主要分为两种方法:(1)基于错误分布统计,如 EpiNano^[106]预测 N⁶-甲基腺苷(N⁶-methyladenosine, m⁶A),ELIGOS^[107]预测 m⁶A 和 5-甲氧基尿苷(5-Methoxy-UTP, 5moU);(2)基于电流信号,如 Tombo^[97]预测 m⁶A 和 5-甲基胞苷(5-Methyl-CTP, m⁵C),MINES 预测 m⁶A。但是,在单分子水平上,单核苷酸分辨率的 RNA 修饰检测尚未得到证实。

5.5 DNA 存储

机器学习技术可以应用于 DNA 存储的多个方面,例如,加快 DNA 存储编码速度,高效准确地随机访问,从而促进大规模 DNA 存储系统的商业化进程。最直观地,Basecalling 最开始就是利用机器学习的技术,如隐马尔可夫模型、RNN、CNN,结合计算机硬件速度的发展,可以提高碱基识别的准确率和效率。此外,机器学习技术还可以扩展纳米孔可检测的分子类型,提高 DNA 存储的信息密度。

机器学习不仅能够助力提升 DNA 存储读取过程中的 Basecalling 效率,而且还能够在一定程度上绕过 Basecalling,直接对原始信号进行深度分析。2020 年,Doroschak 等^[109]提出了 Porcupine,一种终端用户分子标记系统,仅需无核酸酶水和一台 MinION 设备就可以在几秒内

检测 DNA 分子标签。Porcupine 分子标签由分子比特(molbits)组成,共有 96 种,因此可以编码 96 bits 的信息。为了迅速检测标签,Porcupine 直接使用纳米孔测序的原始电流信号来区分分子比特,无须 Basecalling,从而加快识别速度。如图 5(d)所示,为了提高分子比特的分类准确率,Porcupine 使用 ONT 的 Scruppie 模拟分子比特的电流信号,通过遗传算法使得 96 种分子比特的电流信号尽可能正交,提高后续分子比特的识别准确率。通过预训练的卷积神经网络模型,Porcupine 成功存储了 32 bits 的分子标签(64 bits 的纠错码,可纠 9 位错),并且最快可以在 7 s 内成功识别标签,实现了实时性。Porcupine 用它的可容错性和实时性证明了 DNA 存储和纳米孔检测技术结合后的潜在商用价值。

机器学习同样可以应用于计算模拟之中,以提升 DNA 存储系统的存储容量。通常,在 DNA 存储体系结构中,会利用 DNA 引物(DNA primer),并通过 DNA 杂交原理实现对储存信息的随机访问。因此,在大规模的 DNA 存储系统中,需要相当多的 DNA 引物。然而,为了降低错误访问的概率,必须减少引物序列之间的相似性,这会导致引物设计十分复杂且困难,从而严重阻碍了大规模 DNA 存储的发展。因此,精确控制和预测 DNA 杂交过程对设计大规模 DNA 存储系统至关重要。2021 年,Buterez^[110]提出了一种应用于预测 DNA 杂交任务的机器学习方法。他使用 nupack 生成了一个 250 万大小的杂交数据集,评估了多个机器学习模型(RNN、CNN、CNN lite、RoBERTa),发现机器学习模型不仅能够更准确地预测 DNA 杂交,还能减少运行时间。

6 总 结

作为一项极具前瞻性和革命性的信息存储方案, DNA 存储技术利用生物分子的特性,尤

其是 DNA 分子的巨大存储容量和潜在的超长寿命, 为解决未来数据爆炸式增长带来的存储问题提供了可能。然而, 在从实验室研究走向实际应用的过程中, 仍存在诸多技术挑战。

首先, 传统的测序方法, 如 Sanger 测序^[71]或二代测序技术^[72], 在速度、成本及操作便携性上无法满足 DNA 存储的实际应用需求。纳米孔检测技术在此背景下应运而生, 并展现出显著优势。国内外商业巨头都积极研发并推出了具有竞争力的纳米孔测序产品。例如: 在国外, ONT 公司推出的 MinION 便携式测序仪, 通过与 Apple 公司搭载 M3 芯片的新一代 iMac 或 MacBook 等设备连接, 能够实现现场实时的数据读取和分析, 推动“任何人在任何地方测序”^[111]; 在国内, 齐碳科技已推出首款国产纳米孔基因测序仪 QNome-9604。这些都预示着 DNA 存储领域在全球范围内正迎来快速发展。

尽管当前的纳米孔测序技术在推进 DNA 存储领域已取得显著进展, 但仍面临一些性能瓶颈问题。对于满足大规模 DNA 存储数据的高效处理需求而言, 高通量测序能力尤为重要。在这方面, 尽管 MinION 设备凭借其原理演示和卓越的便携性脱颖而出, 但 ONT 推出的诸如 GridION 和 PromethION 2/24/28 等高通量测序系统更是吸引了广泛的关注。随着对大量数据实时读取要求的增长, 更高吞吐量的测序仪器将在 DNA 存储解决方案中扮演更为关键的角色。此外, R9.4 版本纳米孔的测序准确率约为 85%, 测序速度可达 450 base/s。然而, R9.4 对长均聚物的测序能力仍有待提高, 尤其是在信号识别上的局限性^[76,84]。随着 R10 系列纳米孔的推出, 特别是 R10.4 版本及其引入的 duplex sequencing 技术, 有望提升对长片段和复杂序列测序的准确性, 但这还需要进一步验证和完善^[112-113]。

伴随着人工智能技术的快速发展, 机器学习算法的应用正逐步颠覆纳米孔检测技术和 DNA

存储的传统框架。通过深度学习和模式识别技术, 不仅可以优化生物分子检测的多样性和复杂性, 而且能够开发出更为精确的 Basecalling 算法, 提升 DNA 存储编解码效率。更重要的是, 机器学习有望通过精准解析蛋白质三维结构(如 AlphaFold^[114]所展现的能力), 以及优化生物纳米孔的设计机制, 从根本上改良纳米孔检测技术的性能指标, 进而有力推动 DNA 存储技术向更高层次演进与发展。

DNA 存储技术尚未普及的一个主要原因是缺乏直观易用的用户接口和工具链, 要求操作者具备一定的生物实验技能和数据分析能力。为了优化 DNA 存储技术, 在图 6 中, 本综述提出一种用户友好的 DNA 存储软件平台, 用户通过该平台与高度集成化的 DNA 存储设备进行交互: (1) 图形化界面为用户提供可视化接口; (2) DNA 文件系统进行 DNA 数据的组织与管理; (3) 数据压缩与自动化编码系统通过高效的数据编码算法, 将数据转为 DNA 存储的格式; (4) 智能分析与纠错系统集成先进的错误检测和校正机制, 提高数据恢复的准确度, 并提供修复建议; (5) 安全与隐私保护系统确保 DNA 存储的数据安全性; (6) 技术支持解答用户在使用过程中遇到的问题。

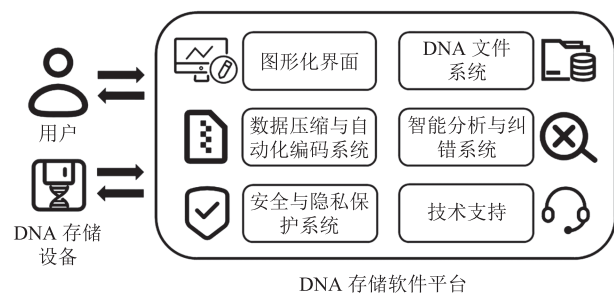


图 6 用户友好的 DNA 存储软件平台

Fig. 6 A user-friendly DNA storage software platform

展望未来十年, 随着纳米孔测序技术的不断革新, 机器学习在该领域的深入渗透, 以及更先进的用户友好型的 DNA 存储系统构建, DNA 存

储技术有可能实现重大突破, 真正迈向商业化应用阶段。相信在不久的将来, 一个崭新的 DNA 存储时代将会拉开帷幕。

参 考 文 献

- [1] Reinsel D, Gantz J, Rydning J. The digitization of the world: from edge to core [EB/OL]. (2020-05-21) [2024-02-25]. <https://www.seagate.com/files/www-content/our-story/trends/files/dataage-idc-report-final.pdf>.
- [2] Williams ED, Ayres RU, Heller M. The 1.7 kilogram microchip: energy and material use in the production of semiconductor devices [J]. *Environmental Science & Technology*, 2002, 36(24): 5504-5510.
- [3] Church GM, Gao Y, Kosuri S. Next-generation digital information storage in DNA [J]. *Science*, 2012, 337(6102): 1628.
- [4] Allentoft ME, Collins M, Harker D, et al. The half-life of DNA in bone: measuring decay kinetics in 158 dated fossils [J]. *Proceedings of the Royal Society B: Biological Sciences*, 2012, 279(1748): 4724-4733.
- [5] Bhat WA. Bridging data-capacity gap in big data storage [J]. *Future Generation Computer Systems*, 2018, 87: 538-548.
- [6] Jaeger L, Chworos A. The architectonics of programmable RNA and DNA nanostructures [J]. *Current Opinion in Structural Biology*, 2006, 16(4): 531-543.
- [7] Hao YY, Li Q, Fan CH, et al. Data storage based on DNA [J]. *Small Structures*, 2021, 2(2): 2000046.
- [8] Merrifield B. Solid phase synthesis [J]. *Science*, 1986, 232(4748): 341-347.
- [9] Belitsky JM, Nguyen DH, Wurtz NR, et al. Solid-phase synthesis of DNA binding polyamides on oxime resin [J]. *Bioorganic & Medicinal Chemistry*, 2002, 10(8): 2767-2774.
- [10] Reuter JA, Spacek DV, Snyder MP. High-throughput sequencing technologies [J]. *Molecular Cell*, 2015, 58(4): 586-597.
- [11] Kasianowicz JJ, Brandin E, Branton D, et al. Characterization of individual polynucleotide molecules using a membrane channel [J]. *Proceedings of the National Academy of Sciences of the United States of America*, 1996, 93(24): 13770-13773.
- [12] Kasianowicz J, Walker B, Krishnasastri M, et al. Genetically engineered pores as metal ion biosensors [J]. *MRS Online Proceedings Library (OPL)*, 1993, 330: 217.
- [13] Bayley H, Cremer PS. Stochastic sensors inspired by biology [J]. *Nature*, 2001, 413(6852): 226-230.
- [14] Clarke J, Wu HC, Jayasinghe L, et al. Continuous base identification for single-molecule nanopore DNA sequencing [J]. *Nature Nanotechnology*, 2009, 4(4): 265-270.
- [15] Cherf GM, Lieberman KR, Rashid H, et al. Automated forward and reverse ratcheting of DNA in a nanopore at 5-Å precision [J]. *Nature Biotechnology*, 2012, 30(4): 344-348.
- [16] Manrao EA, Derrington IM, Laszlo AH, et al. Reading DNA at single-nucleotide resolution with a mutant MspA nanopore and phi29 DNA polymerase [J]. *Nature biotechnology*, 2012, 30(4): 349-353.
- [17] Cao C, Ying YL, Hu ZL, et al. Discrimination of oligonucleotides of different lengths with a wild-type aerolysin nanopore [J]. *Nature Nanotechnology*, 2016, 11(8): 713-718.
- [18] Lieberman KR, Cherf GM, Doody MJ, et al. Processive replication of single DNA molecules in a nanopore catalyzed by phi29 DNA polymerase [J]. *Journal of the American Chemical Society*, 2010, 132(50): 17961-17972.
- [19] Wanunu M, Dadosh T, Ray V, et al. Rapid electronic detection of probe-specific microRNAs using thin nanopore sensors [J]. *Nature Nanotechnology*, 2010, 5(11): 807-814.

- [20] Movileanu L, Schmittschmitt JP, Scholtz JM, et al. Interactions of peptides with a protein pore [J]. *Biophysical Journal*, 2005, 89(2): 1030-1045.
- [21] Yusko EC, Bruhn BR, Eggenberger OM, et al. Real-time shape approximation and fingerprinting of single proteins using a nanopore [J]. *Nature Nanotechnology*, 2017, 12(4): 360-367.
- [22] Bell NAW, Keyser UF. Specific protein detection using designed DNA carriers and nanopores [J]. *Journal of the American Chemical Society*, 2015, 137(5): 2035-2041.
- [23] Squires A, Atas E, Meller A. Nanopore sensing of individual transcription factors bound to DNA [J]. *Scientific Reports*, 2015, 5: 11643.
- [24] Jain M, Olsen HE, Paten B, et al. The Oxford Nanopore MinION: delivery of nanopore sequencing to the genomics community [J]. *Genome Biology*, 2016, 17: 239.
- [25] Neher E, Sakmann B. Single-channel currents recorded from membrane of denervated frog muscle fibres [J]. *Nature*, 1976, 260(5554): 799-802.
- [26] Ying YL, Cao C, Hu YX, et al. A single biomolecule interface for advancing the sensitivity, selectivity and accuracy of sensors [J]. *National Science Review*, 2018, 5(4): 450-452.
- [27] Arroyo JO, Kukura P. Non-fluorescent schemes for single-molecule detection, imaging and spectroscopy [J]. *Nature Photonics*, 2016, 10(1): 11-17.
- [28] Wang YH, Zhao Y, Bollas A, et al. Nanopore sequencing technology, bioinformatics and applications [J]. *Nature Biotechnology*, 2021, 39(11): 1348-1365.
- [29] Oxford Nanopore Technologies. Flow cells [EB/OL]. [2024-02-25]. <https://nanoporetech.com/how-it-works/flow-cells-and-nanopores>.
- [30] Li JL, Stein D, McMullan C, et al. Ion-beam sculpting at nanometre length scales [J]. *Nature*, 2001, 412(6843): 166-169.
- [31] Storm AJ, Chen JH, Ling XS, et al. Fabrication of solid-state nanopores with single-nanometre precision [J]. *Nature Materials*, 2003, 2(8): 537-540.
- [32] Gilboa T, Zreben A, Girsault A, et al. Optically-monitored nanopore fabrication using a focused laser beam [J]. *Scientific Reports*, 2018, 8(1): 9765.
- [33] Yamazaki H, Hu R, Zhao Q, et al. Photothermally assisted thinning of silicon nitride membranes for ultrathin asymmetric nanopores [J]. *ACS Nano*, 2018, 12(12): 12472-12481.
- [34] Meller A. Dynamics of polynucleotide transport through nanometre-scale pores [J]. *Journal of Physics: Condensed Matter*, 2003, 15(17): R581.
- [35] Yan SH, Li XT, Zhang PK, et al. Direct sequencing of 2'-deoxy-2'-fluoroarabinonucleic acid (FANA) using nanopore-induced phase-shift sequencing (NIPSS) [J]. *Chemical Science*, 2019, 10(10): 3110-3117.
- [36] Jia WD, Hu CZ, Wang YQ, et al. Programmable nano-reactors for stochastic sensing [J]. *Nature Communications*, 2021, 12(1): 5811.
- [37] Song LZ, Hobaugh MR, Shustak C, et al. Structure of staphylococcal α -hemolysin, a heptameric transmembrane pore [J]. *Science*, 1996, 274(5294): 1859-1865.
- [38] Niederweis M. Mycobacterial porins—new channel proteins in unique outer membranes [J]. *Molecular Microbiology*, 2003, 49(5): 1167-1177.
- [39] Parker MW, Buckley JT, Postma JPM, et al. Structure of the *Aeromonas* toxin proaerolysin in its water-soluble and membrane-channel states [J]. *Nature*, 1994, 367(6460): 292-295.
- [40] Valpuesta JM, Carrascosa JL. Structure of viral connectors and their function in bacteriophage assembly and DNA packaging [J]. *Quarterly Reviews of Biophysics*, 1994, 27(2): 107-155.
- [41] Mueller M, Grauschopf U, Maier T, et al. The structure of a cytolytic α -helical toxin pore reveals its assembly mechanism [J]. *Nature*, 2009, 459(7247): 726-730.

- [42] Subbarao GV, van den Berg B. Crystal structure of the monomeric porin OmpG [J]. *Journal of Molecular Biology*, 2006, 360(4): 750-759.
- [43] Mindell JA, Zhan H, Huynh PD, et al. Reaction of diphtheria toxin channels with sulfhydryl-specific reagents: observation of chemical reactions at the single molecule level [J]. *Proceedings of the National Academy of Sciences of the United States of America*, 1994, 91(12): 5272-5276.
- [44] Boersma AJ, Bayley H. Continuous stochastic detection of amino acid enantiomers with a protein nanopore [J]. *Angewandte Chemie*, 2012, 124(38): 9744-9747.
- [45] Li MY, Ying YL, Li S, et al. Unveiling the heterogenous dephosphorylation of DNA using an aerolysin nanopore [J]. *ACS Nano*, 2020, 14(10): 12571-12578.
- [46] Deamer D, Akesson M, Branton D. Three decades of nanopore sequencing [J]. *Nature Biotechnology*, 2016, 34(5): 518-524.
- [47] Movileanu L, Howorka S, Braha O, et al. Detecting protein analytes that modulate transmembrane movement of a polymer chain within a single protein pore [J]. *Nature Biotechnology*, 2000, 18(10): 1091-1095.
- [48] Thakur AK, Movileanu L. Real-time measurement of protein-protein interactions at single-molecule resolution using a biological nanopore [J]. *Nature Biotechnology*, 2019, 37(1): 96-101.
- [49] Wei RS, Tampé R, Rant U. Stochastic sensing of proteins with receptor-modified solid-state nanopores [J]. *Biophysical Journal*, 2012, 102(3): 429a.
- [50] Fahie MA, Yang B, Mullis M, et al. Selective detection of protein homologues in serum using an OmpG nanopore [J]. *Analytical Chemistry*, 2015, 87(21): 11143-11149.
- [51] Bell NAW, Keyser UF. Digitally encoded DNA nanostructures for multiplexed, single-molecule protein sensing with nanopores [J]. *Nature Nanotechnology*, 2016, 11(7): 645-651.
- [52] Yu LN, Kang XQ, Li FJ, et al. Unidirectional single-file transport of full-length proteins through a nanopore [J]. *Nature Biotechnology*, 2023, 41(8): 1130-1139.
- [53] Sauciuc A, Morozzo della Rocca B, Tadema MJ, et al. Translocation of linearized full-length proteins through an engineered nanopore under opposing electrophoretic force [J/OL]. *Nature Biotechnology*. (2023-09-18)[2024-02-25]. <https://www.nature.com/articles/s41587-023-01954-x>.
- [54] Ouldali H, Sarthak K, Ensslen T, et al. Electrical recognition of the twenty proteinogenic amino acids using an aerolysin nanopore [J]. *Nature Biotechnology*, 2020, 38(2): 176-181.
- [55] Wang KF, Zhang SY, Zhou X, et al. Unambiguous discrimination of all 20 proteinogenic amino acids and their modifications by nanopore [J]. *Nature Methods*, 2023, 21(1): 92-101.
- [56] Zhang Y, Yi YK, Li ZY, et al. Peptide sequencing based on host-guest interaction-assisted nanopore sensing [J]. *Nature Methods*, 2024, 21(1): 102-109.
- [57] Wang Y, Zheng DL, Tan QL, et al. Nanopore-based detection of circulating microRNAs in lung cancer patients [J]. *Nature Nanotechnology*, 2011, 6(10): 668-674.
- [58] Rozevsky Y, Gilboa T, van Kooten XF, et al. Quantification of mRNA expression using single-molecule nanopore sensing [J]. *ACS Nano*, 2020, 14(10): 13964-13974.
- [59] Guo JM, Amini S, Lei Q, et al. Robust and long-term cellular protein and enzymatic activity preservation in biomineralized mammalian cells [J]. *ACS Nano*, 2022, 16(2): 2164-2175.
- [60] Baoutina A, Bhat S, Partis L, et al. Storage stability of solutions of DNA standards [J]. *Analytical Chemistry*, 2019, 91(19): 12268-12274.
- [61] Kohll AX, Antkowiak PL, Chen WD, et al. Stabilizing synthetic DNA for long-term data storage with earth alkaline salts [J]. *Chemical*

- Communications, 2020, 56(25): 3613-3616.
- [62] Organick L, Ang SD, Chen YJ, et al. Random access in large-scale DNA data storage [J]. *Nature Biotechnology*, 2018, 36(3): 242-248.
- [63] Lopez R, Chen YJ, Ang SD, et al. DNA assembly for nanopore data storage readout [J]. *Nature Communications*, 2019, 10: 2933.
- [64] Chen WG, Han MZ, Zhou JT, et al. An artificial chromosome for data storage [J]. *National Science Review*, 2021, 8(5): nwab028.
- [65] Sun FJ, Dong YM, Ni M, et al. Mobile and self-sustained data storage in an extremophile genomic DNA [J]. *Advanced Science*, 2023, 10(10): 2206201.
- [66] Liu HJ, Wang JB, Song SP, et al. A DNA-based system for selecting and displaying the combined result of two input variables [J]. *Nature Communications*, 2015, 6: 10089.
- [67] Ge ZL, Gu HZ, Li Q, et al. Concept and development of framework nucleic acids [J]. *Journal of the American Chemical Society*, 2018, 140(51): 17808-17819.
- [68] Chen KK, Kong JL, Zhu JB, et al. Digital data storage using DNA nanostructures and solid-state nanopores [J]. *Nano Letters*, 2018, 19(2): 1210-1215.
- [69] Cao C, Krapp LF, Al Ouahabi A, et al. Aerolysin nanopores decode digital information stored in tailored macromolecular analytes [J]. *Science Advances*, 2020, 6(50): eabc2661.
- [70] Tabatabaei SK, Pham B, Pan C, et al. Expanding the molecular alphabet of DNA-based data storage systems with neural network nanopore readout processing [J]. *Nano Letters*, 2022, 22(5): 1905-1914.
- [71] Sanger F. Determination of nucleotide sequences in DNA [J]. *Science*, 1981, 214(4526): 1205-1210.
- [72] Bentley DR, Balasubramanian S, Swerdlow HP, et al. Accurate whole human genome sequencing using reversible terminator chemistry [J]. *Nature*, 2008, 456(7218): 53-59.
- [73] Eid J, Fehr A, Gray J, et al. Real-time DNA sequencing from single polymerase molecules [J]. *Science*, 2009, 323(5910): 133-138.
- [74] Greener JG, Kandathil SM, Moffat L, et al. A guide to machine learning for biologists [J]. *Nature Reviews Molecular Cell Biology*, 2022, 23(1): 40-55.
- [75] LeCun Y, Bengio Y, Hinton G. Deep learning [J]. *Nature*, 2015, 521(7553): 436-444.
- [76] Wick RR, Judd LM, Holt KE. Performance of neural network basecalling tools for Oxford Nanopore sequencing [J]. *Genome Biology*, 2019, 20: 129.
- [77] Rang FJ, Kloosterman WP, de Ridder J. From squiggle to basepair: computational approaches for improving nanopore sequencing read accuracy [J]. *Genome Biology*, 2018, 19(1): 90.
- [78] Magi A, Semeraro R, Mingrino A, et al. Nanopore sequencing data analysis: state of the art, applications and challenges [J]. *Briefings in Bioinformatics*, 2018, 19(6): 1256-1272.
- [79] Senol Cali D, Kim JS, Ghose S, et al. Nanopore sequencing technology and tools for genome assembly: computational analysis of the current state, bottlenecks and future directions [J]. *Briefings in Bioinformatics*, 2019, 20(4): 1542-1559.
- [80] Amarasinghe SL, Su SA, Dong XY, et al. Opportunities and challenges in long-read sequencing data analysis [J]. *Genome Biology*, 2020, 21: 30.
- [81] David M, Dursi LJ, Yao D, et al. Nanocall: an open source basecaller for Oxford Nanopore sequencing data [J]. *Bioinformatics*, 2017, 33(1): 49-55.
- [82] Boža V, Brejová B, Vinař T. DeepNano: deep recurrent neural networks for basecalling in MinION nanopore reads [J]. *PLoS One*, 2017, 12(6): e0178751.
- [83] Teng HT, Cao MD, Hall MB, et al. Chiron: translating nanopore raw signal directly into nucleotide sequence using deep learning [J].

- GigaScience, 2018, 7(5): giy037.
- [84] Zeng JW, Cai HM, Peng H, et al. Causalcall: nanopore basecalling using a temporal convolutional network [J]. *Frontiers in Genetics*, 2020, 10: 1332.
- [85] Dittforth S, Ozturk D, Mueller M. Benchmarking the Oxford Nanopore Technologies basecallers on AWS [EB/OL]. (2023-05-18)[2024-02-25]. <https://aws.amazon.com/blogs/hpc/benchmarking-the-oxford-nanopore-technologies-basecallers-on-aws/>.
- [86] Fu SH, Wang AQ, Au KF. A comparative evaluation of hybrid error correction methods for error-prone long reads [J]. *Genome Biology*, 2019, 20(1): 26.
- [87] Lima L, Marchet C, Caboche S, et al. Comparative assessment of long-read error correction software applied to Nanopore RNA-sequencing data [J]. *Briefings in Bioinformatics*, 2020, 21(4): 1164-1181.
- [88] Koren S, Walenz BP, Berlin K, et al. Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation [J]. *Genome Research*, 2017, 27(5): 722-736.
- [89] Salmela L, Walve R, Rivals E, et al. Accurate self-correction of errors in long reads using de Bruijn graphs [J]. *Bioinformatics*, 2017, 33(6): 799-806.
- [90] Au KF, Underwood JG, Lee L, et al. Improving PacBio long read accuracy by short read alignment [J]. *PLoS One*, 2012, 7(10): e46679.
- [91] Goodwin S, Gurtowski J, Ethe-Sayers S, et al. Oxford Nanopore sequencing, hybrid error correction, and *de novo* assembly of a eukaryotic genome [J]. *Genome Research*, 2015, 25(11): 1750-1756.
- [92] Salmela L, Rivals E. LoRDEC: accurate and efficient long read error correction [J]. *Bioinformatics*, 2014, 30(24): 3506-3514.
- [93] Bao E, Lan LX. HALC: High throughput algorithm for long read error correction [J]. *BMC Bioinformatics*, 2017, 18: 204.
- [94] Wang LT, Qu L, Yang LS, et al. NanoReviser: an error-correction tool for nanopore sequencing based on a deep learning algorithm [J]. *Frontiers in Genetics*, 2020, 11: 900.
- [95] Li H. Minimap and minimap: fast mapping and *de novo* assembly for noisy long sequences [J]. *Bioinformatics*, 2016, 32(14): 2103-2110.
- [96] Loman NJ, Quick J, Simpson JT. A complete bacterial genome assembled *de novo* using only nanopore sequencing data [J]. *Nature Methods*, 2015, 12(8): 733-735.
- [97] Stoiber M, Quick J, Egan R, et al. *De novo* identification of DNA modifications enabled by genome-guided nanopore signal processing [Z/OL]. *BioRxiv Preprint*, BioRxiv: 094672, 2016.
- [98] Simpson JT, Workman RE, Zuzarte PC, et al. Detecting DNA cytosine methylation using nanopore sequencing [J]. *Nature Methods*, 2017, 14(4): 407-410.
- [99] Rand AC, Jain M, Eizenga JM, et al. Mapping DNA methylation with high-throughput nanopore sequencing [J]. *Nature Methods*, 2017, 14(4): 411-413.
- [100] Liu Q, Fang L, Yu GL, et al. Detection of DNA base modifications by deep recurrent neural network on Oxford Nanopore sequencing data [J]. *Nature Communications*, 2019, 10(1): 2449.
- [101] Ni P, Huang N, Zhang Z, et al. DeepSignal: detecting DNA methylation state from Nanopore sequencing reads using deep-learning [J]. *Bioinformatics*, 2019, 35(22): 4586-4595.
- [102] Liu Q, Georgieva DC, Egli D, et al. NanoMod: a computational tool to detect DNA modifications using Nanopore long-read sequencing data [J]. *BMC Genomics*, 2019, 20: 78.
- [103] Yuen ZWS, Srivastava A, Daniel R, et al. Systematic benchmarking of tools for CpG methylation detection from nanopore sequencing [J]. *Nature Communications*, 2021, 12: 3438.
- [104] Liu Y, Rosikiewicz W, Pan ZW, et al. DNA methylation-calling tools for Oxford Nanopore

- sequencing: a survey and human epigenome-wide evaluation [J]. *Genome Biology*, 2021, 22(1): 295.
- [105] Fang G, Munera D, Friedman DI, et al. Genome-wide mapping of methylated adenine residues in pathogenic *Escherichia coli* using single-molecule real-time sequencing [J]. *Nature Biotechnology*, 2012, 30(12): 1232-1239.
- [106] Liu HL, Begik O, Lucas MC, et al. Accurate detection of m⁶A RNA modifications in native RNA sequences [J]. *Nature Communications*, 2019, 10(1): 4079.
- [107] Jenjaroenpun P, Wongsurawat T, Wadley TD, et al. Decoding the epitranscriptional landscape from native RNA sequences [J]. *Nucleic Acids Research*, 2021, 49(2): 620.
- [108] Lorenz DA, Sathe S, Einstein JM, et al. Direct RNA sequencing enables m⁶A detection in endogenous transcript isoforms at base-specific resolution [J]. *RNA*, 2020, 26(1): 19-28.
- [109] Doroschak K, Zhang KR, Queen M, et al. Rapid and robust assembly and decoding of molecular tags with DNA-based nanopore signatures [J]. *Nature Communications*, 2020, 11(1): 5454.
- [110] Buterez D. Scaling up DNA digital data storage by efficiently predicting DNA hybridisation using deep learning [J]. *Scientific Reports*, 2021, 11(1): 20517.
- [111] Sanghera G. Oxford Nanopore meets Apple's M3 silicon chip, hailing a new era of distributed genome sequencing [EB/OL]. (2023-10-31)[2024-02-25]. <https://nanoporetech.com/about-us/news/blog-oxford-nanopore-meets-apples-m3-silicon-chip-hailing-new-era-distributed-genome>.
- [112] Tytgat O, Gansemans Y, Weymaere J, et al. Nanopore sequencing of a forensic STR multiplex reveals loci suitable for single-contributor STR profiling [J]. *Genes*, 2020, 11(4): 381.
- [113] Huang YT, Liu PY, Shih PW. Homopolish: a method for the removal of systematic errors in nanopore sequencing by homologous polishing [J]. *Genome Biology*, 2021, 22(1): 95.
- [114] Jumper J, Evans R, Pritzel A, et al. Highly accurate protein structure prediction with AlphaFold [J]. *Nature*, 2021, 596(7873): 583-589.