

Review

Not peer-reviewed version

Understanding Pathways in Bioinformatics, Genomics, and Health Applications

[Diptarup Mallick](#)*

Posted Date: 19 January 2026

doi: 10.20944/preprints202601.1343.v1

Keywords: bioinformatics; genomics; deep learning; state space models (SSM); precision medicine; reproducibility; comparative genomics; tokenization; large language models (LLM); long-range dependencies; telomere-to-telomere (T2T)



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a [Creative Commons CC BY 4.0 license](#), which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Review

Understanding Pathways in Bioinformatics, Genomics, and Health Applications

Diptarup Mallick

Post Graduate Department of Zoology, Barasat Government College, India; diptarupmallick3@gmail.com

Abstract

The rapid expansion of high-throughput sequencing technologies and the completion of telomere-to-telomere (T2T) assemblies have transformed genomics into a data-driven discipline, shifting the research focus from data generation to large-scale computational discovery. This literature review synthesizes foundational and emergent pathways in bioinformatics, genomics, and their integration into health applications. We examine the critical role of genomic reproducibility and benchmarking in establishing clinical trust, alongside mathematical models for comparative genomics, such as the Double-Cut-and-Join (DCJ) distance. A significant portion of this review is dedicated to methodological shifts in representation learning, specifically evaluating the impact of Byte-Pair Encoding (BPE) tokenization on genomic language models and the dominance of repetitive elements in sequence vocabularies. Furthermore, we explore the evolution of deep learning architectures, contrasting traditional convolutional and recurrent neural networks with recent advancements in State Space Models (SSMs). These emergent architectures, such as Caduceus and Mamba, demonstrate linear-time complexity and superior performance in capturing long-range regulatory dependencies across ultra-long genomic sequences. Finally, we discuss how these computational innovations converge to support the goals of precision medicine. By mapping these trajectories, this review provides a comprehensive overview of the technical and theoretical challenges inherent in modeling the complexity of the human genome for clinical and biological insights.

Keywords: bioinformatics; genomics; deep learning; state space models (SSM); precision medicine; reproducibility; comparative genomics; tokenization; large language models (LLM); long-range dependencies; telomere-to-telomere (T2T)

1. Introduction

The exponential growth of biomedical data in the last two decades has fundamentally transformed the landscape of bioinformatics, genomics, and their applications in health sciences [1], [6]. High-throughput sequencing technologies, the advent of telomere-to-telomere (T2T) assemblies [7], and powerful computational methods have collectively enabled unprecedented insights into the structure, function, and variation of genomes [8]. These advances have shifted the paradigm from data generation to data-driven discovery, where the challenge lies in transforming massive, heterogeneous, and complex datasets into actionable knowledge for both basic science and medical applications [9], [10].

Bioinformatics stands at the intersection of biology, computer science, and statistics, leveraging computational approaches to analyze, interpret, and visualize biological data [11]. Genomics, as a central focus within bioinformatics, deals with the comprehensive study of genomes, encompassing both their static sequence properties and dynamic regulatory mechanisms [12]. The integration of these disciplines with health applications, particularly in the era of precision medicine, has opened new avenues for disease diagnosis, prognosis, and therapy [13], [14].

However, this integration is not without challenges. Issues such as reproducibility of genomic analyses [1], scalability of computational models to ultra-long sequences [5], [15], optimal tokenization strategies for genomic language models [3], efficient algorithms for comparative

genomics [2], and the effective application of deep learning architectures [4], [16] all present active and evolving research frontiers. This literature review synthesizes foundational and recent advances along these pathways, drawing on selected key works to map the current landscape and future trajectories in bioinformatics, genomics, and health.

2. Methodology

2.1. Literature Search Strategy

A comprehensive search was performed across major electronic databases, including PubMed, IEEE Xplore, Google Scholar, and arXiv. The inclusion of preprint servers (specifically arXiv and bioRxiv) was prioritized to capture the most recent technological shifts in State Space Models (SSMs) and Telomere-to-Telomere (T2T) analysis, reflecting the current state-of-the-art as of 2024–2025.

The search utilized combinations of the following primary and secondary keywords:

- Primary: Bioinformatics, Genomics, Deep Learning, State Space Models, Precision Medicine.
- Secondary: Genomic reproducibility, Double-Cut-and-Join (DCJ), Byte-Pair Encoding (BPE), Long-range genomic dependencies, T2T assembly, Caduceus, Mamba.

2.2. Selection Criteria

To ensure high scientific rigor and relevance, studies were selected based on the following criteria:

- Temporal Relevance: Preference was given to studies published between 2012 and 2025 to capture the rise of high-throughput sequencing and the deep learning revolution.
- Thematic Alignment: Works were included if they addressed specific methodological challenges such as sequence tokenization, algorithmic complexity in comparative genomics, or the scalability of neural architectures to ultra-long DNA sequences.
- Impact and Reliability: Selection favored peer-reviewed journals (e.g., Nature, Science, Bioinformatics) and high-impact benchmarking papers from established consortia such as GIAB and MAQC.

2.3. Data Extraction and Thematic Synthesis

The selected literature was categorized into four thematic pillars to provide a logical progression:

- Foundational Robustness: Investigating the "reproducibility crisis" in bioinformatics and the mathematical frameworks for evolutionary distance.
- Representation Learning: Evaluating how DNA is tokenized and processed by Large Language Models (LLMs).
- Architectural Evolution: Contrasting traditional architectures (CNN, RNN, Transformers) with emergent linear-time models (SSMs).
- Clinical Integration: Synthesizing how computational advances translate into precision medicine and diagnostic tools.

2.4. Quality Assessment

Each reference was evaluated for its contribution to the field. Methodological papers were assessed based on their benchmarking rigor (e.g., use of the Genomics Long-Range Benchmark), while theoretical papers were evaluated for their mathematical soundness in modeling genome rearrangements and evolutionary pathways. A total of 40 key references were ultimately selected to form the basis of this review.

3. Conceptual Foundations and Methodological Advances in Bioinformatics and Genomics

3.1. The Role of Reproducibility in Genomic Research

Reproducibility is a cornerstone of scientific inquiry, ensuring that findings are robust, generalizable, and translatable to real-world applications [17]. In genomics, the concept of reproducibility is particularly nuanced due to the multi-stage nature of data generation and analysis processes. As Icer Baykal et al. argue, reproducibility in genomics can be classified along several axes, including methods reproducibility and “genomic reproducibility”—the ability of bioinformatics tools to yield consistent results across technical replicates [1].

Technical replicates arise from multiple sequencing runs or library preparations of the same biological sample. While methods reproducibility assesses the deterministic behavior of computational tools, genomic reproducibility addresses the impact of stochastic elements introduced during sequencing [18]. This distinction is crucial for medical genomics, where clinical decisions may hinge on the robustness of variant calling [19].

Furthermore, Icer Baykal et al. emphasize that bioinformatics tools can both mitigate and introduce unwanted variation [1]. For example, normalization techniques can remove batch effects, but algorithmic biases—such as reference bias in read alignment—may confound analyses [20]. To address these challenges, community-wide benchmarking efforts like the Genome in a Bottle (GIAB) consortium [21] and the adoption of standardized workflow languages such as Nextflow or Snakemake [22] have become critical for evaluating the reliability of genomic analyses.

3.2. Comparative Genomics: Quantifying Evolutionary Distances and Reconstruction

Comparative genomics seeks to understand evolutionary relationships and functional conservation by comparing whole genomes [23]. A central computational concept is the rearrangement distance—the minimal number of large-scale rearrangements (reversals, translocations, fissions, and fusions) required to transform one genome into another. The Double-Cut-and-Join (DCJ) model has emerged as a powerful abstraction for modeling such rearrangements [2], [24].

However, as Aganezov and Alekseyev observe, extending this concept to multiple genomes introduces complexity [2]. The median score is NP-hard to compute, leading researchers to use approximations like the “triangle score” [2]. These findings have practical implications for phylogeny reconstruction and ancestral genome inference [25]. The gap between easily computed pairwise distances and the biologically meaningful median score underscores the need for algorithmic innovation in understanding genome evolution [26].

3.3. Tokenization and Representation Learning in Genomic Language Models

The rise of large language models (LLMs) has inspired analogous approaches in genomics, where DNA sequences serve as the substrate for unsupervised representation learning [27]. A critical consideration is the choice of tokenization strategy. While k-mer tokenization was traditional, Byte-Pair Encoding (BPE) has recently been adopted for genomic data [3], [28].

Popova et al. evaluate BPE tokenization in the context of T2T primate genomes, finding that BPE is highly sensitive to repetitive elements [3]. Their analysis reveals that BPE vocabularies are often more reflective of species-specific repeats than evolutionary relationships [3]. This suggests that while BPE effectively compresses sequences, it may obscure functionally relevant variation, necessitating domain-specific adaptations like repeat masking or hybrid k-mer strategies [29], [30].

3.4. Deep Learning Architectures in Bioinformatics

Deep learning has revolutionized bioinformatics by extracting patterns from high-dimensional biomedical data [16]. As Min et al. review, deep learning architectures are now routinely applied across omics, biomedical imaging, and signal processing [4].

3.5. Architecture Taxonomy and Applications

The taxonomy of deep learning models in bioinformatics comprises deep neural networks (DNNs), convolutional neural networks (CNNs), and recurrent neural networks (RNNs) [4]. CNNs, with their ability to capture spatial locality, have achieved state-of-the-art results in sequence-based motif discovery and biomedical imaging [31]. RNNs and Transformers are adept at modeling sequential dependencies, such as splice junction prediction or protein structure forecasting [32], [33].

The successful deployment of deep learning has been facilitated by open-source libraries like TensorFlow and PyTorch [34]. However, practical challenges remain, including the handling of imbalanced data, interpretability of learned representations (the "black box" problem), and the need for multimodal data integration [4], [35].

4. Advanced Computational Models for Long-Range Genomic Analysis

4.1. Modeling Long-Range Dependencies

A defining characteristic of eukaryotic genomes is long-range regulatory interactions, where enhancers modulate gene expression over hundreds of kilobases [36]. Conventional models struggle with these ultra-long sequences due to quadratic memory constraints. Popov et al. address this by benchmarking State Space Models (SSMs)—architectures with linear complexity—on long-range genomic tasks [5].

Their findings demonstrate that SSMs, specifically the Caduceus and Hawk architectures, match or surpass transformer performance in variant effect prediction while generalizing to input sequences 10-100 times longer than those seen during training [5]. This allows for processing sequences of 1 million base pairs on a single GPU, enabling the modeling of entire genomic regions [37].

4.2. Benchmarking and Evaluation

The Genomics Long-Range Benchmark (GLRB) encompasses tasks such as gene expression prediction and chromatin feature identification [5]. Key findings suggest that SSMs are uniquely capable of zero-shot extrapolation, maintaining stable performance up to 1 million base pairs [5]. These results highlight the importance of architectural innovation in capturing the structural variation and regulatory landscapes of the genome [38].

4.3. Integrating Bioinformatics Pathways into Health Applications

The ultimate goal of these computational pathways is the realization of precision medicine [39]. The ability to reproducibly and accurately analyze genomic data is essential for clinical translation [1]. Furthermore, the integration of multi-omics data—combining genomics, transcriptomics, and proteomics—represents the frontier of holistic clinical modeling [40].

5. Conclusions

This literature review has traced the pathways through which bioinformatics, genomics, and computational health applications are evolving. From foundational concepts of reproducibility and evolutionary distance, through methodological advances in representation learning and deep learning architectures, to the integration of these approaches in precision medicine, the field is marked by rapid innovation and complex challenges.

Key themes emerging from the reviewed literature include the critical importance of reproducibility and benchmarking, the need for domain-specific adaptations in computational methods, the scalability of models to ultra-long and heterogeneous data, and the imperative of interpretability in health applications. While significant progress has been made, the full realization of bioinformatics' potential will require continued methodological refinement, interdisciplinary collaboration, and a commitment to open science and data sharing.

As the volume, variety, and velocity of biomedical data continue to grow, the pathways outlined in this review will serve as both a roadmap and a call to action for researchers, clinicians, and policymakers working at the intersection of computation, genomics, and health.

Funding: Not applicable.

Data Availability: Not applicable.

Ethical Statement: None.

Conflict of Interest: Not applicable.

References

1. Icer Baykal, P. et al. Genomic reproducibility in the bioinformatics era. arXiv 2023, arXiv:2308.09558.
2. Aganezov, S.; Alekseyev, M.A. On pairwise distances and median score of three genomes under DCJ. BMC Bioinformatics 2012, 13, S1.
3. Popova, M. et al. When repeats drive the vocabulary: a Byte-Pair Encoding analysis of T2T primate genomes. arXiv 2025, arXiv:2505.08918.
4. Min, S. et al. Deep learning in bioinformatics. Brief. Bioinform. 2017, 18, 851–869.
5. Popov, M. et al. Leveraging State Space Models in Long Range Genomics. arXiv 2025, arXiv:2504.06304.
6. Stephens, Z.D. et al. Big Data: Astronomical or Genomical? PLoS Biol. 2015, 13, e1002195.
7. Nurk, S. et al. The complete sequence of a human genome. Science 2022, 376, 44–53.
8. Shendure, J. et al. DNA sequencing at 40: past, present and future. Nature 2017, 550, 345–353.
9. Lander, E.S. Initial impact of the sequencing of the human genome. Nature 2011, 470, 187–197.
10. Baker, M. 1,500 scientists lift the lid on reproducibility. Nature 2016, 533, 452–454.
11. Luscombe, N.M. et al. What is bioinformatics? A proposed definition and overview of the field. Methods Inf. Med. 2001, 40, 346–358.
12. Birney, E. The impact of genomics on 21st century medicine. Cold Spring Harb. Mol. Case Stud. 2019, 5, a004317.
13. Collins, F.S.; Varmus, H. A new initiative on precision medicine. N. Engl. J. Med. 2015, 372, 793–795.
14. Topol, E.J. High-performance medicine: the convergence of human and artificial intelligence. Nat. Med. 2019, 25, 44–56.
15. Gu, A.; Dao, T. Mamba: Linear-time sequence modeling with selective state spaces. arXiv 2023, arXiv:2312.00752.
16. LeCun, Y. et al. Deep learning. Nature 2015, 521, 436–444.
17. Sandve, G.K. et al. Ten simple rules for reproducible computational research. PLoS Comput. Biol. 2013, 9, e1003285.
18. Mangul, S. et al. Systematic visualization of the reproducibility of published genomic surveys. Nat. Methods 2019, 16, 11–12.
19. Poplin, R. et al. A universal SNP and small-indel variant caller using deep neural networks. Nat. Biotechnol. 2018, 36, 983–987.
20. Li, H. Toward better understanding of artifacts in variant calling from high-coverage samples. Bioinformatics 2014, 30, 2843–2851.
21. Zook, J.M. et al. Integrating human sequence data sets provides a benchmark of West African ancestry. Sci. Data 2019, 6, 287.
22. Di Tommaso, P. et al. Nextflow enables reproducible computational workflows. Nat. Biotechnol. 2017, 35, 316–319.

23. Altschul, S. et al. Basic local alignment search tool. *J. Mol. Biol.* 1990, 215, 403–410.
24. Yancopoulos, S. et al. Efficient distance calculation and finite chromosome phylogeny under unrestricted genome rearrangements. *Bioinformatics* 2005, 21, 3340–3346.
25. Bourque, G. Comparative genomics and genome evolution. *Curr. Opin. Genet. Dev.* 2009, 19, 507–512.
26. Fertin, G. et al. *Combinatorics of Genome Rearrangements*. 1st ed.; MIT Press: Cambridge, MA, USA, 2009; pp. 1–330.
27. Ji, Y. et al. DNABERT: pre-trained Bidirectional Encoder Representations from Transformers model for DNA-language in genome. *Bioinformatics* 2021, 37, 2112–2120.
28. Dalla-Torre, H. et al. The Nucleotide Transformer: Building and evaluating robust foundation models for human genomics. *bioRxiv* 2023, 2023.01.11.523679.
29. Nguyen, E. et al. HyenaDNA: Long-range genomic sequence modeling at single nucleotide resolution. In *Advances in Neural Information Processing Systems*; Curran Associates: Red Hook, NY, USA, 2023; Vol. 36.
30. Thomas, C. et al. Evo: DNA foundation modeling from molecular to genome scale. *bioRxiv* 2024, 2024.02.27.582234.
31. Alipanahi, B. et al. Predicting the sequence specificities of DNA- and RNA-binding proteins by deep learning. *Nat. Biotechnol.* 2015, 33, 831–838.
32. Jumper, J. et al. Highly accurate protein structure prediction with AlphaFold. *Nature* 2021, 596, 583–589.
33. Eraslan, G. et al. Deep learning: new computational modelling techniques for genomics. *Nat. Rev. Genet.* 2019, 20, 389–403.
34. Paszke, A. et al. PyTorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*; Curran Associates: Red Hook, NY, USA, 2019; Vol. 32.
35. Wang, Y. et al. The applications of deep learning in multi-omics data integration. *Brief. Bioinform.* 2021, 22, bbab154.
36. Lieberman-Aiden, E. et al. Comprehensive mapping of long-range interactions reveals folding principles of the human genome. *Science* 2009, 326, 289–293.
37. Schiff, P. et al. Caduceus: Bi-directional Equivariant Long-range DNA Sequence Modeling. *arXiv* 2024, arXiv:2403.03230.
38. Zhou, J.; Troyanskaya, O.G. Predicting effects of noncoding variants with deep learning-based sequence model. *Nat. Methods* 2015, 12, 931–934.
39. Hamburg, M.A.; Collins, F.S. The Path to Personalized Medicine. *N. Engl. J. Med.* 2010, 363, 301–304.
40. Hasin, Y. et al. Multi-omics strategies and data integration. *Genome Biol.* 2017, 18, 83

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.