Generative Biology

This manuscript (<u>permalink</u>) was automatically generated from <u>In-Vivo-Group/generative-biology@4ff558b</u> on June 14, 2024.

Authors

- Alexander J. Titus [™]
 - **(D** <u>0000-0002-0145-9564</u> **· (C** <u>alexandertitus</u>

In Vivo Group, Washington, DC, USA; International Computer Science Institute, Berkeley, CA, USA · Funded by Grant TBD

- Matthew E. Walsh

Center for Health Security, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA; Department of Environmental Health and Engineering, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA

Abstract

The rapid pace of progress in generative artificial intelligence (AI) techniques like deep learning, reinforcement learning, and transformer neural networks is transforming the life sciences and biomedicine. This living review paper provides an updatable, comprehensive overview and analysis of the latest literature on generative biology – the application of cutting-edge generative AI methods to accelerate insights and innovation across the life sciences and healthcare. All authors are welcome to contribute to this review via pull requests.

The review synthesizes key developments in using generative models for de novo biomedical discovery, design, and decision support. It examines techniques and applications including deep learning on omics data for personalized medicine, generative chemistry for drug development, protein structure prediction for molecular engineering, image synthesis for pathology, language models for clinical decision support, robotic simulation for prosthetics, and generative networks for cell programming.

The review highlights representative studies and benchmarks in each area while contextualizing progress, limitations, emerging best practices, and directions for future work. It also discusses social and ethical challenges raised by generative biology applications, such as compounding bias, system opacity, and dual-use risks, alongside proposed solutions.

As a living review, this paper will be continually updated as the field rapidly advances to provide researchers and practitioners with an up-to-date reference on the state of the art in employing generative AI to accelerate biomedicine for the collective good.

Executive Summary

The goal with be a 2 page TL;DR of the review after v1 is complete. Need more content.

Introduction

This is the start of the Generative Biology living review!

Computers, Algorithms and the Internet

1950s and 1960s: Early computers and algorithms

Computers were used in the early 1950s for population genetics calculations [1]. Notably, the inception of computational modeling in biology dates to the origins of computer science itself. British mathematician and logician Alan Turing, often referred to as "the father of computing", used primitive computers to implement a model of biological morphogenesis (the emergence of pattern and shape in living organisms) in 1952 [2]. At about the same time, a computer called MANIAC was used for measuring speculative genetic codes; it was originally built for weaponry research at the Los Alamos National Laboratory in New Mexico [3].

Computers were used for the study of protein structure by the 1960s, and other increasingly diverse analyses. These developments marked the rise of the computational biology field, stemming from research focused on protein crystallography, in which scientists found computers indispensable for carrying out laborious Fourier analyses to determine the three-dimensional structure of proteins [4,5].

In addition to advances in determination of protein structures through crystallography, the first sequence of protein, insulin, was published [6,7]. More efficient protein sequencing methods, such as the Edman degradation technique [8], enabled sequencing 15 different proteins over a decade [9]. COMPROTEIN, one of the first bioinformatics softwares developed in the early 1960s, was designed to overcome the limitations of Edman sequencing [10]. In an effort to simplify the handling of protein sequence data for the COMPROTEIN software, a one-letter amino acid code was developed [11]. This one-letter code was first used in the Atlas of Protein Sequence and Structure [12], the first biological sequence database, laying the groundwork for paleogenetic studies.

Development of methods to compare protein sequences followed. The Needleman-Wunsch algorithm [13], the first dynamic programming algorithm developed for pairwise protein sequence alignments, was introduced in the 1970s. Multiple sequence alignment (MSA) algorithms followed in the early 1980s. Progressive sequence alignment was introduced by Feng and Doolittle in 1987 [14]. The MSA software CLUSTAL, a simplification of the Feng-Doolittle algorithm [15] was developed in 1988. It is still used and maintained to this day [16].

1970s: From protein to DNA analysis

The deciphering of all 64 triplet codons of the genetic code in 196817 fueled a desire to efficiently determine the sequence of DNA that existed into the 1970s. This desire led to the development of cost-efficient DNA sequencing methods, such as the Maxam-Gilbert and Sanger sequencing techniques in the mid-1970s [6,7,17]. With this new ability to generate DNA sequence data, a paradigm shift from protein analysis to DNA analysis occurred in the late 1970s. Concurrently,

concerns over recombinant DNA research led to safety protocols established during the 1975 Asilomar conference [18].

New DNA sequencing techniques resulted in significantly more data to be analyzed, a task at which computation could help. The first software dedicated to analyzing Sanger sequencing reads was published in 1979 [19]. DNA sequences began to be utilized in phylogenetic inference with pioneering methods like maximum likelihood for inferring phylogenetic trees from DNA sequences [20]. Several bioinformatics tools and statistical methods were developed following this work. The adoption of Bayesian statistics in molecular phylogeny in the 1990s was inspired by this [21] and is still commonly used in biology today [22]. Yet, numerous computational limitations needed to be overcome during the latter half of the 1970s to expand the utilization of computing in the life sciences, especially in DNA analysis. The subsequent decade proved instrumental in addressing these challenges.

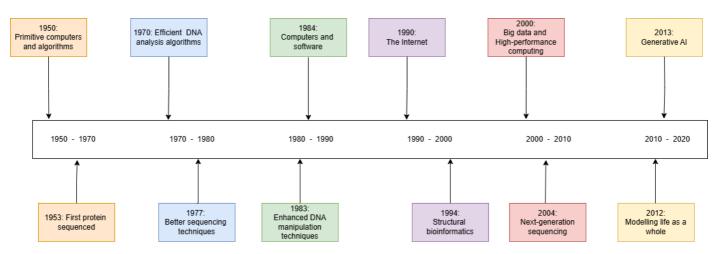


Figure 1: The history of parallel advancements in computing and the life sciences: A timeline of major milestones.

1980s: Simultaneous advances in computing and biology

Parallel advancements in biology and computing propelled bioinformatics forward during the 1980s and 1990s. Molecular techniques like gene targeting and amplification, using enzymes like restriction endonucleases and DNA ligases, laid the groundwork for genetic engineering [18]. The polymerase chain reaction (PCR) transformed gene amplification, while innovations like Taq polymerase and thermal cyclers optimized the process [23].

Computing accessibility surged with microcomputers like the Commodore PET, Apple II, and Tandy TRS-80, along with bioinformatics software like the GCG software suite [24] and DNASTAR [25], another sequence manipulation suite capable of assembling and analyzing Sanger sequencing data. Other sequence manipulation suites were developed to run on CP/M, Apple II, and Macintosh computers [26] in the years 1984 and 1985. Free code copies of this software were offered on demand by some developers. This propelled an upcoming software-sharing movement in the programming world [27,28].

The free software movement, led by the GNU project, promoted open-source bioinformatics tools. Major sequence databases (EMBL, GenBank, DDBJ) standardized data formatting and enabled global sharing. Bioinformatics journals, like CABIOS, which is now known as Bioinformatics (Oxford, England) accentuated computational methods' importance. Desktop workstations with Unix-like systems and scripting languages aided bioinformatics analyses, and scripting languages simplified tool development.

1990s: The genomics era and web-based bioinformatics

The genomics era began in the mid-1990s with the complete sequencing of the Haemophilus influenzae genome [29], initiating genome-scale analyses. This milestone was followed by the publication of the human genome at the beginning of the 21st century, which served as the definitive catalyst for the genomic era [30]. This transformative event spurred the design and development of several specialized Perl-based software to assemble whole-genome sequencing reads: PHRAP [31], Celera Assembler [32] among others.

Tim Berners-Lee's pioneering work at CERN in the early 1990s resulted in the World Wide Web, transforming global communication and ushering in an era of unprecedented access to information. With the advent of the internet, researchers gained a powerful platform to share and access vast amounts of biological data efficiently. This facilitated collaborative efforts in biology and genomics, leading to the establishment of foundational databases such as the EMBL Nucleotide Sequence Data Library [33] and the GenBank database became the responsibility of the NCBI [34] in 1992. Also, the famous NCBI website came online in 1994, featuring the efficient pairwise alignment tool BLAST [35]. After that, the world saw the birth of major databases we still rely on today: Genomes (1995), PubMed (1997), and Human Genome (1999) [36,37,38].

The proliferation of web-based resources transformed access to bioinformatics tools, democratizing their availability and usability for researchers worldwide. Through the development of web platforms, bioinformatics tools became more user-friendly and accessible. This shift enabled researchers to interact with sophisticated analytical tools without needing extensive computational expertise or access to specialized hardware. Consequently, the widespread adoption of web-based bioinformatics resources facilitated broader participation in genomic and molecular research, accelerating scientific discovery and collaboration on a global scale. Graphical web servers emerged as a convenient alternative to traditional UNIX-based systems, simplifying data analysis without the need for complex installations. The continued relevance of servers for scientific purposes is exemplified by the AlphaFold Server which uses the latest AlphaFold 3 model [39, J. et al. Accurate structure prediction of biomolecular interactions with AlphaFold 3. Nature 1–3 (2024) doi:10.1038/s41586-024-07487-w], released in 2024, to provide highly accurate biomolecular structure predictions in a unified platform.

The internet facilitated the dissemination of scientific research through online publications, challenging traditional print-based methods. Early initiatives like BLEND [40] paved the way for internet-based scientific publishing by shedding insights into the potentials and obstacles associated with using the internet for scientific publications. This study paved the way for leveraging the Internet for both data set storage and dissemination, leading up to the establishment of preprint servers like arXiv (est. 1991) [41] and bioRxiv (est. 2013) [42] which changed the way scientific findings are shared and accessed. These platforms democratized access to scientific knowledge by enabling researchers to share their work rapidly and openly, facilitating interdisciplinary collaborations and the cross-pollination of ideas.

The experimental determination of the first three-dimensional structure of a protein, specifically, myoglobin, occurred in 1958 via X-ray diffraction [4]. However, earlier groundwork by Pauling and Corey with the publication of two articles in 1951 that reported the prediction of α -helices and β -sheets [43] laid the foundation for predicting protein structures. Similar to advances in other biological sciences, the utilization of computers has made it feasible to conduct calculations aimed at predicting the secondary and tertiary structure of proteins, with varying levels of confidence. This capability has been notably enhanced by the development of fold recognition algorithms, also known as threading algorithms [44,45]. However, proteins are dynamic entities, requiring advanced biophysical models to describe their interactions and movements accurately. Force fields have been formulated to describe the interactions among atoms, enabling the introduction of tools for modeling the molecular dynamics of proteins during the 1990s [46]. Used to study the behavior and interactions of atoms and molecules over time, molecular dynamics simulations calculate the positions and velocities of atoms based on physical principles. Despite the theoretical advancements and availability

of tools, executing molecular dynamics simulations remained challenging in practice due to the substantial computational resources they demanded.

Graphical processing Units (GPUs) have made molecular dynamics more accessible [47], with applications extending to other bioinformatics fields requiring intensive computation. However, the internet's role in data dissemination, coupled with increasing computational power, has led to the proliferation of 'Big Data' in bioinformatics.

2000-2010: High-throughput sequencing and big data

Second-generation sequencing technologies democratized high-throughput bioinformatics. For example '454' pyrosequencing, a high-throughput DNA sequencing technique played a significant role in advancing genomics research by enabling rapid and cost-effective sequencing of DNA samples, particularly for applications such as whole-genome sequencing [48], but computational challenges arose with increased data volumes. Decreasing sequencing costs resulted in more data being generated, emphasizing data organization and accessibility. Specialized repositories and standardization efforts were needed to ensure data interoperability. High-performance computing adaptation became vital to address the increased amounts of data within bioinformatics projects. The surge in bioinformatics projects, accompanied by a vast influx of data, prompted adjustments from funding bodies to accommodate the demand for high-performance computing resources and collaborative initiatives.

While basic computer setups suffice for some projects, others demand complex infrastructures and substantial expertise. Government-sponsored entities like <u>Compute Canada</u>, <u>New York State's High-Performance Computing Program</u>, <u>The European Technology Platform for High-Performance Computing</u>, and <u>National Center for High-Performance Computing</u> served researchers' computational needs. Companies like Amazon, Microsoft, and Google, among many others, offer bioinformatics and life sciences services, emphasizing the field's importance.

Table 1. Organizations providing High-Performance Computing Resources for Bioinformatics and Life Sciences

Organization	Computing Resources
Compute Canada	Provides high-performance computing resources and support services to researchers and innovators across Canada. They offer supercomputers, cloud platforms, data storage, and training programs to advance scientific research and innovation in various fields.
New York State's High-Performance Computing Program	Provides researchers, businesses, and educational institutions with access to high-performance computing (HPC) resources and expertise to support their computational research and development efforts.
The European Technology Platform for High-Performance Computing	Fosters collaboration among industry, research, and academic stakeholders to advance high-performance computing (HPC) technology in Europe.
National Center for High-Performance Computing	Facility for high-performance computing (HPC) resources including large-scale computational science and engineering, cluster and grid computing, middleware development, visualization and virtual reality, data storage, networking, and HPC-related training.

Organization	Computing Resources
National Center for Supercomputing Applications	Offers high-performance computing resources such as the Blue Waters supercomputer, provides advanced data storage solutions, data analysis, and visualization tools, and supports interdisciplinary research in fields such as astrophysics, climate modeling, and genomics.
Oak Ridge Leadership Computing Facility	Provides supercomputing resources, such as the Summit supercomputer, for scientific research, offers support services including software development, data storage, and visualization, and facilitates research in various fields including climate science, biology, and materials science.
Swiss National Supercomputing Centre	Provides high-performance computing systems including the Piz Daint supercomputer, offers cloud computing services, data management, and user support, and facilitates scientific research in areas such as climate modeling, physics, and life sciences.
Barcelona Supercomputing Center	Provides access to MareNostrum, one of the most powerful supercomputers in Europe, offers resources for high-performance computing, data storage, and computational sciences, and supports research in fields including bioinformatics, computational biology, and engineering.
Japan's RIKEN Center for Computational Science	Houses the Fugaku supercomputer, one of the world's fastest supercomputers, provides resources for computational science, data processing, and artificial intelligence, and supports research in fields such as life sciences, materials science, and disaster prevention.
National Supercomputing Centre Singapore	Provides high-performance computing resources and support services, offers data storage, cloud computing, and software development services, and supports research in fields including bioinformatics, environmental modeling, and smart cities.

Community computing platforms democratized participation and expanded bioinformatics research's reach. Platforms like BOINC [49] enabled broad participation in bioinformatics. Experts can submit computing tasks to BOINC, while non-experts and science enthusiasts can volunteer their computer resources to process these tasks. Several life sciences projects are available through BOINC, including protein-ligand docking, malaria simulations, and protein folding [49].

2010-Today: The present and future

The integration of computers into biology has ushered in a new era of research possibilities, allowing for increasingly complex studies. While before, the focus was on individual genes or proteins, advancements today enable the analysis of entire genomes or proteomes [50]. This shift toward a holistic approach in biology is evident in disciplines like genomics, proteomics, and glycomics, which have limited interconnection between them.

The next leap at the intersection of computing and the life sciences lies in modeling entire living organisms and their environments simultaneously, integrating all molecular categories. This has already been achieved in a whole cell model of Mycoplasma genitalium, in which all its genes, products and their known metabolic interactions have been reconstructed [51]. Driven by advancements in measurement techniques, improved computational performance and artificial intelligence (AI) techniques, whole-cell modeling is increasingly becoming realistic and feasible. In contrast to traditional bottom-up approaches relying on molecular interaction networks, a predictive model has been developed for genome-wide phenotypes of budding yeast using deep learning [52]. The main applications of whole-cell modeling have been in producing useful substances and discovering drugs, such as antimicrobials [53,54,55,56] since whole-cell modeling was first directed

toward unicellular organisms. Meanwhile, models of cultured human cells have also been developed, which have found applications in cell differentiation and medical research [57]. The possibility of modeling entire multicellular organisms may not be far off, considering the rapid pace of technological and computational advancements like artificial intelligence (AI).

Advances in Generative AI for Proteins

Introduction

- Background on proteins as key molecular machines in biology
- Promise of generative AI to accelerate protein discovery and engineering
- Overview of scope covering recent advances in last 1-2 years

Structure Prediction

- AlphaFold2 as a breakthrough method for structure prediction
- Novel model architecture and training methodology
- Examples of new biological insights from predicted structures

Function Prediction

- Using predicted structures to infer protein functions
- Structure-based identification of catalytic and binding sites
- Case studies of novel enzyme functions discovered

Designing Novel Proteins

- Generative models for designing functional protein sequences
- Leveraging structural constraints for optimized protein engineering
- Applications in industrial enzymes, therapeutics, biomaterials

Interaction Prediction

- Modeling protein-protein interactions with graph networks
- Structure-based prediction of protein-drug bindings
- Applications in drug discovery and toxicity screening

Outlook

- Challenges and next steps in improving accuracy
- Hybrid physics- and data-driven approaches
- Ethical considerations in synthetic protein design

Advances in Generative AI for RNA

Introduction

- Background on key roles of RNA in biology
- Promise of generative models to advance RNA research

Scope focused on latest advances in RNA prediction and design

Structure Prediction

- · Transformer-based prediction of RNA folding
- Novel architectures for incorporating chemical constraints
- Improved accuracy on complex structures like ribosomes

Function Prediction

- Inferring RNA functions from sequence and structure
- Identifying motifs, domains, and atomic binding sites
- Applications in understanding long noncoding RNAs

Interaction Prediction

- Graph neural networks for RNA-protein interactions
- Structure-augmented modeling of splice sites
- Predicting RNA base editing targets

Design of RNA Therapeutics

- Generative models for optimizing siRNA, antisense design
- Reinforcement learning for chemical modification patterns
- · Progress in computational RNA-targeted drug design

Outlook

- Key challenges in prediction of long RNA structures
- Design of RNA for self-assembly and scaffolding
- Ethical use of synthetic RNA technologies

Advances in Generative AI for DNA

Introduction

- Background on the role of DNA as a carrier of genetic information
- Promise of generative models to advance DNA research
- Scope focused on the latest advances in DNA prediction and design

Sequence Modeling

Transformer architectures for modeling DNA sequences

The HyenaDNA paper introduces a new genomic foundation model called HyenaDNA that can process DNA sequences at single nucleotide resolution with ultralong context lengths up to 1 million base pairs - a 500x increase over previous transformer models. HyenaDNA uses a parameter-efficient convolutional architecture that allows it to scale subquadratically with sequence length, enabling the use of full genome-scale context. On a range of regulatory genomics prediction tasks, HyenaDNA

matches or exceeds the performance of previous state-of-the-art models while using 1500x fewer parameters and 3200x less pretraining data. HyenaDNA also demonstrates the ability to perform challenging species classification using the full mutational profile visible at 1 million base pairs of context. The authors explore new training techniques to enable ultralong sequence modeling as well as prompt-based tuning methods for rapidly adapting to new tasks without updating pre-trained weights [58].

Pretraining on large genomic datasets

Applications in variant calling and annotation

Regulation Prediction

- Graph neural networks for modeling 3D genome architecture
- Predicting enhancer-promoter interactions and expression
- Design of synthetic promoters and enhancers

Genome Editing

- · Generative models for CRISPR guide design
- Contextual prediction of on-target editing efficacy
- Modeling of off-target effects during optimization

DNA Data Generation

- Variational autoencoders for realistic DNA sequences
- Generating paired genomic-transcriptomic data
- Applications in training genome interpretation models

Outlook

Challenges in predicting long-range chromatin interactions

The development of transformer architectures like HyenaDNA that can leverage genome-scale context creates new opportunities for understanding long-range chromatin interactions, gene regulation, and intercellular networks. However, significant challenges remain in improving model accuracy and uncertainty quantification. Hybrid physics- and data-driven approaches may help address these gaps. Safety considerations around synthetic genome design also warrant further research to ensure responsible innovation as these generative capabilities advance. Overall, the future looks bright for generative models that can capture both local mutations and global patterns critical to biology [58].

Responsible design of synthetic genomes

Ethical considerations for human genome editing

Generative AI for Autonomous Experimentation

Introduction

- Promise of generative models to accelerate scientific discovery
- Rise of autonomous labs and robotics for automated experimentation
- Overview of generative Al's role in self-driving research

Closed-Loop Systems

- Integrating computational hypothesis generation with robotic wet lab testing
- Reinforcement learning pipelines for autonomous optimization
- Case studies in materials science, drug discovery

Automated Experiment Design

- Using generative models to design novel compounds, genes
- Leveraging simulations to predict experimental outcomes
- Robotic execution of designed experiments

Adaptive Sampling

- · Active learning to iteratively select most informative experiments
- Bayesian optimization powered by neural networks
- · Applications in probing molecular design spaces

Real-Time Learning

- Deploying models on lab edge devices
- Online learning from experimental data streams
- · Improving models and experiment plans on-the-fly

Outlook

- Key challenges around model accuracy and integration
- The future of data-driven, self-driving laboratories
- · Risks of full automation and need for human oversight

Conclusion

- Generative AI as a powerful tool for autonomous experimentation
- Accelerating discovery alongside human researchers
- Responsible implementation will maximize benefit

Generative AI and the Biosecurity Landscape

Introduction

- Background on biosecurity threats from natural, accidental, and intentional pathogens
- Rise of generative AI as a dual use technology for biodefense and misuse

Enhanced Risks

- · Automated bioweapon design with generative models
- Relatively low computing needs to generate dangerous agents
- Challenges detecting artificially generated sequences/organisms

Enhanced Response Capabilities

- Generative models for vaccine and therapeutic design
- High-throughput testing of countermeasures with synthetic data
- Al for early detection of emergent pathogens and outbreaks

Recommendations

- Increased oversight for generative model development/release
- Expanding biosecurity legislation and regulations
- Fostering open research and global cooperation

Outlook

- Trajectory toward increasingly powerful generative biological capabilities
- Need for preventative ethics research and guidance
- Maintaining public trust and avoiding overreaction

Conclusion

- Balancing generative Al's benefits and risks in biology
- Importance of thoughtful governance and responsible innovation
- · Staying ahead of the curve on biosecurity

Policy Responses to Generative AI in Biology

Introduction

- Background on rise of powerful generative models for biology
- Overview of risks like bioweapons, environmental damage
- Need for governance to ensure responsible development

Self-Governance by Developers

- Voluntary guidelines on ethical AI by corporations
- Limiting access to certain capabilities like human editing
- Issues with self-regulation and transparency

Governmental Regulations

- New biosecurity regulations on certain AI technologies
- Restrictions on use of synthetic biology IP
- Challenges with fast pace of technology change

International Governance

- Proposals for global observatory to monitor risks
- Treaties restricting development of bioweapons
- Difficulty achieving consensus and compliance

Public Deliberation and Ethics

- Activities to involve broader stakeholders
- Gathering public attitudes on acceptable uses of generative bio Al
- Informing policy with deliberative democracy

Outlook

- Likely increase in debate and policy activity in this space
- Balancing innovation and security will be a key challenge
- Importance of thoughtful multidisciplinary discourse

Conclusion

- No easy policy solutions, but inaction also carries risks
- Policy should enable innovation but promote responsible use
- This will require sustained public deliberation and coordination

References

1. Computers in the study of evolution

Science Progress (1967) https://pubmed.ncbi.nlm.nih.gov/4859964

2. The chemical basis of morphogenesis

Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences (1952-08-14) https://royalsocietypublishing.org/doi/10.1098/rstb.1952.0012 DOI: 10.1098/rstb.1952.0012

3. Scientific uses of the MANIAC

HL Anderson

Journal of Statistical Physics (1986-06-01) https://doi.org/10.1007/BF02628301

DOI: 10.1007/bf02628301

A three-dimensional model of the myoglobin molecule obtained by x-ray analysis 4.

JC Kendrew, G Bodo, HM Dintzis, RG Parrish, H Wyckoff, DC Phillips Nature (1958-03-08) https://pubmed.ncbi.nlm.nih.gov/13517261 DOI: 10.1038/181662a0

Structure of myoglobin: A three-dimensional Fourier synthesis at 2 A. resolution 5.

JC Kendrew, RE Dickerson, BE Strandberg, RG Hart, DR Davies, DC Phillips, VC Shore Nature (1960-02-13) https://pubmed.ncbi.nlm.nih.gov/18990802

DOI: 10.1038/185422a0

The amino-acid sequence in the glycyl chain of insulin. 1. The identification of lower 6. peptides from partial hydrolysates

F Sanger, EOP Thompson

Biochemical Journal (1953-02) https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1198157/

The amino-acid sequence in the glycyl chain of insulin. 2. The investigation of peptides 7. from enzymic hydrolysates

F Sanger, EOP Thompson

Biochemical Journal (1953-02) https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1198158/

A method for the determination of amino acid sequence in peptides 8.

P Edman

Archives of Biochemistry (1949-07) https://pubmed.ncbi.nlm.nih.gov/18134557

The origins of bioinformatics 9.

JB Hagen

Nature Reviews. Genetics (2000-12) https://pubmed.ncbi.nlm.nih.gov/11252753

DOI: 10.1038/35042090

10. Comprotein: a computer program to aid primary protein structure determination

Margaret Oakley Dayhoff, Robert S Ledley

Proceedings of the December 4-6, 1962, fall joint computer conference on - AFIPS '62 (Fall) (1962) http://portal.acm.org/citation.cfm?doid=1461518.1461546

DOI: 10.1145/1461518.1461546

https://febs.onlinelibrary.wiley.com/toc/14321033/5/2 11.

12. Collecting, comparing, and computing sequences: the making of Margaret O. Dayhoff's Atlas of Protein Sequence and Structure, 1954-1965

Bruno J Strasser

Journal of the History of Biology (2010) https://pubmed.ncbi.nlm.nih.gov/20665074

DOI: 10.1007/s10739-009-9221-0

13. A general method applicable to the search for similarities in the amino acid sequence of two proteins

SB Needleman, CD Wunsch

Journal of Molecular Biology (1970-03) https://pubmed.ncbi.nlm.nih.gov/5420325

DOI: 10.1016/0022-2836(70)90057-4

14. Progressive sequence alignment as a prerequisite to correct phylogenetic trees

DF Feng, RF Doolittle

Journal of Molecular Evolution (1987) https://pubmed.ncbi.nlm.nih.gov/3118049

DOI: 10.1007/bf02603120

15. **CLUSTAL: a package for performing multiple sequence alignment on a microcomputer**

DG Higgins, PM Sharp

Gene (1988-12-15) https://pubmed.ncbi.nlm.nih.gov/3243435

DOI: <u>10.1016/0378-1119(88)90330-7</u>

16. Clustal Omega, accurate alignment of very large numbers of sequences

Fabian Sievers, Desmond G Higgins

Methods in Molecular Biology (Clifton, N.J.) (2014) https://pubmed.ncbi.nlm.nih.gov/24170397

DOI: <u>10.1007/978-1-62703-646-7_6</u>

17. A new method for sequencing DNA

AM Maxam, W Gilbert

Proceedings of the National Academy of Sciences of the United States of America (1977-02)

https://pubmed.ncbi.nlm.nih.gov/265521

DOI: 10.1073/pnas.74.2.560

18. Summary statement of the Asilomar conference on recombinant DNA molecules.

P Berg, D Baltimore, S Brenner, RO Roblin, MF Singer

Proceedings of the National Academy of Sciences of the United States of America (1975-06) https://www.ncbi.nlm.nih.gov/pmc/articles/PMC432675/

19. A strategy of DNA sequencing employing computer programs.

R Staden

Nucleic Acids Research (1979-06-11) https://www.ncbi.nlm.nih.gov/pmc/articles/PMC327874/

20. Evolutionary trees from DNA sequences: a maximum likelihood approach

J Felsenstein

Journal of Molecular Evolution (1981) https://pubmed.ncbi.nlm.nih.gov/7288891

DOI: 10.1007/bf01734359

21. Probability distribution of molecular evolutionary trees: a new method of phylogenetic inference

B Rannala, Z Yang

Journal of Molecular Evolution (1996-09) https://pubmed.ncbi.nlm.nih.gov/8703097

DOI: 10.1007/bf02338839

22. A biologist's guide to Bayesian phylogenetic analysis

Fabrícia F Nascimento, Mario Dos Reis, Ziheng Yang

Nature Ecology & Evolution (2017-10) https://pubmed.ncbi.nlm.nih.gov/28983516

DOI: 10.1038/s41559-017-0280-x

23. The Nobel Prize in Chemistry 1993

NobelPrize.org

https://www.nobelprize.org/prizes/chemistry/1993/mullis/lecture/

24. A comprehensive set of sequence analysis programs for the VAX

J Devereux, P Haeberli, O Smithies

Nucleic Acids Research (1984-01-11) https://pubmed.ncbi.nlm.nih.gov/6546423

DOI: 10.1093/nar/12.1part1.387

25. DNASTAR's Lasergene Sequence Analysis Software

Timothy G Burland

Bioinformatics Methods and Protocols (1999) https://doi.org/10.1385/1-59259-192-2:71

ISBN: 9781592591923

26. Apple II PASCAL programs for molecular biologists

B Malthiery, B Bellon, D Giorgi, B Jacq

Nucleic Acids Research (1984-01-11) https://pubmed.ncbi.nlm.nih.gov/6320099

DOI: 10.1093/nar/12.1part2.569

27. The GNU Manifesto - GNU Project - Free Software Foundation

https://www.gnu.org/gnu/manifesto.en.html

28. https://www.researchgate.net/publication/221307757 The Free Software Movement and the GNULinux Operating System

29. Whole-genome random sequencing and assembly of Haemophilus influenzae Rd

RD Fleischmann, MD Adams, O White, RA Clayton, EF Kirkness, AR Kerlavage, CJ Bult, JF Tomb, BA Dougherty, JM Merrick

Science (New York, N.Y.) (1995-07-28) https://pubmed.ncbi.nlm.nih.gov/7542800

DOI: 10.1126/science.7542800

30. The sequence of the human genome

JC Venter, MD Adams, EW Myers, PW Li, RJ Mural, GG Sutton, HO Smith, M Yandell, CA Evans, RA Holt, ... X Zhu

Science (New York, N.Y.) (2001-02-16) https://pubmed.ncbi.nlm.nih.gov/11181995

DOI: 10.1126/science.1058040

31. Consed: a graphical tool for sequence finishing

D Gordon, C Abajian, P Green

Genome Research (1998-03) https://pubmed.ncbi.nlm.nih.gov/9521923

DOI: <u>10.1101/gr.8.3.195</u>

32. A whole-genome assembly of Drosophila

EW Myers, GG Sutton, AL Delcher, IM Dew, DP Fasulo, MJ Flanigan, SA Kravitz, CM Mobarry, KH Reinert, KA Remington, ... JC Venter

Science (New York, N.Y.) (2000-03-24) https://pubmed.ncbi.nlm.nih.gov/10731133

DOI: 10.1126/science.287.5461.2196

33. The EMBL data library

CM Rice, R Fuchs, DG Higgins, PJ Stoehr, GN Cameron

Nucleic Acids Research (1993-07-01) https://pubmed.ncbi.nlm.nih.gov/8332519

DOI: 10.1093/nar/21.13.2967

34. **GenBank**

Dennis A Benson, Mark Cavanaugh, Karen Clark, Ilene Karsch-Mizrachi, David J Lipman, James Ostell, Eric W Sayers

Nucleic Acids Research (2013-01) https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3531190/
DOI: 10.1093/nar/gks1195

35. BLAST: at the core of a powerful and diverse set of sequence analysis tools

Scott McGinnis, Thomas L Madden

Nucleic Acids Research (2004-07-01) https://www.ncbi.nlm.nih.gov/pmc/articles/PMC441573/
DOI: 10.1093/nar/gkh435

36. Genomes OnLine Database (GOLD): a monitor of genome projects world-wide

Axel Bernal, Uy Ear, Nikos Kyrpides

Nucleic Acids Research (2001-01-01) https://www.ncbi.nlm.nih.gov/pmc/articles/PMC29859/

37. PubMed

PubMed

https://pubmed.ncbi.nlm.nih.gov/

- 38. https://academic.oup.com/nar/article/26/1/94/2379498
- 39. Abramson

40. The BLEND system Programme for the study of some 'electronic journals'*

B Shackel

Ergonomics (1982-04) http://www.tandfonline.com/doi/abs/10.1080/00140138208924954
DOI: 10.1080/00140138208924954

- 41. arXiv.org e-Print archive https://arxiv.org/
- 42. bioRxiv.org the preprint server for Biology https://www.biorxiv.org/

43. Configurations of Polypeptide Chains With Favored Orientations Around Single Bonds: Two New Pleated Sheets

L Pauling, RB Corey

Proceedings of the National Academy of Sciences of the United States of America (1951-11) https://pubmed.ncbi.nlm.nih.gov/16578412

DOI: 10.1073/pnas.37.11.729

44. Sixty-five years of the long march in protein secondary structure prediction: the final stretch?

Yuedong Yang, Jianzhao Gao, Jihua Wang, Rhys Heffernan, Jack Hanson, Kuldip Paliwal, Yaoqi Zhou

Briefings in Bioinformatics (2018-05-01) https://pubmed.ncbi.nlm.nih.gov/28040746

DOI: 10.1093/bib/bbw129

45. Computational methods for protein structure prediction and modeling

New York, N.Y.: Springer

(2007) http://archive.org/details/computationalmet0000unse_u4q5

ISBN: 9780387333212

46. Molecular dynamics simulations: advances and applications

Adam Hospital, Josep Ramon Goñi, Modesto Orozco, Josep L Gelpí

Advances and Applications in Bioinformatics and Chemistry: AABC (2015-11-19)

https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4655909/

DOI: 10.2147/aabc.s70333

47. To Milliseconds and Beyond: Challenges in the Simulation of Protein Folding

Thomas J Lane, Diwakar Shukla, Kyle A Beauchamp, Vijay S Pande *Current opinion in structural biology* (2013-02)

https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3673555/

DOI: 10.1016/j.sbi.2012.11.002

48. Next-Generation Sequencing and Its Application: Empowering in Public Health Beyond Reality

Nidhi Gupta, Vijay K Verma

Microbial Technology for the Welfare of Society (2019-09-13)

https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7122948/

DOI: <u>10.1007/978-981-13-8844-6 15</u>

49. BOINC: A Platform for Volunteer Computing

David P Anderson

Journal of Grid Computing (2020-03-01) https://doi.org/10.1007/s10723-019-09497-9

DOI: <u>10.1007/s10723-019-09497-9</u>

50. Structural proteomics by NMR spectroscopy

Joon Shin, Woonghee Lee, Weontae Lee

Expert Review of Proteomics (2008-08) https://pubmed.ncbi.nlm.nih.gov/18761469

DOI: 10.1586/14789450.5.4.589

51. A whole-cell computational model predicts phenotype from genotype

Jonathan R Karr, Jayodita C Sanghvi, Derek N Macklin, Miriam V Gutschow, Jared M Jacobs, Benjamin Bolival, Nacyra Assad-Garcia, John I Glass, Markus W Covert

Cell (2012-07-20) https://pubmed.ncbi.nlm.nih.gov/22817898

DOI: 10.1016/j.cell.2012.05.044

52. Using deep learning to model the hierarchical structure and function of a cell

Jianzhu Ma, Michael Ku Yu, Samson Fong, Keiichiro Ono, Eric Sage, Barry Demchak, Roded Sharan, Trey Ideker

Nature Methods (2018-04) https://www.nature.com/articles/nmeth.4627

DOI: 10.1038/nmeth.4627

53. Why Build Whole-Cell Models?

Javier Carrera, Markus W Covert

Trends in Cell Biology (2015-12) https://pubmed.ncbi.nlm.nih.gov/26471224

DOI: 10.1016/j.tcb.2015.09.004

54. The future of whole-cell modeling

Derek N Macklin, Nicholas A Ruggero, Markus W Covert

Current Opinion in Biotechnology (2014-08) https://pubmed.ncbi.nlm.nih.gov/24556244

DOI: 10.1016/j.copbio.2014.01.012

55. Computer-Aided Whole-Cell Design: Taking a Holistic Approach by Integrating Synthetic With Systems Biology

Lucia Marucci, Matteo Barberis, Jonathan Karr, Oliver Ray, Paul R Race, Miguel de Souza Andrade, Claire Grierson, Stefan Andreas Hoffmann, Sophie Landon, Elibio Rech, ... Christopher Woods

Frontiers in Bioengineering and Biotechnology (2020-08-07)

https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7426639/

DOI: 10.3389/fbioe.2020.00942

56. Accelerated discovery via a whole-cell model

Jayodita C Sanghvi, Sergi Regot, Silvia Carrasco, Jonathan R Karr, Miriam V Gutschow, Benjamin Bolival, Markus W Covert

Nature methods (2013-12) https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3856890/

DOI: 10.1038/nmeth.2724

57. A blueprint for human whole-cell modeling

Balázs Szigeti, Yosef D Roth, John AP Sekar, Arthur P Goldberg, Saahith C Pochiraju, Jonathan R Karr

Current opinion in systems biology (2018-02)

https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5966287/

DOI: <u>10.1016/j.coisb.2017.10.005</u>

58. HyenaDNA: Long-Range Genomic Sequence Modeling at Single Nucleotide Resolution

Eric Nguyen, Michael Poli, Marjan Faizi, Armin Thomas, Callum Birch-Sykes, Michael Wornow, Aman Patel, Clayton Rabideau, Stefano Massaroli, Yoshua Bengio, ... Chris Ré arXiv (2023) https://doi.org/gs3v5j

DOI: 10.48550/arxiv.2306.15794