

Bioinformatics

B**Mark A. Ragan***The University of Queensland, Australia*

INTRODUCTION

Over the past 30 years, bioinformatics has emerged as new discipline at the interface of molecular bioscience with mathematics, computer science and information technology. Bioinformatics is driven by data arising from high-throughput technologies in molecular bioscience including DNA and genome sequencing, gene expression analysis, protein and RNA structure characterisation, and bio-imaging. To enable biological discovery, bioinformatics draws on and extends technologies for data capture, management, integration and mining, computing, and communication technology including the Internet. The rise of genomics, from the first bacterial and model-organism projects to the Human Genome Projects and the thousands of genome projects that have followed, has been a key driver for bioinformatics, and in turn has enabled these projects to be completed and their results applied. Genomics, however, was never an end unto itself, but rather was intended to enable the understanding of complex biological systems. Bioinformatics continues to evolve in support of its constituent domains and, increasingly, their integration into genome-scale molecular systems biology.

This article presents bioinformatics first from the perspective of computer science and IT, then from the perspective of bioscience. In practice these perspectives often merge, making bioinformatics a rich, vibrant area of multidisciplinary research and application.

BACKGROUND

The term *bioinformatics* was introduced in 1970 in reference to the study of informatic processes in biological systems (Hogeweg, 2011). In this original usage, *bioinformatics* encompassed “how living systems gather, process, store and use information” (Nurse, 2008). Never widely adopted, this usage was superseded in the late 1980s when bioinformatics, as presently understood, emerged as a new field at the interface of molecular bioscience with computer science and information technology (Dickson, 1987). Today bioinformatics builds on mathematics, statistics and algorithmics, and finds applications across the biosciences particularly in genomics, proteomics, structural biology and molecular systems biology. Biology is increasingly an information science, with bioinformatics a key enabling technology.

Other disciplines have developed at the bioscience - computer science - IT interface, and there is little consensus on where boundaries should be drawn among them. Bioinformatics is sometimes said to focus on the development and application of methods and software tools to acquire, manage, analyse and/or visualise biological data, whereas *computational biology* refers more to the application of these methods and tools to theoretical or applied biological questions (Huerta *et al.*, 2000). *Biomathematics* or *mathematical biology* involves the development or use of mathematical modeling or simulation, while *biostatistics* emphasises experimental design and statistical analysis. *Mo-*

lecular systems biology focuses on the inference or analysis of networks of genes, proteins and/or other cellular molecules, while *synthetic biology* applies these technologies to design and engineer new biological functions or organisms.

BIOINFORMATICS FROM THE PERSPECTIVE OF COMPUTER SCIENCE AND INFORMATION TECHNOLOGY

One way of exploring the interface between molecular bioscience and IT is to track experimental data from their generation, capture and retrieval, through their aggregation and dissemination *via* international data services, to their subsequent analysis. Here I deconstruct data analysis into data models, algorithms, analytical methods and software, workflows and visualisation. This trajectory is common to scientific data, although bioinformatics is notable for its culture of open data, well-established data formats and standards, and data reuse facilitated by large international repositories with associated data services.

Data Generation, Storage and Retrieval

Instruments and experiments generate diverse data types in molecular bioscience. Capturing these primary data and the associated metadata, and managing their storage and retrieval, are primary activities in bioinformatics. The quantities of data generated by DNA-sequencing platforms, in particular, are such that raw data are no longer archived; rather, bioinformatic methods are used to assess quality and extract summaries. Data formats are specific to experimental technologies and, to some extent, instrument manufacturers. In some areas of molecular bioscience, standards have been developed to ensure that data can be interpreted unambiguously and, in principle, the experiment can be reproduced. For example, the MIAME (Minimum Information About a Mi-

croarray Experiment) standard (Brazma *et al.*, 2001) specifies how experimental design, laboratory protocols, biological samples, microarray platforms, and raw and processed data must be described, and recommends the use of certain data formats and ontologies. The MIBBI project promotes the coherence of minimum-standards checklists across nearly 40 areas of bioscience and biomedicine (Taylor *et al.*, 2008).

Public Data Resources

Newly generated biomolecular data (*e.g.* DNA and protein sequences, protein structures, gene-expression data) are submitted to public data repositories, where they are assigned unique persistent identifiers; all major bioscience journals require new data to be so identified. The main international collection centres are the US National Centre for Biotechnology Information (NCBI), the EMBL European Bioinformatics Institute (EBI), and the DNA Data Bank of Japan (DDBJ). Individually and in collaboration with each other, these centres carry out further quality control on incoming data, conduct research in bioinformatics, promote best practice in bioscience data management and analysis, and provide comprehensive online data services (*e.g.* search, retrieval, integrative analyses over multiple data sources, and links to journal articles and patents) which are cost-free at the “point of use” for the international research community. Other public data resources serve specific areas of molecular bioscience, *e.g.* protein structure. Increasingly, the largest projects in molecular bioscience maintain their own public data resources, *e.g.* The Cancer Genome Atlas and the International Cancer Genome Consortium. Both large and small data sources are reviewed in the annual Database issue of the journal *Nucleic Acids Research*.

Data Formats and Models

The need for data integration and re-use has driven the development of standard data presentation

10 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

www.igi-global.com/chapter/bioinformatics/183756

Related Content

Using Metaheuristics as Soft Computing Techniques for Efficient Optimization

Sergio Nesmachnow (2015). *Encyclopedia of Information Science and Technology, Third Edition* (pp. 7390-7399).

www.irma-international.org/chapter/using-metaheuristics-as-soft-computing-techniques-for-efficient-optimization/112436

Financial Risk Intelligent Early Warning System of a Municipal Company Based on Genetic Tabu Algorithm and Big Data Analysis

Hui Liu (2022). *International Journal of Information Technologies and Systems Approach* (pp. 1-14).

www.irma-international.org/article/financial-risk-intelligent-early-warning-system-of-a-municipal-company-based-on-genetic-tabu-algorithm-and-big-data-analysis/307027

Accident Causation Factor Analysis of Traffic Accidents using Rough Relational Analysis

Caner Erdenand Numan Çelebi (2016). *International Journal of Rough Sets and Data Analysis* (pp. 60-71).

www.irma-international.org/article/accident-causation-factor-analysis-of-traffic-accidents-using-rough-relational-analysis/156479

The Concept of the Shapley Value and the Cost Allocation Between Cooperating Participants

Alexander Kolker (2018). *Encyclopedia of Information Science and Technology, Fourth Edition* (pp. 2095-2107).

www.irma-international.org/chapter/the-concept-of-the-shapley-value-and-the-cost-allocation-between-cooperating-participants/183923

State of the Art and Future Trends of Datacenter Networks

George Michelogiannakis (2015). *Encyclopedia of Information Science and Technology, Third Edition* (pp. 1883-1892).

www.irma-international.org/chapter/state-of-the-art-and-future-trends-of-datacenter-networks/112593