# Interpreting and integrating big data in the life sciences

Serghei Mangul[1][2]

[1]Department of Computer Science, University of California Los Angeles, Los Angeles, CA, USA

[2]Institute for Quantitative and Computational Biosciences, University of California Los Angeles, Los Angeles, CA, USA

Email: smangul@ucla.edu

## Abstract (200 words)

Recent advances in omics technologies have led to the broad applicability of computational techniques across various domains of life science and medical research. These technologies provide an unprecedented opportunity to collect omics data from hundreds of thousands of individuals and to study gene-disease association without the aid of prior assumptions about the trait biology. Despite the many advantages of modern omics technologies, interpretations of big data produced by such technologies require advanced computational algorithms. Below I outline key challenges that biomedical researches are facing when interpreting and integrating big omics data. I discuss the reproducibility aspect of big data analysis in the life sciences and review current practices in reproducible research. Finally, I explain the skills which biomedical researchers need to acquire in order to independently analyze big omics data.
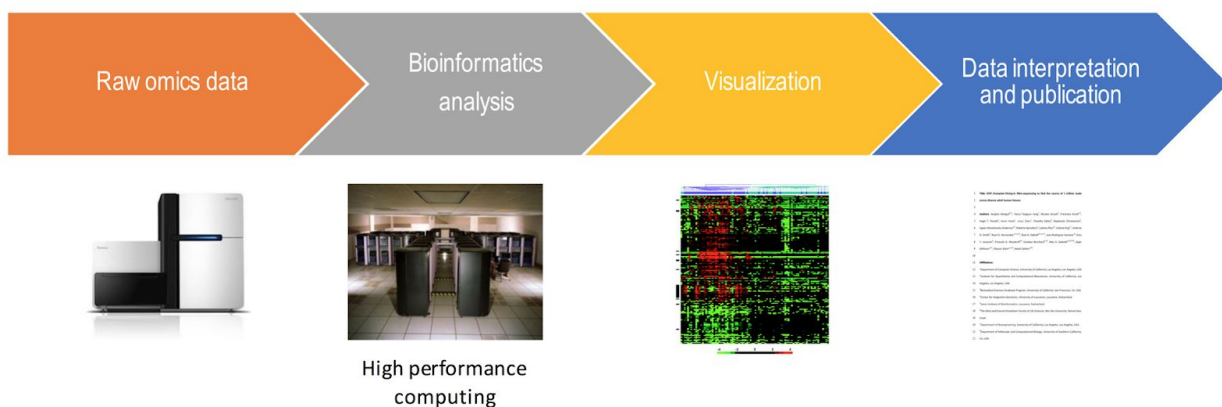
## Introduction

Recent advances in omics technologies have led to the broad applicability of modern high-throughput technologies across various domains of life science and medical research. These technologies are capable of generating big-data sets across large-scale clinical cohorts allows connecting complex diseases to relevant genomic features. However, analysis of big data requires the use of sophisticated bioinformatics algorithms capable of differentiating technical noise from the biological signal in the data. The analysis of big data in life sciences typically starts with the analysis of raw data and concludes with data visualization and interpretation of patterned data produced by analyses (**Figure 1**). While these techniques expand analytical opportunities, a researcher must have adequate computational skills to properly use bioinformatics algorithms. Learning the computational skills required s for analysis and interpretation of big omics data can be challenging to many life science and medical researchers. At present, the bioinformatics community mechanisms to researchers in effectively learning to analyze big omics data.

Ability to leverage various types of omics datasets from large-scale clinical cohorts is essential to studying the functional mechanisms underlying the connections between genetics, immune system, and disease etiology. For example, the availability of rich omics data generated by the TCGA consortium[1] provides an unprecedented opportunity for the discovery of how genetic variation affects the development, progression, and drug response of cancerous

tumors[2]. Additionally, bioinformatics methods for large-scale clinical cohorts promise to identify novel markers prognostic of disease risk across a variety of diseases, including cancer.

Despite the increasing size and complexity of datasets in the biological and medical sciences, many biomedical researchers today lack sufficient computational skills to analyze the large-scale data they generate. At present, the digital gap in contemporary biology limits the potential of these biomedical researchers to creatively explore their data. Further, the digital gap limits the collaborative potential of these biomedical researchers and computational scientists[3]. Training life sciences in computational techniques can potentially narrow this digital gap, expand the skills of biomedical researchers, and improve the ability of biomedical researchers to leverage the consultation services they seek with computational scientists.

**Figure 1.** Workflow of big data analysis and interpretation in the life sciences.

**Bioinformatics methods in life science research**

During the past decade, the rapid advancement of omics technologies has led to the development of an enormous amount and diversity of bioinformatics algorithms across various fields of modern biology[4]. hen applied to high-dimensional clinical datasets, bioinformatics algorithms identification of novel disease subtypes and discovery of novel markers which may be prognostic of disease risk. Such bioinformatics algorithms are usually encapsulated as computational software tools[5]. The majority of bioinformatics tools are designed for UNIX operating system, which requires a user to operate the tools using command line--without the benefit of a graphical user interface (GUI).

**Barriers in interpreting and integrating big data in the life sciences**

A major barrier in interdisciplinary studies is a lack of a common communication style between researchers trained and working in increasingly specialized academic disciplines. Today's big-data projects require that life science researchers either learn how to use command-line tools or outsource their data analysis to computational experts. Active engagement in analyzing data generated by life science researchers is essential to advancing

interdisciplinary research in the biological sciences. However, biologists and medical researchers often lack formal training in the use of computational techniques. Scholars have developed teaching models aimed to support biomedical researchers transition from using a graphical interface (e.g., Microsoft Excel) to UNIX command line[3]. Biomedical researchers often possess various levels of computational skills; workshops require adjusting the teaching pace for trainees who have different levels of computational background. Additionally, the ability of trainees to switch from a graphical interface with the assistance of a mouse to a command line interface with no mouse support varies across the life science researchers.

Training a life science research to use computational techniques poses unique challenges and requires a special approach. For example, such training does not require cultivating a deep understanding of fundamental computer science principles. Instead, such training is limited to applied learning, requiring quick assimilation of introduced techniques and acquired skills within the context of the research project of interest. In general, flexible pedagogy is preferred. For example, instead of introducing the fundamental computational concepts formally, such concepts are introduced on an as-needed basis, are combined with numerous hands-on examples, and rely on the instructor's guidance to consolidate the learner's newly-acquired knowledge.
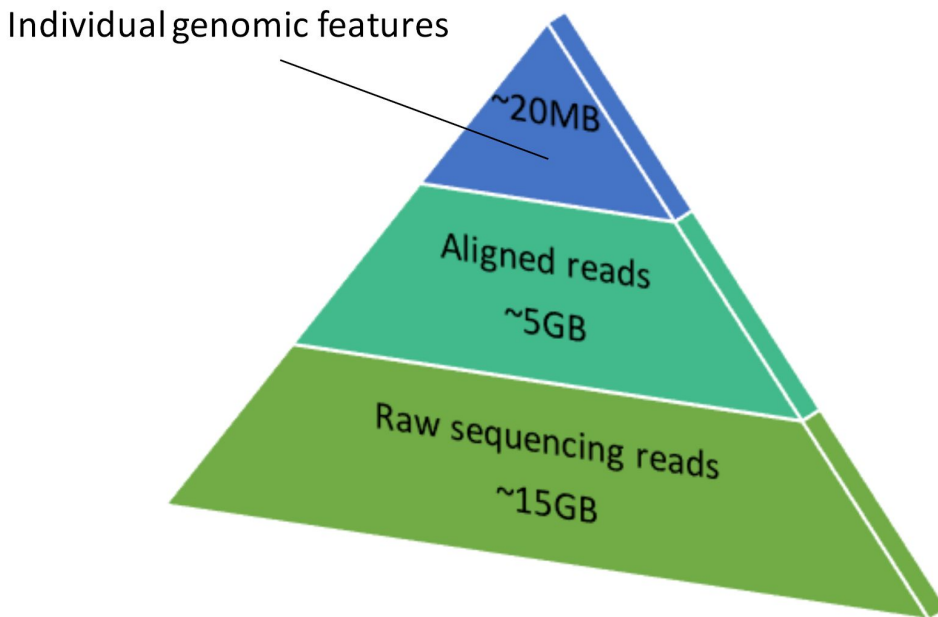
Biomedical researchers who have completed the training in computational skills often continue to engage with the command line and with training instructors long after the training session has formally ended. Researchers who have completed the workshops are also able to

better engage with computer scientists when seeking consultation services for their projects. The training model I developed [3] has helped life science researchers to learn the "language of computing," a skill that allows them to better understand what analyses can be accomplished with UNIX, and how to ask for specific types of help from computational scientists.

Alternatively, biomedical researchers can delegate large-scale data analyses to bioinformatics cores. However, outsourcing analyses present several challenges. Many complex issues arise during the analysis of big omics data which are difficult to predict in advance. In such cases, life science and biomedical research can optimize the analysis if they are adequately trained and remain involved in the analytical workflow. In addition, research groups utilizing the core often want to move the project in different directions from what was originally proposed. Another approach is to develop a GUI that allows researchers with a limited computational background to easily create, run, and troubleshoot analytical pipelines. While useful to researchers with a limited computational background, these interfaces may have limited computational capacity compared to high-performance clusters and might not suitable for analysis of big omics data generated from some large-scale clinical cohorts.

**Computational resources required to analyze the big omics data**

Computational resources required to analyze the big omics data differs significantly depending on the step of the analysis. Analyses dealing with raw data typically require a significant amount of computational resources to perform essential tasks and space to store the output data. Analyzing big omics data can typically be performed only on high-performance clusters. For example, the analysis of differentially expressed genes based on RNA-Seq data starts with raw reads produced by the sequencing machines and concludes with a list of differentially expressed genes with corresponding gene expression fold changes. The computational resources and amount of storage needed to analyze and store the results of analysis significantly various across various step of the analysis (**Figure 2**). As a result, the steps of the analysis requiring significant computational resources and space to store the data can be only performed on high-performance clusters. Other steps require a smaller amount of resources and thus can perform locally on the personal computer or laptop. Typically, the analysis performed on the local machine does not require the knowledge of command line skills. Such analyses typically involve various statistical analyses and visualization steps, and these tasks can be performed using the widely popular statistical language $R$[6]. For example, once the gene expression levels were obtained from RNA-Seq data on the high-performance cluster, one can transfer them to a personal computer and locally perform differential expression analysis using available $R$ packages[7].

**Figure 2.** The amount of space needed to store the results of the sequencing analysis.

**Computational reproducibility in life science research**

An astonishing number of bioinformatics software tools are designed to accommodate increasingly bigger, more complex, and more specialized bio-datasets are developed each year[4]. With the increasing importance and popularity of computational and data-enabled approaches among biomedical researchers, it becomes ever more critical to ensure that the developed software is usable[8] and the Uniform Resource Locator (URL), through which the software tool is accessible, is archivally stable. Consistently usable and accessible software provides a foundation for the reproducibility of published biomedical research, defined as the ability to

replicate published findings by running the same computational tool on data generated by the study[8,9].

Open data, open software, and reproducible research are important aspects of big data analysis in the life sciences[10]. Reproducing previously published results can be made possible, by releasing all research objects, such as raw data, and publically available, archivally stable, and installable computer code. However, a lack of strict implementation or enforcement of journal policies for resource sharing harms rigor and reproducibility as some authors refuse to share the data[11] or source code.

Despite these challenges, consistently usable and accessible software provides a necessary foundation for rigorous and reproducible data-intensive biomedical research[11,12]. In addition, the usability — or, 'user-friendliness' — of software tools is important, and it can affect its scientific utility. Currently, an estimated 74% of computational software resources are accessible through URLs published in the original paper[13]. Many developed tools are difficult to install and some are impossible to install[13]. Kumar and Dudley warn that poorly maintained or implemented tools will hinder progress in "big data"-driven fields[14].

any journals now require that omics data generated by the published study should be shared when the paper is released, an important step forward toward improving computational reproducibility in our field. However, the bioinformatics community still lacks the comprehensive policies on precisely how openly shared code used to perform the analysis and

generate the figures should be. In a promising effort from *eLife* journal, editors suggest that *R* code used to generate figures should be shared together with the figure[15].

**Discussion**

High-throughput technologies have changed the landscape of training, research, and education in biomedical fields[16]. Big data generated by those technologies across large-scale clinical cohorts can potentially enable a researcher to connect complex diseases to relevant omics features. As our knowledge of scientifically-validated disease-trait matches increase, new opportunities emerge for development of novels diagnostic and therapeutic tools. However, analysis of big data requires the use of sophisticated bioinformatics algorithms which are often packaged as command-line-driven software tools. A researcher who wishes to use such tools must acquire specific computational skillsets, which are not included in the traditional life science curriculum at major Universities. With the increasing size and complexity of big omics datasets in the biological and medical sciences, researchers are facing a growing dilemma of devoting the time to acquire key computational skills or outsource the analyses to computational researchers.

At present, biomedical researchers are not involved in computational training on a large-scale worldwide. In this review, I provide evidence that the computational training model - that is, when life science research groups receive training and resources to analyze the data that they generate, is a more sustainable approach. The computational training model of life science researchers, when successfully applied across many research institutions worldwide,

has the potential to change the landscape of contemporary biomedical research, training, and education. If a critical mass of biomedical researchers obtains computational skills sufficient for analyzing big data, computational training will more likely become an integral part of analysis curricula at these institutions.

The computational training model offers benefits for both individual researchers and the scientific community. Life science and biomedical researchers gain a competitive skill when learning to conduct analysis in a command-line setting. Today's omics data generates file sizes too large to be opened on a personal computer. These novice computational researchers often must perform their analyses on a high-performance cluster with command line tools and, in the process, the researchers become familiar with programming and basic system administration tasks. Such valuable skills could be leveraged to further the researcher's projects and career.

One important outcome of a comprehensively implemented computational training model is to improve reproducibility in the big data-driven fields of life science and medical research. The standard for rigorous and reproducible analysis is an emerging topic with multiple initiatives across research groups. The scientific community has identified current challenges to ensuring reproducibility of interpreting and integrating big data analysis in the life sciences[11]. For example, *eLife* journal raised the bar of reproducibility, challenging the traditional static representation of data and results of the analysis (usually in the form of PDF or HTML formats). Instead, eLife now suggests a code-based publication, which enables data and analysis to be fully reproducible by the reader[17,18].

## References

1.  Cancer Genome Atlas Research Network *et al.* The Cancer Genome Atlas Pan-Cancer analysis project. *Nat. Genet.* **45**, 1113–1120 (2013).

2.  Thorsson, V. *et al.* The Immune Landscape of Cancer. *Immunity* **48**, 812–830.e14 (2018).

3.  Mangul, S., Martin, L. S., Hoffmann, A., Pellegrini, M. & Eskin, E. Addressing the Digital Divide in Contemporary Biology: Lessons from Teaching UNIX. *Trends Biotechnol.* **35**, 901–903 (2017).

4.  Wren, J. D. Bioinformatics programs are 31-fold over-represented among the highest impact scientific papers of the past two decades. *Bioinformatics* **32**, 2686–2691 (2016).

5.  Altschul, S. *et al.* The anatomy of successful computational biology software. *Nat. Biotechnol.* **31**, 894–897 (2013).

6.  Cornillon, P.-A. *et al. R for Statistics*. (CRC Press, 2012).

7.  Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).

8.  List, M., Ebert, P. & Albrecht, F. Ten Simple Rules for Developing Usable Software in Computational Biology. *PLoS Comput. Biol.* **13**, e1005265 (2017).

9.  Beaulieu-Jones, B. K. & Greene, C. S. Reproducibility of computational workflows is automated using continuous analysis. *Nature Biotechnology* **35**, 342–346 (2017).

10. Murphy, F. Open access, open data, FAIR Data and their implications for life sciences researchers. *Emerging Topics in Life Sciences* **2**, 759–762 (2018).

11. Stodden, V., Seiler, J. & Ma, Z. An empirical analysis of journal policy effectiveness for computational reproducibility. *Proc. Natl. Acad. Sci. U. S. A.* **115**, 2584–2589 (2018).

12. Beaulieu-Jones, B. K. & Greene, C. S. Reproducibility of computational workflows is automated using continuous analysis. *Nat. Biotechnol.* **35**, 342–346 (2017).

13. Mangul, S. *et al.* A comprehensive analysis of the usability and archival stability of omics computational tools and resources. (2018). doi:10.1101/452532

14. Kumar, S. & Dudley, J. Bioinformatics software for biologists in the genomics era. *Bioinformatics* **23**, 1713–1717 (2007).

15. Bauer, P. C. Writing a Reproducible Paper in R Markdown. *SSRN Electronic Journal* doi:10.2139/ssrn.3175518

16. Hayden, E. C. Genome researchers raise alarm over big data. *Nature* (2015). doi:10.1038/nature.2015.17912

17. Perkel, J. M. Pioneering 'live-code' article allows scientists to play with each other's results. *Nature* (2019). doi:10.1038/d41586-019-00724-7

18. Burgess, S. Reproducible research article collection. (2018). doi:10.14293/s2199-1006.1.sor-uncat.clsuuhc.v1