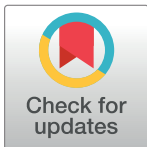


COMMUNITY PAGE

Enabling precision medicine via standard communication of HTS provenance, analysis, and results

Gil Alterovitz^{1,2,3}, Dennis Dean⁴, Carole Goble⁵, Michael R. Crusoe⁶, Stian Soiland-Reyes⁵, Amanda Bell⁷, Anais Hayes⁷, Anita Suresh⁸, Anjan Purkayastha^{7,9}, Charles H. King^{7,10}, Dan Taylor¹¹, Elaine Johanson¹², Elaine E. Thompson¹², Eric Donaldson¹², Hiroki Morizono^{13,14}, Hsinyi Tsang^{15,16}, Jeet K. Vora⁷, Jeremy Goecks¹⁷, Jianchao Yao¹⁸, Jonas S. Almeida¹⁹, Jonathon Keeney⁷, KanakaDurga Addepalli¹⁶, Konstantinos Krampis^{20,21}, Krista M. Smith¹⁰, Lydia Guo²², Mark Walderhaug¹⁰, Marco Schito²³, Matthew Ezewudo²³, Nuria Guimera²⁴, Paul Walsh²⁵, Robel Kahsay⁷, Srikanth Gottipati²⁶, Timothy C. Rodwell⁸, Toby Bloom²⁷, Yuching Lai²⁴, Vahan Simonyan¹⁰, Raja Mazumder^{7,10*}



OPEN ACCESS

Citation: Alterovitz G, Dean D, Goble C, Crusoe MR, Soiland-Reyes S, Bell A, et al. (2018) Enabling precision medicine via standard communication of HTS provenance, analysis, and results. *PLoS Biol* 16(12): e3000099. <https://doi.org/10.1371/journal.pbio.3000099>

Published: December 31, 2018

Copyright: This is an open access article, free of all copyright, and may be freely reproduced, distributed, transmitted, modified, built upon, or otherwise used by anyone for any lawful purpose. The work is made available under the [Creative Commons CC0](https://creativecommons.org/licenses/by/4.0/) public domain dedication.

Funding: US Food and Drug Administration had supported this study (grant number FDA - BAA-12-00118) (website: https://www.fbo.gov/index?s=opportunity&mode=form&tab=core&id=5c9bc2459a39b01a39c6b6f6b8cda930&_cvview=0). BioExcel CoE, a project funded by the European Commission Horizon 2020 Framework Programme under contract H2020-EINFRA-2015-1-675728. Received by CG and SSR. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

Abbreviations: API, application programming interface; BCO, BioCompute Object; BD2K, Big Data

1 Harvard/MIT Division of Health Sciences and Technology, Harvard Medical School, Boston, Massachusetts, United States of America, **2** Computational Health Informatics Program, Boston Children’s Hospital, Boston, Massachusetts, United States of America, **3** Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Boston, Massachusetts, United States of America, **4** Seven Bridges, Cambridge, Massachusetts, United States of America, **5** School of Computer Science, The University of Manchester, Manchester, United Kingdom, **6** Common Workflow Language Project, Vilnius, Lithuania, **7** The Department of Biochemistry & Molecular Medicine, The George Washington University Medical Center, Washington, DC, United States of America, **8** Foundation for Innovative New Diagnostics (FIND), Geneva, Switzerland, **9** OpenBox Bio, Vienna, Virginia, United States of America, **10** The McCormick Genomic and Proteomic Center, The George Washington University, Washington, DC, United States of America, **11** Internet 2, Washington, DC, United States of America, **12** US Food and Drug Administration, Silver Spring, Maryland, United States of America, **13** Center for Genetic Medicine, Children’s National Medical Center, Washington, DC, United States of America, **14** The Department of Genomics and Precision Medicine, The George Washington University, School of Medicine and Health Sciences, Washington, DC, United States of America, **15** Center for Biomedical Informatics and Information Technology, National Cancer Institute, National Institutes of Health, Gaithersburg, Maryland, United States of America, **16** Attain, McClean, Virginia, United States of America, **17** Computational Biology Program, Oregon Health & Science University, Portland Oregon, United States of America, **18** MRL IT, Merck & Co., Boston, Massachusetts, United States of America, **19** Stony Brook University, School of Medicine and College of Engineering and Applied Sciences, Stony Brook, New York, United States of America, **20** Department of Biological Sciences, Hunter College of The City University of New York, New York, New York, United States of America, **21** Institute for Computational Biomedicine, Weill Cornell Medical College, New York, New York, United States of America, **22** Wellesley College, Wellesley, Massachusetts, United States of America, **23** Critical Path Institute, Tucson, Arizona, United States of America, **24** DDL Diagnostic Laboratory, Rijswijk, Netherlands, **25** NSilico Life Science, Nova Center, Belfield Innovation Park, University College Dublin, Dublin Ireland, **26** OTSUKA Pharmaceutical Development & Commercialization, Princeton, New Jersey, United States of America, **27** New York Genome Center, New York, New York, United States of America

* mazumder@gwu.edu

Abstract

A personalized approach based on a patient’s or pathogen’s unique genomic sequence is the foundation of precision medicine. Genomic findings must be robust and reproducible, and experimental data capture should adhere to findable, accessible, interoperable, and reusable (FAIR) guiding principles. Moreover, effective precision medicine requires standardized reporting that extends beyond wet-lab procedures to computational methods. The

to Knowledge; CGI, computer graphic imaging; CWL, Common Workflow Language; dbGaP, Database of Genotypes and Phenotypes; EHR, electronic health record; EMBL-EBI, European Molecular Biology Laboratory-European Bioinformatics Institute; Env, environmental; FAIR, findable, accessible, interoperable, and reusable; FHIR, Fast Healthcare Interoperability Research; FDA, Food and Drug Administration; GA4GH, Global Alliance for Genomics and Health; GWAS, genome-wide association studies; HIVE, High-performance Integrated Virtual Environment; HPC, high-performance computing; HL7, Health Level 7; HTS, high-throughput sequencing; IEEE, Institute of Electrical and Electronics Engineers; IO, input/output; ISO, International Organization for Standardization; JSON, JavaScript Object Notation; LD Hub, Linkage Disequilibrium Hub; NCBI, National Center for Biotechnology Information; NCI, The National Cancer Institute; NGS, Next-Generation Sequencing; NIH, National Institutes of Health; ORCID, Open Researcher and Contributor ID; OS, operating system; PAV, Provenance, Authoring and Versioning ontology; Prereq, prerequisite; PROV, provenance specification; PROV-O, PROV ontology; ReSeqTB-UVP, The Relational Sequencing Tuberculosis-Unified Variant Pipeline; RO, research object; TCGA, The Cancer Genome Atlas; URI, uniform resource identifier; UVP, Unified Variant Pipeline; W3C, World Wide Web Consortium; WGS, whole genome sequencing; Xref, external reference.

BioCompute framework (<https://w3id.org/biocompute/1.3.0>) enables standardized reporting of genomic sequence data provenance, including provenance domain, usability domain, execution domain, verification kit, and error domain. This framework facilitates communication and promotes interoperability. Bioinformatics computation instances that employ the BioCompute framework are easily relayed, repeated if needed, and compared by scientists, regulators, test developers, and clinicians. Easing the burden of performing the aforementioned tasks greatly extends the range of practical application. Large clinical trials, precision medicine, and regulatory submissions require a set of agreed upon standards that ensures efficient communication and documentation of genomic analyses. The BioCompute paradigm and the resulting BioCompute Objects (BCOs) offer that standard and are freely accessible as a GitHub organization (<https://github.com/biocompute-objects>) following the “Open-Stand.org principles for collaborative open standards development.” With high-throughput sequencing (HTS) studies communicated using a BCO, regulatory agencies (e.g., Food and Drug Administration [FDA]), diagnostic test developers, researchers, and clinicians can expand collaboration to drive innovation in precision medicine, potentially decreasing the time and cost associated with next-generation sequencing workflow exchange, reporting, and regulatory reviews.

Introduction

Precision medicine requires the seamless production and consumption of genomic information. The National Center for Biotechnology Information’s (NCBI) Database of Genotypes and Phenotypes (dbGaP) [1] and ClinVar [2] illustrate the benefits of genomic data sharing structures such as genome-wide association studies (GWAS). Linkage Disequilibrium Hub (LD Hub), a centralized database of GWAS results for diseases and/or traits [3], is another example of success. Although the importance of data sharing is established, recording, reporting, and sharing of analysis protocols are often overlooked. Standardized genomic data generation empowers clinicians, researchers, and regulatory agencies to evaluate the reliability of biomarkers generated from complex analyses. Trustworthy results are increasingly critical as genomics play a larger role in clinical practice. In addition, fragmented approaches to reporting impede the advancement of genomic data analysis techniques.

The price of high-throughput sequencing (HTS) decreased from US\$20 per base in 1990 to less than US\$0.01 per base in 2011 [4]. Lower costs and greater accessibility resulted in a proliferation of data and corresponding analyses that in turn advanced the field of bioinformatics. Novel drug development and precision medicine research stand to benefit from innovative, reliable, and accurate -omics-based (i.e., genomics, transcriptomics, proteomics) investigation [5]. However, the availability of HTS has outpaced existing practices for reporting on the protocols used in data analysis.

Fast Healthcare Interoperability Resources (FHIR) [6,7] and the Global Alliance for Genomics and Health (GA4GH) [8] capture and communicate genomic information within specific community domains. The Common Workflow Language (CWL) [9] and research objects (ROs) [10] capture repeatable and reproducible workflows in a domain agnostic manner. The BioCompute framework (<https://w3id.org/biocompute/1.3.0>) combines these standards via a BioCompute Object (BCO) [11] to report the provenance of genomic sequencing data in the context of regulatory review and research. A BCO is designed to satisfy findable, accessible,

interoperable, and reusable (FAIR) data principles [12], ensuring that data and pipelines are available for evaluation, validation, and verification [11,13,14]. The BCO also meets the National Institutes of Health (NIH) strategic plan for data science [15], which states that the quality of clinical data should be maintained at all stages of the research cycle, from generation through the entire analysis process. These characteristics ensure that the BioCompute framework is applicable in any context in which scientists are required to report on data provenance, including large clinical trials or the development of a knowledge base. In the following text, we describe how the BioCompute framework (see Fig 1) leverages and harmonizes FHIR, GA4GH, CWL, and RO to create a unified standard for the collection and reporting of genomic data.

Background

At a recent Academy of Medical Sciences symposium on preclinical biomedical research, participants identified several measures for improving reproducibility. These include greater openness and transparency, defined reporting guidelines, and better utilization of standards and quality control measures [16]. Compromised reproducibility dissipates resources and hinders progress in the life sciences, as highlighted by several other publications [17,18].

Researchers utilize two types of reproducibility—method reproducibility and result reproducibility [19]. The first type, also called repeatability, is defined by providing detailed experimental methods so others can repeat the procedure exactly. The second type, also called replicability, is defined by achieving the same results as the original experiment by closely adhering to the methods. Method reproducibility depends on a comprehensive description of research procedures. CWL achieves this aim for computational methods by encouraging scientists to adhere to a common language. This facilitates better methods for identification of errors and locating deviances. ROs for workflows are built on this concept by using metadata manifests that describe the experimental context, including packaging of the method, provenance logs, and the associated codes and data [20].

Universally reproducible data is an aspirational goal, but challenges remain. Without widely adopted repeatability and analytical standards for Next-Generation Sequencing (NGS)/HTS studies, regulatory agencies, researchers, and industry cannot effectively collaborate to validate results and drive the emergence of new fields [21]. Irreproducible results can cause major delays and be a substantial expenditure for applicants submitting work for regulatory review and may even be viewed as not trustworthy by third-party verification groups or regulatory agencies. A large number of workflow management systems and bioinformatics platforms have been developed to overcome this barrier, each with their own unique method to record computational workflows, pipelines, versions, and parameters, but these efforts remain disorganized and require harmonization to be fully effective. BioCompute enables clear communication of “what was done” and “why it was done” by tracking provenance and documenting processes in a standard format irrespective of the platform or the programming language or even the tool used.

Provenance of data

In order to be reproducible, the origin and history of research data must be maintained. Provenance refers to a datum’s history starting from the original source, namely, its lineage. A lineage graph can show the source of a datum in a database, data movement between databases or computational processes, or data generated from a computational process. Complementary to data lineage is a process audit. This is a historical trail that provides snapshots of intermediary states, values for configurations and parameters, and traceability of stepwise analytical

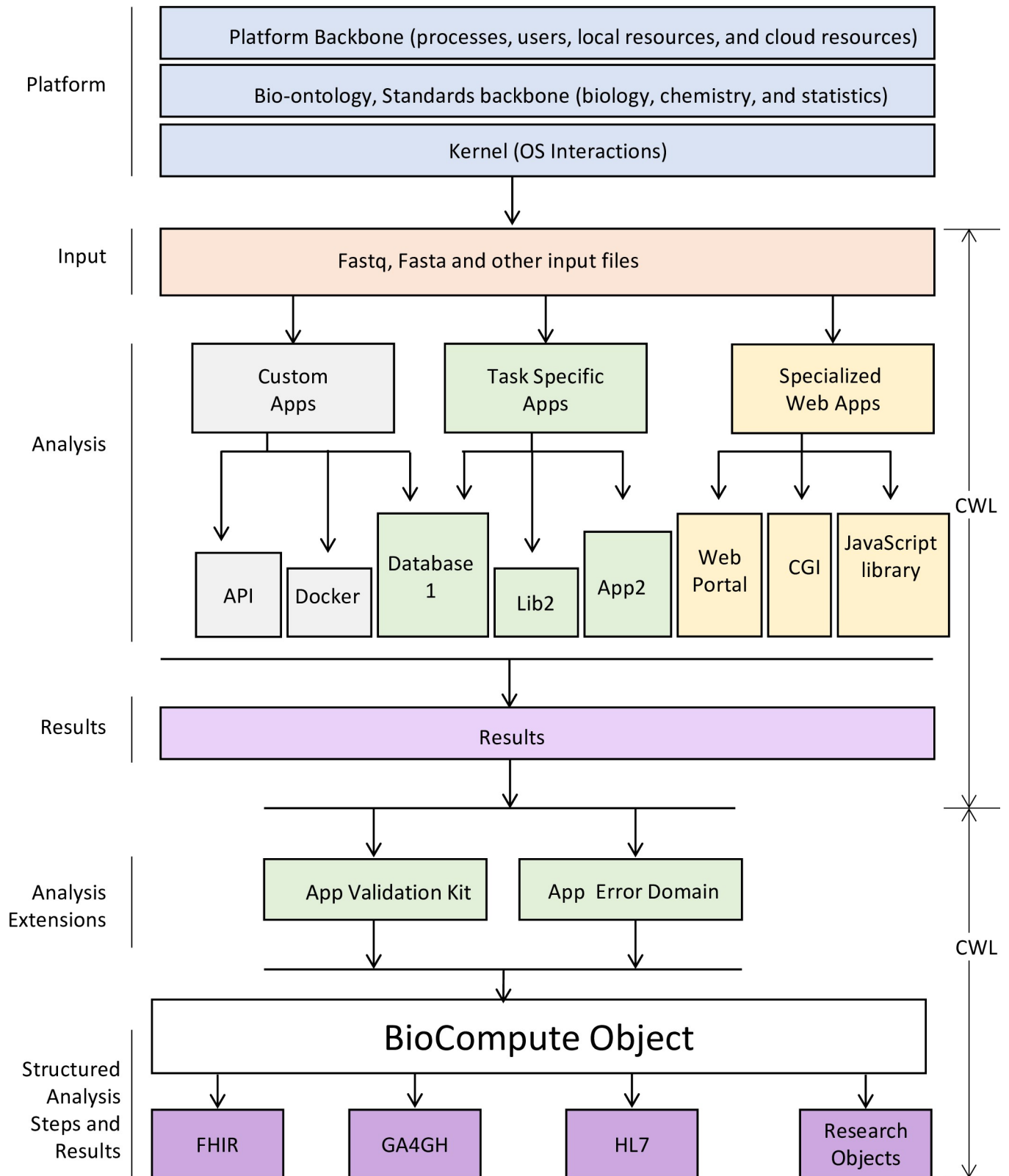


Fig 1. Schematic of BCO as a framework for advancing regulatory science by incorporating existing standards and introducing additional concepts that include digital signature, usability domain, validation kit, and error domain. API, application programming interface; app, application; BCO, BioCompute Object; CGI, computer graphic imaging; FHIR, Fast Healthcare Interoperability Research; GA4GH, Global Alliance for Genomics and Health; HL7, Health Level 7; OS, operating system.

<https://doi.org/10.1371/journal.pbio.3000099.g001>

processing [22]. Such audit trails enable an independent reviewer to effectively evaluate a computational investigation. Both types of records gather critical provenance information to ensure accuracy and validity of experimental results which can be presented as is or collected into databases. Modern computational workflows produce large amounts of fine-grained but useless trace records, whereas modern web developments facilitate easy data transformation and copying. Combing and organizing the large volume of resulting material is a daunting challenge. In the molecular biology field alone, there are hundreds of public databases that curate, filter and annotate data to make them more useful. Only a handful of these retain the “source” data; the remainder consist of secondary views of the source data or views of other publications’ views [23]. Databases are challenged to accurately collect lineage and process records while also maintaining granularity and “black-box” steps [24,25].

Provenance tracing issues have far-reaching effects on scientific work. Advancement depends on confidence in each of the following—accuracy and validity of the data, process used, and knowledge generated by research. Establishing trust is especially difficult when reporting a complex, multistep process involving aggregation, modeling, and analysis [26]. Computational investigations require collaboration with adjacent and disparate fields to effectively analyze a large volume of information. Effective collaboration requires a solution beyond open data to establish open science. Provenance must be preserved and reported to promote transparency and reproducibility in complex analyses [27]. Standards must be established to reliably communicate genomic data between databases and individual scientists.

An active community has engaged in provenance standardization to achieve these aims [28], culminating in the World Wide Web Consortium (W3C) provenance specification (PROV) [29]. PROV Ontology (PROV-O) is used by FHIR and ROs and is based on the concept of generating an entity target via an agent’s activity (see Fig 2). Workflow management systems capture analytic processes, while bioinformatic platforms capture analysis details. In combination, they provide a record of data provenance. Few, if any, systems and platforms offer a consistent method to accurately capture all facets of the various roles assumed by an agent who manipulates digital artifacts. BCO encourages adoption of standards such as PROV-O (<https://www.w3.org/TR/prov-o/>) and Open Researcher and Contributor ID (ORCID; <https://orcid.org/>) by defining how to recreate a complete history of what was computed, how it was computed, by whom it was computed, and why it was computed. Also known as the provenance domain, this section of BCOs incorporate the Provenance, Authoring and Versioning ontology (PAV, namespace <http://purl.org/pav/>) to capture “just enough” information to track how data are authored, curated, retrieved, and processed among many specific designations of an “agent.”

Key considerations for communication of provenance, analysis, and results

Workflow management systems

Scientific workflows have emerged as a model for representing and managing complex scientific computations [26]. Each step in a workflow specifies a process or computation to be executed, linked with other steps by the data flow and dependencies. In addition, workflows describe the mechanisms to carry out the steps in a distributed computing environment [26,30]. Documentation of both analytic processes and provenance information is essential for useful reporting [31].

Workflow management systems coordinate the sequential components in an analysis pipeline [30]. They also enable researchers to generate pipelines that can be executed locally on institutional servers and remotely on the cloud [32]. Cloud infrastructure, high-performance

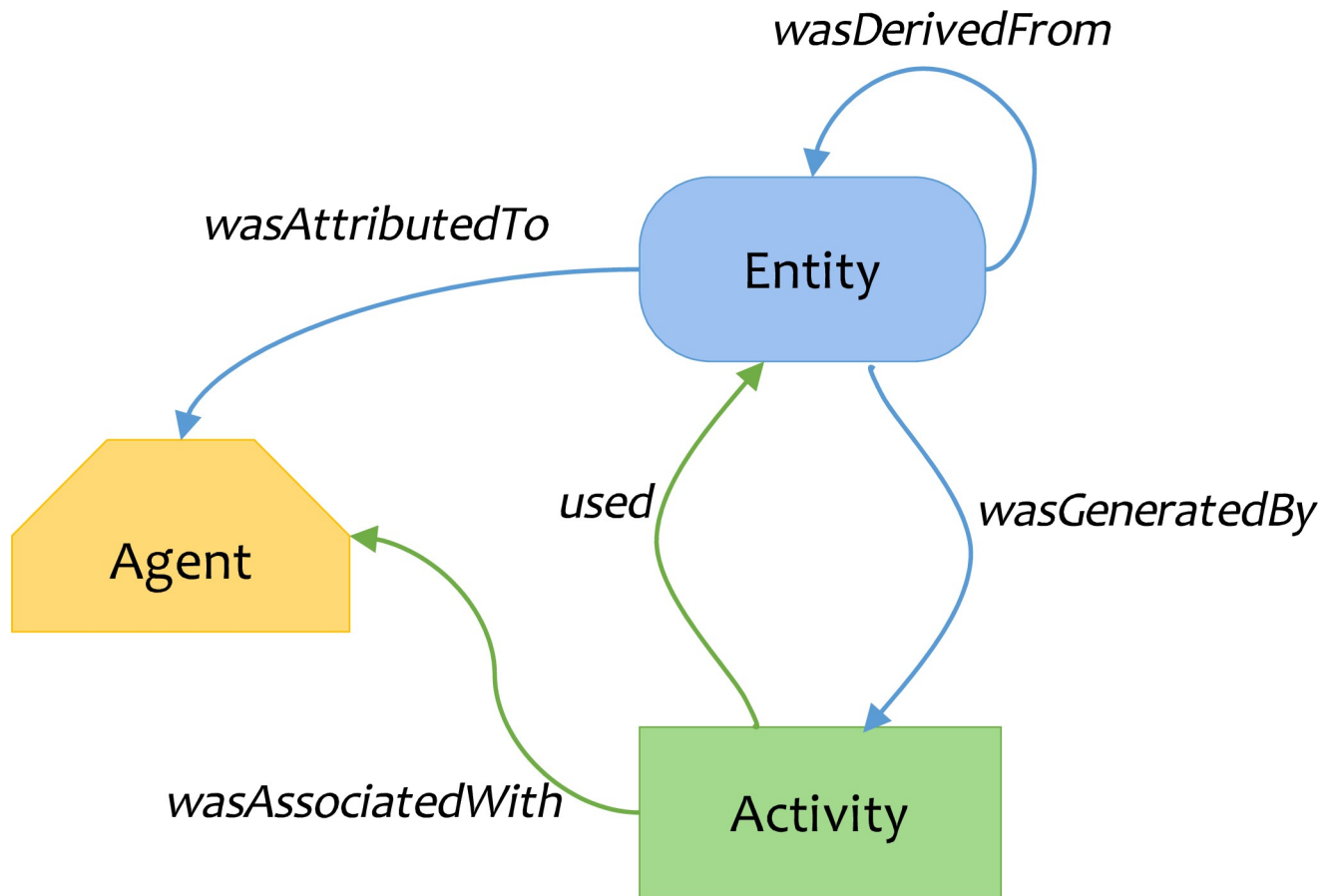


Fig 2. W3C PROV data model overview, used in Fast Healthcare Interoperability Research (FHIR) and research object (RO). Adapted from <http://www.w3.org/TR/prov-primer/>.

<https://doi.org/10.1371/journal.pbio.3000099.g002>

computing (HPC) systems, and Big Data cluster-computation frameworks enhance data reproducibility and portability (see [S1 Text](#)). Workflow-centric ROs with executable components are an extension of these management systems [5,20]. Extensive reviews of workflow systems currently in use for bioinformatics have already been published [32–35], and we are not recommending any one over the others. Currently, workflow management systems capture provenance information but rarely in the PROV standard (<https://www.w3.org/TR/prov-overview/>). Therefore, BCOs rely on existing standards that are themselves based on PROV standards, like CWL to manage pipeline details, and ROs and FHIR to unify and enhance interoperability.

Bioinformatics platforms

HTS technology is increasingly relevant in the clinical setting, with a growing need to store, access, and compute more sequencing reads and other biomedical data [36]. The increase in computational requirements has directed the scientists in this space to a call for standard usage methods on integrated computing infrastructure, including storage and computational nodes. This kind of standardization will minimize transfer costs and remove the bottlenecks found in both downstream analyses and community communication of computational analysis results [37]. For bioinformatics platforms, communication requirements include (a) recording all

analysis details such as parameters and input data sets and (b) sharing analysis details so that others can understand and reproduce analyses.

HTC environments deliver large amounts of processing capacity over long periods of time. These are ideal environments for long-term computation projects, such as those performed for genomic research [38]. Most HTC platforms utilize distributed cloud-computing environments to support extra-large data set storage and computation, as well as hosting tools and workflows for many biological analyses. Cloud-based infrastructures also reduce the “data silo” phenomenon by converting data into reproducible formats that facilitate communication (see [S1 Text](#)). Additionally, the National Cancer Institute (NCI) has initiated the Cloud Pilots project in order to test a distributed computing approach for the multilevel, large-scale data sets available on The Cancer Genome Atlas (TCGA) [39]. Several of the high-throughput [26], cloud-based platforms that have been developed, including High-performance Integrated Virtual Environment (HIVE) [37,40] and Galaxy [41]—along with commercial platforms from companies like DNAnexus (dnanexus.com) and Seven Bridges Genomics (sevenbridges.com)—have participated in the development of BioCompute. This participation ensures that while using these bioinformatics platforms, users would not need to keep track of all of the information needed to create a BCO. Such information will be automatically or semiautomatically collected during the creation and running of a workflow.

The genomic community has come to acknowledge the necessity of data sharing and communication to facilitate reproducibility and standardization [42,43]. Data sharing is crucial in everything from long-term clinical treatments to public-health emergency response [44]. Extending bioinformatics platforms to include data provenance, standard workflow computation, and encoding results with available standards through BCO implementation will greatly support the exchange of genomic data analysis methods for regulatory review.

Regulatory supporting standards

Assessment of data submitted in a regulatory application requires clear communication of data provenance, computational workflows, and traceability. A regulatory reviewer must be able to verify that sequencing was done appropriately and that pipelines and parameters were applied correctly, and assess the final results. They must have the tools to critically evaluate the validity of results such as allelic difference or variant call. Because of these requirements, review of any clinical trial or any submission supported with HTS results requires considerable time and expertise. Inclusion of a BCO with a regulatory submission would help to ensure that data provenance is unambiguous and that the bioinformatics workflow is fully documented [11,15,23,45,46].

To truly understand and compare computational tests, a standard method (like BCO) requires tools to capture progress and to communicate the workflow and input/output data. As the regulatory field progresses, methods have been developed and are being continually refined to capture workflows and exchange data electronically [26]. See [Fig 3](#) for BioCompute extensions to HTS analysis that support data provenance and reproducibility.

BCOs and their harmonizing efforts

The BioCompute paradigm and BCOs were conceptualized to capture the specific details of HTS computational analyses. The primary objectives of a BCO is to (a) harmonize and communicate HTS data and computational results as well as (b) encourage interoperability and the reproducibility of bioinformatics protocols. Harmonizing HTS computational analyses is especially applicable to clarify the results from clinical trials and other genomic related data for regulatory submissions.

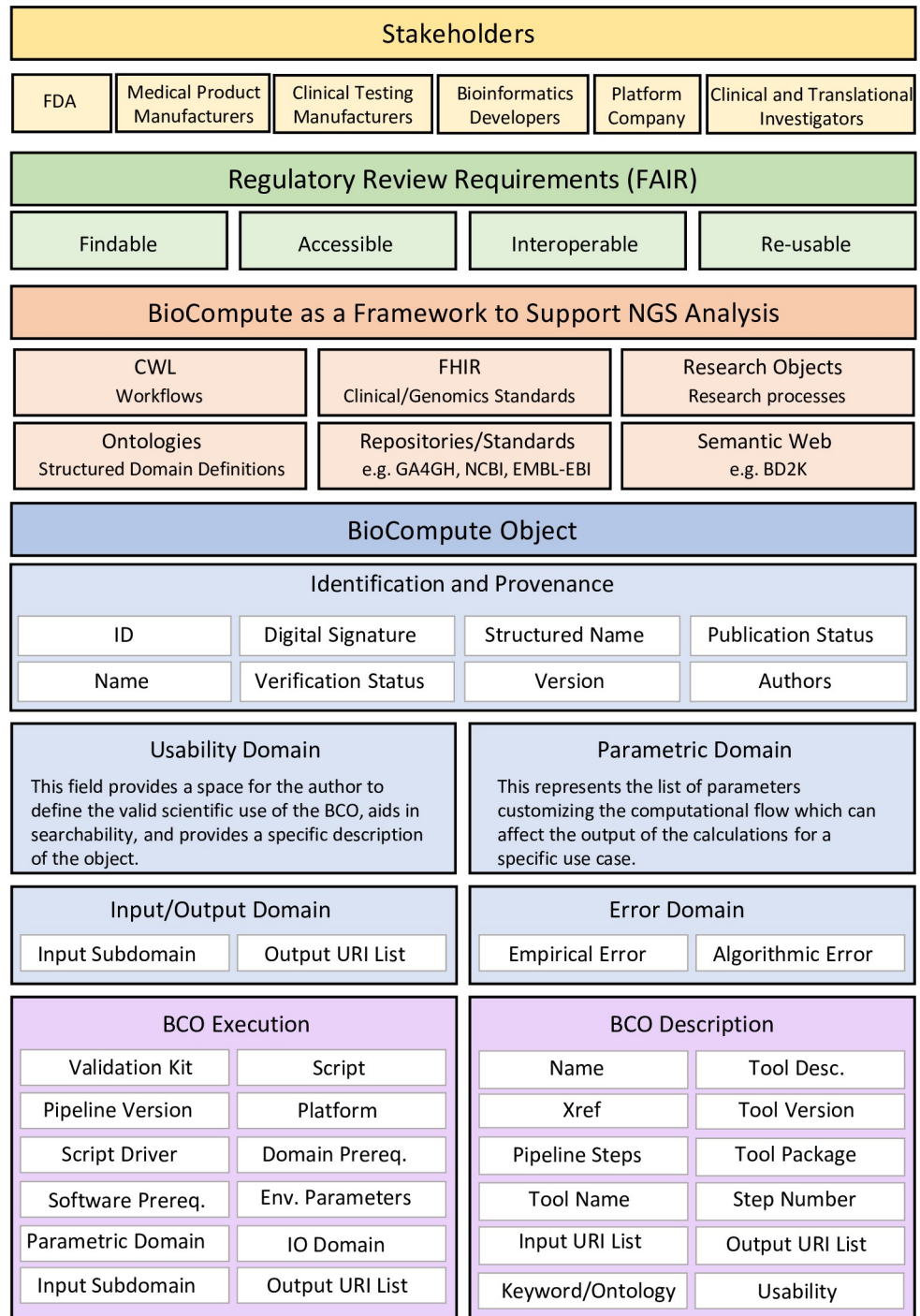


Fig 3. Generic HTS platform schematic with proposed BCO integrations and extensions. BCO, BioCompute Object; BD2K, Big Data to Knowledge; Desc., description; EMBL-EBI, European Molecular Biology Laboratory-European Bioinformatics Institute; Env., environmental; FDA, Food and Drug Administration; FHIR, Fast Healthcare Interoperability Research; GA4GH, Global Alliance for Genomics and Health; ID, identification; IO, input/output; NCBI, National Center for Biotechnology Information; NGS, Next-Generation Sequencing; Prereq., prerequisite; PROV, provenance specification; RO, research object; URI, uniform resource identifier; W3C, World Wide Web Consortium; Xref, external reference.

<https://doi.org/10.1371/journal.pbio.3000099.g003>

The BioCompute paradigm is novel in its combination of existing standards with the methodologies and tools to evaluate an experiment both programmatically and empirically. The BCO takes a snapshot of an experiment's entire computational procedure adhering to FAIR data guidelines. A BCO is findable and publicly accessible through the BCO portal, interoperable by maintaining the computational context and data provenance, which also makes the computational experiment reusable. Using this snapshot of the experiment, which can include the range of acceptable experimental results in the verification kit, allows any other user to run the exact experiment and produce the same results. Additionally, through the use of provenance and usability domains, a reviewer can quickly decide whether the underlying scientific principles merit approval or whether further review or reanalysis is required.

BCO specification

The BCO specification provides details about the BCO structure. BCOs are represented in JavaScript Object Notation (JSON) formatted text. The JSON format was chosen because it is both human and machine readable as well as easy to write and store. Top-level BCO fields include BCO identifier, type, digital signature, and specification version. The rest of the information is organized into several domains. Below is a brief summary of the type of information present in the various domains. For additional details, please see the latest version of the BCO specification (<https://w3id.org/biocompute/releases>).

Provenance domain includes fields such as structured name, version, inheritance, status, contributors, license, and creation and modification dates. Usability domain provides a clear description of the intended use of the BCO. Extension domain allows users to add content related to other standards that enhances interoperability. Description domain provides details such as keywords, external references, and human-readable descriptions and sequence of pipeline steps. The execution domain is machine-readable information that can be used to run the entire pipeline. Parametric domain provides information on parameters that were changed from default, and the input/output domain provides links to input and output files. The error domain includes details that can be used to verify whether a particular BCO has been used as intended and that any errors are within the acceptable range. Error domain along with verification kit (such as *in silico* generated read files with known mutations) of the BCO allows verification of a workflow in different bioinformatics platforms, for example.

BCO implementation

The BCO stores information relating to every package and every script—including nuanced data that are often not reported, like version number—in a human- and machine-readable format. A typical workflow analyzing HTS data may have a dozen iterations as different software, packages, and scripts are used—each with their own parameters. Through trial and error, various analyses and parameters are refined, which may yield new insights. If new insights are discovered, a snapshot of the state that produced the result in question can be stored as a BCO. To further increase the probability of successful replication, the BCO can even be verified before it is sent out for replication. In this step, for example, a third-party can verify that all input parameters across the entire analysis pipeline will generate robust and reliable output within the usability domain described in the BCO.

The BCO is therefore far more efficient than any existing means of communicating HTS analysis information. For example, a researcher has discovered a specific variant that corresponds to a population being studied in her own country, and she is interested in learning the relevance of this variant in other populations. By using a BCO, her entire analysis can be quickly understood and repeated by her international colleagues, working with their own data

sets, with high confidence in the ability to compare all of the results. If the researcher then uses her data for a clinical trial that is subsequently submitted to the FDA for regulatory review, she can be confident that all of the necessary details are included to successfully repeat her analysis. The final result is a transparent, efficient process that may substantially reduce the amount of time, money, and potential confusion involved with clarifying any details that might otherwise have been overlooked.

In order to form a more cooperative community we have migrated all of the BioCompute development to GitHub (<https://github.com/biocompute-objects>) and are following the OpenStand.org principles for collaborative open standards development (<https://open-stand.org/about-us/principles/>). The BioCompute specification has been published as a GitHub repository so that comments and issues can be addressed using the GitHub issue-tracking system. The GitHub organization is setup to have a separate repository for each of the use case examples or implementations. Each of these will be able to link back to the specification, encouraging parallelized development.

BCO use case: The relational sequencing tuberculosis-unified variant pipeline

Several BCO examples are available in repositories under the BioCompute organization on GitHub (<https://github.com/biocompute-objects/>). As an example, the Unified Variant Pipeline (UVP) BCO is described. The complex UVP BCO captures a validated whole genome sequencing (WGS) analysis pipeline, for The Relational Sequencing Tuberculosis (ReSeqTB) knowledge base. This knowledge base is a one-stop resource for curated *Mycobacterium tuberculosis* genotypic and phenotypic data that have been standardized and aggregated [47]. It is designed as a global resource to rapidly predict *M. tuberculosis* antimicrobial resistance from raw sequence data. The UVP is an Illumina-based consensus NGS pipeline comprising several bioinformatic tools with defined thresholds designed to annotate and produce a list of mutations (SNPs and indels) in comparison to a reference *M. tuberculosis* genome. The ReSeqTB platform also includes a database containing phenotypic metadata from culture-based testing and molecular probe-based genotypic data, all housed in a cloud environment enabling multitiered user access. In addition, a web-based application containing public data is now freely available (www.reseqtb.org).

A publicly available BCO was developed and released for the UVP to standardize and communicate the process for inputting drug resistance profiles using sequence-based technologies. The UVP BCO for ReSeqTB delineates the key aspects or domains in both human- and machine- readable JSON file format. The current version of the UVP BCO can be found in the public BCO GitHub (<https://github.com/biocompute-objects/UVP-BCO>). As multiple bioinformatic pipelines currently exist for *M. tuberculosis*, the BioCompute paradigm will allow users, including regulators, to identify variables and the ways in which they differ between pipelines, while assessing sources of error, in order to standardize and validate HTS results required for patient management.

Discussion and conclusions

Robust and reproducible data analysis is key to successful personalized medicine and genomic initiatives. Researchers, clinicians, administrators, and patients are all tied by the information in electronic health records (EHRs) and databases. Current systems rely on data stored with incomplete provenance records and in different computing languages. This has created a cumbersome and inefficient healthcare environment.

The initiatives discussed in this manuscript seek to make data and analyses communicable, repeatable, and reproducible to facilitate collaboration and information sharing from data

producers to data users. Without an infrastructure like BCO, increased HTS creates silos of unusable data, making standardized regulation of reproducibility more difficult. To clear the bottleneck for downstream analysis, the provenance (or origin) of data along with the analysis details (e.g., parameters, workflow versions) must be tracked to ensure accuracy and validity. The development of high-throughput, cloud-based infrastructures (such as DNAnexus, Galaxy, HIVE, and Seven Bridges Genomics) enables users to capture data provenance and store the analyses in infrastructures that allow easy user interaction and creation of BCOs according to BCO specifications described above and in the specification document (https://github.com/biocompute-objects/BCO_Specification).

Platform-independent provenance has largely been ignored in HTS. Emerging standards enable both representation of genomic information and linking of provenance information. By harmonizing across these standards, provenance information can be captured across both clinical and research settings, extending into the conceptual, experimental methods and the underlying computational workflows. There are several use cases of such work, including submission for FDA diagnostic evaluations, the original use case for the BCO. Such standards also enable robust and reproducible science and facilitate open science between collaborators, and the development of these standards is meant to satisfy the needs of downstream consumers of genomic information.

The need to reproducibly communicate HTS computational analyses and results has led to collaboration among disparate industry groups. Through outreach activities—including conferences and workshops—awareness of the importance of standardization, tracking, and reproducibility methods has improved [9,48]. Standards like FHIR and ROs capture the underlying data provenance to be shared in frameworks like GA4GH, enabling collaboration around reproducible data and analyses. New computing standards like CWL increase the scalability and reproducibility of data analysis. The BioCompute paradigm acts as a harmonizing umbrella to facilitate human and machine communication, increasing interoperability in fields that use genomic data. Detailed BioCompute specifications (available at: https://github.com/biocompute-objects/BCO_Specification/) can be used to generate BCOs by any bioinformatics platform that pulls underlying data and analysis provenance into its infrastructure. Ongoing BCO pilots are currently working to streamline the flow, with the goal of providing users with effortlessly reproducible bioinformatics analyses. As BCOs aim to simplify the review of data that are essential for FDA approval, these pilots mirror clinical trials involving HTS data for FDA submissions.

Subsequently, the incorporation of a verification test kit (a test for the range of input values that will still produce the same output values) that evaluates the integrity of the BCO will enable pipelines described by a BCO to be implemented on other servers. The verification test kit would consist of simulated or real HTS data sets in which the expected results are known. Each unique implementation of the BCO can then be tested using this verification kit to ensure that it reports the expected results. This work will further the development of the error domain, which describes the observed deviations from the expected results. We can arrive at values in this domain by repeatedly running the pipeline on the verification kit data and testing the analysis methods deployed. In this way, BioCompute can fulfill the goal of communicating exactly what was done and also what it means scientifically.

Community involvement has grown to more than 300 contributors, participants, and collaborators from public institutions (including NCI, the FDA, and others), universities (including George Washington University, University of Manchester, Harvard, and others), and several private sector partners. The BioCompute effort has resulted in two publications, three workshops, and FDA submissions. Efforts to integrate the BioCompute standard into existing platforms mentioned above are underway. The standard has also moved to Working Group

status within the Institute of Electrical and Electronics Engineers (IEEE; <http://sites.ieee.org/sagroups-2791/>) and is expected to be an International Organization for Standardization (ISO)-recognized standard. Finally, publicly accessible databases of crowd-sourced BCOs that allow researchers or clinicians to reproduce workflows in a variety of contexts are also planned.

Supporting information

S1 Text. Additional details on projects, concepts, and examples mentioned in the main text of the publication.

(DOCX)

Acknowledgments

We would like to recognize all the speakers and participants of the 2017 HTS-CSRS workshop who facilitated the discussion on standardizing HTS computations and analyses. The workshop was cosponsored by FDA and GW. The comments and input during the workshop were processed and integrated into the BCO specification document and this manuscript. The participants of the 2017 HTS-CSRS workshop are listed here: <https://osf.io/h59uh/wiki/2017%20HTS-CSRS%20Workshop/>. The contributions of the authors are an informal communication and represent their own views.

References

1. Mailman MD, Feolo M, Jin Y, Kimura M, Tryka K, et al. (2007) The NCBI dbGaP database of genotypes and phenotypes. *Nat Genet* 39: 1181–1186. <https://doi.org/10.1038/ng1007-1181> PMID: 17898773
2. Landrum MJ, Lee JM, Riley GR, Jang W, Rubinstein WS, et al. (2014) ClinVar: public archive of relationships among sequence variation and human phenotype. *Nucleic Acids Res* 42: D980–985. <https://doi.org/10.1093/nar/gkt1113> PMID: 24234437
3. Zheng J, Erzurumluoglu AM, Elsworth BL, Kemp JP, Howe L, et al. (2017) LD Hub: a centralized database and web interface to perform LD score regression that maximizes the potential of summary level GWAS data for SNP heritability and genetic correlation analysis. *Bioinformatics* 33: 272–279. <https://doi.org/10.1093/bioinformatics/btw613> PMID: 27663502
4. Sawyer E (2017) High Throughput Sequencing and Cost Trends. *Nature Education*.
5. Committee on the Review of Omics-Based Tests for Predicting Patient Outcomes in Clinical Trials; Board on Health Care Services; Board on Health Sciences Policy; Institute of Medicine; Micheel CM, Nass SJ, Omenn GS, editors. (2012). *Evolution of Translational Omics: Lessons Learned and the Path Forward*. Washington (DC). PMID: 24872966
6. Beredimas N, Kilintzis V, Chouvarda I, Maglaveras N (2015) A reusable ontology for primitive and complex HL7 FHIR data types. *Conf Proc IEEE Eng Med Biol Soc* 2015: 2547–2550. <https://doi.org/10.1109/EMBC.2015.7318911> PMID: 26736811
7. Alterovitz G, Warner J, Zhang P, Chen Y, Ullman-Cullere M, et al. (2015) SMART on FHIR Genomics: facilitating standardized clinico-genomic apps. *J Am Med Inform Assoc* 22: 1173–1178. <https://doi.org/10.1093/jamia/ocv045> PMID: 26198304
8. Lawler M, Siu LL, Rehm HL, Chanock SJ, Alterovitz G, et al. (2015) All the World's a Stage: Facilitating Discovery Science and Improved Cancer Care through the Global Alliance for Genomics and Health. *Cancer Discov* 5: 1133–1136. <https://doi.org/10.1158/2159-8290.CD-15-0821> PMID: 26526696
9. Peter Amstutz MRC, Nebojša Tijić (editors), Brad Chapman, John Chilton, Michael Heuer, Andrey Kartashov, Dan Leehr, Hervé Ménager, Maya Nedeljkovich, Matt Scales, Stian Soiland-Reyes, Luka Stojanovic (2016) *Common Workflow Language, Specification*, Common Workflow Language working group.
10. Bechhofer S, Buchan I, De Roure D, Missier P, Ainsworth J, et al. (2013) Why linked data is not enough for scientists. *Future Generation Computer Systems-the International Journal of Grid Computing and Esience* 29: 599–611.

11. Simonyan V, Goecks J, Mazumder R (2017) Biocompute Objects-A Step towards Evaluation and Validation of Biomedical Scientific Computations. *PDA J Pharm Sci Technol* 71: 136–146. <https://doi.org/10.5731/pdajpst.2016.006734> PMID: 27974626
12. Wilkinson MD, Dumontier M, Aalbersberg IJ, Appleton G, Axton M, et al. (2016) The FAIR Guiding Principles for scientific data management and stewardship. *Sci Data* 3: 160018. <https://doi.org/10.1038/sdata.2016.18> PMID: 26978244
13. Manolio TA, Brooks LD, Collins FS (2008) A HapMap harvest of insights into the genetics of common disease. *J Clin Invest* 118: 1590–1605. <https://doi.org/10.1172/JCI34772> PMID: 18451988
14. Boyd SD (2013) Diagnostic applications of high-throughput DNA sequencing. *Annu Rev Pathol* 8: 381–410. <https://doi.org/10.1146/annurev-pathol-020712-164026> PMID: 23121054
15. NIH (2018) NIH Strategic Plan For Data Science. In: Research OoE, editor.
16. Bishop D (2015) Reproducibility and reliability of biomedical research. The Academy of Medical Sciences.
17. Pusztai L, Hatzis C, Andre F (2013) Reproducibility of research and preclinical validation: problems and solutions. *Nat Rev Clin Oncol* 10: 720–724. <https://doi.org/10.1038/nrclinonc.2013.171> PMID: 24080600
18. Samuel Reich E (2011) Cancer trial errors revealed. *Nature* 469: 139–140. <https://doi.org/10.1038/469139a> PMID: 21228842
19. Goodman SN, Fanelli D, Ioannidis JP (2016) What does research reproducibility mean? *Sci Transl Med* 8: 341ps312.
20. Belhajjame K, Zhao J, Garijo D, Gamble M, Hettne K, et al. (2015) Using a suite of ontologies for preserving workflow-centric research objects. *Journal of Web Semantics* 32: 16–42.
21. Kjer KM, Gillespie JJ, Ober KA (2007) Opinions on multiple sequence alignment, and an empirical comparison of repeatability and accuracy between POY and structural alignment. *Syst Biol* 56: 133–146. <https://doi.org/10.1080/10635150601156305> PMID: 17366144
22. Bose R, Frew J (2005) Lineage retrieval for scientific data processing: A survey. *Acm Computing Surveys* 37: 1–28.
23. Buneman P, Khanna S. & Wang-Chiew T (2001) Why and Where: A Characterization of Data Provenance. In *Database Theory*. Springer Lecture Notes in Computer Science: pp. 87–93.
24. Freire J, Bonnet, P. & Shasha, D. (2012) Computational Reproducibility: State-of-the-art, Challenges, and Database Research Opportunities. *SIGMOD Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data*: pp. 593–596.
25. Alper P. Enhancing and Abstracting Scientific Workflow Provenance for Data Publishing; 2013.
26. Gil Y. DE, Ellisman M., Fahringer T., Fox G., Gannon D., Goble C., Livny M. Moreau L., Myers J (2007) Examining the Challenges of Scientific Workflows. *IEEE Computer Society Computing Practices*: 9.
27. Reichman OJ, Jones MB, Schildhauer MP (2011) Challenges and opportunities of open data in ecology. *Science* 331: 703–705. <https://doi.org/10.1126/science.1197962> PMID: 21311007
28. Moreau L, Clifford B, Freire J, Futrelle J, Gil Y, et al. (2011) The Open Provenance Model core specification (v1.1). *Future Generation Computer Systems-the International Journal of Grid Computing and Esience* 27: 743–756.
29. Ciccarese P, Soiland-Reyes S, Belhajjame K, Gray AJ, Goble C, et al. (2013) PAV ontology: provenance, authoring and versioning. *J Biomed Semantics* 4: 37. <https://doi.org/10.1186/2041-1480-4-37> PMID: 24267948
30. Goble CA, Bhagat J, Aleksejevs S, Cruickshank D, Michaelides D, et al. (2010) myExperiment: a repository and social network for the sharing of bioinformatics workflows. *Nucleic Acids Res* 38: W677–W682. <https://doi.org/10.1093/nar/gkq429> PMID: 20501605
31. Garijo D, Gil Y, Corcho O (2017) Abstract, link, publish, exploit: An end to end framework for workflow sharing. *Future Generation Computer Systems-the International Journal of Esience* 75: 271–283.
32. Cohen-Boulakia S, Belhajjame K, Collin O, Chopard J, Froidevaux C, et al. (2017) Scientific workflows for computational reproducibility in the life sciences: Status, challenges and opportunities. *Future Generation Computer Systems-the International Journal of Esience* 75: 284–298.
33. Leipzig J (2017) A review of bioinformatic pipeline frameworks. *Brief Bioinform* 18: 530–536. <https://doi.org/10.1093/bib/bbw020> PMID: 27013646
34. Spjuth O, Bongcam-Rudloff E, Hernandez GC, Forer L, Giovacchini M, et al. (2015) Experiences with workflows for automating data-intensive bioinformatics. *Biol Direct* 10: 43. <https://doi.org/10.1186/s13062-015-0071-8> PMID: 26282399

35. Xu J, Thakkar S, Gong B, Tong W (2016) The FDA's Experience with Emerging Genomics Technologies—Past, Present, and Future. *AAPS J* 18: 814–818. <https://doi.org/10.1208/s12248-016-9917-y> PMID: [27116022](https://pubmed.ncbi.nlm.nih.gov/27116022/)
36. Metzker ML (2010) Sequencing technologies—the next generation. *Nat Rev Genet* 11: 31–46. <https://doi.org/10.1038/nrg2626> PMID: [19997069](https://pubmed.ncbi.nlm.nih.gov/19997069/)
37. Simonyan V, Mazumder R (2014) High-Performance Integrated Virtual Environment (HIVE) Tools and Applications for Big Data Analysis. *Genes (Basel)* 5: 957–981.
38. Thain D, Tannenbaum T, Livny M (2005) Distributed computing in practice: the Condor experience. *Concurrency and Computation-Practice & Experience* 17: 323–356.
39. Tomczak K, Czerwinska P, Wiznerowicz M (2015) The Cancer Genome Atlas (TCGA): an immeasurable source of knowledge. *Contemp Oncol (Pozn)* 19: A68–77.
40. Simonyan V, Chumakov K, Dingerdissen H, Faison W, Goldweber S, et al. (2016) High-performance integrated virtual environment (HIVE): a robust infrastructure for next-generation sequence data analysis. *Database (Oxford)* 2016.
41. Afgan E, Baker D, van den Beek M, Blankenberg D, Bouvier D, et al. (2016) The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nucleic Acids Res* 44: W3–W10. <https://doi.org/10.1093/nar/gkw343> PMID: [27137889](https://pubmed.ncbi.nlm.nih.gov/27137889/)
42. Woodcock J, Woosley R (2008) The FDA critical path initiative and its influence on new drug development. *Annu Rev Med* 59: 1–12. <https://doi.org/10.1146/annurev.med.59.090506.155819> PMID: [18186700](https://pubmed.ncbi.nlm.nih.gov/18186700/)
43. Kaye J, Heeney C, Hawkins N, de Vries J, Boddington P (2009) Data sharing in genomics—re-shaping scientific practice. *Nat Rev Genet* 10: 331–335. <https://doi.org/10.1038/nrg2573> PMID: [19308065](https://pubmed.ncbi.nlm.nih.gov/19308065/)
44. Whitty CJ (2017) The contribution of biological, mathematical, clinical, engineering and social sciences to combatting the West African Ebola epidemic. *Philos Trans R Soc Lond B Biol Sci* 372.
45. Buneman P, Khanna S. & Tan W.-C (2000) *Data Provenance: Some Basic Issues*. Springer Foundations of Software Technology and Theoretical Computer Science: pp. 87–93.
46. Kenall A, Harold S, Foote C (2014) An open future for ecological and evolutionary data? *BMC Evol Biol* 14: 66. <https://doi.org/10.1186/1471-2148-14-66> PMID: [24690275](https://pubmed.ncbi.nlm.nih.gov/24690275/)
47. Starks AM, Aviles E, Cirillo DM, Denkinger CM, Dolinger DL, et al. (2015) Collaborative Effort for a Centralized Worldwide Tuberculosis Relational Sequencing Data Platform. *Clin Infect Dis* 61 Suppl 3: S141–146.
48. Hettne KM, Dharuri H, Zhao J, Wolstencroft K, Belhajjame K, et al. (2014) Structuring research methods and data with the research object model: genomics workflows as a case study. *J Biomed Semantics* 5: 41. <https://doi.org/10.1186/2041-1480-5-41> PMID: [25276335](https://pubmed.ncbi.nlm.nih.gov/25276335/)