



IEEE Standard for Bioinformatics Analyses Generated by High-Throughput Sequencing (HTS) to Facilitate Communication

DOI:

[10.1109/IEEESTD.2020.9094416](https://doi.org/10.1109/IEEESTD.2020.9094416)

Document Version

Accepted author manuscript

[Link to publication record in Manchester Research Explorer](#)

Citation for published version (APA):

IEEE P2791 BioCompute Working Group (BCOWG), Mazumder, R. (Ed.), & Simonyan, V. (Ed.) (2020). *IEEE Standard for Bioinformatics Analyses Generated by High-Throughput Sequencing (HTS) to Facilitate Communication: IEEE Std 2791-2020*. IEEE. <https://doi.org/10.1109/IEEESTD.2020.9094416>

Citing this paper

Please note that where the full-text provided on Manchester Research Explorer is the Author Accepted Manuscript or Proof version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version.

General rights

Copyright and moral rights for the publications made accessible in the Research Explorer are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Takedown policy

If you believe that this document breaches copyright please refer to the University of Manchester's Takedown Procedures [<http://man.ac.uk/04Y6Bo>] or contact uml.scholarlycommunications@manchester.ac.uk providing relevant details, so we can investigate your claim.



1 **P2791™/D4**
2 **Draft Standard for Bioinformatics**
3 **Analyses Generated by High-**
4 **Throughput Sequencing (HTS) to**
5 **Facilitate Communication**

6 Sponsor
7
8 **Standards Committee**
9 of the
10 **IEEE Engineering in Medicine and Biology Society**
11

12
13 Approved <Date Approved>
14

15 **IEEE-SA Standards Board**
16

17 Copyright © 2018 by The Institute of Electrical and Electronics Engineers, Inc.
18 Three Park Avenue
19 New York, New York 10016-5997, USA

20 All rights reserved.

21 This document is an unapproved draft of a proposed IEEE Standard. As such, this document is subject to
22 change. USE AT YOUR OWN RISK! IEEE copyright statements SHALL NOT BE REMOVED from draft
23 or approved IEEE standards, or modified in any way. Because this is an unapproved draft, this document
24 must not be utilized for any conformance/compliance purposes. Permission is hereby granted for officers
25 from each IEEE Standards Working Group or Committee to reproduce the draft document developed by
26 that Working Group for purposes of international standardization consideration. IEEE Standards
27 Department must be informed of the submission for consideration prior to any reproduction for
28 international standardization consideration (stds.ipr@ieee.org). Prior to adoption of this document, in
29 whole or in part, by another standards development organization, permission must first be obtained from
30 the IEEE Standards Department (stds.ipr@ieee.org). When requesting permission, IEEE Standards
31 Department will require a copy of the standard development organization's document highlighting the use
32 of IEEE content. Other entities seeking permission to reproduce this document, in whole or in part, must
33 also obtain permission from the IEEE Standards Department.

34 IEEE Standards Department
35 445 Hoes Lane
36 Piscataway, NJ 08854, USA

37

1 **Abstract:** A major goal of this standard is to improve communication of bioinformatics protocols
2 and data in order to facilitate bioinformatics workflow related exchange and communication
3 between regulatory agencies, pharmaceutical companies, bioinformatics platform providers and
4 researchers. Detailed communication helps ensure responsibility, reproducibility, verify
5 bioinformatics protocol, track provenance information and promote interoperability. In addition,
6 this standard also defines the assurance program for evaluating and certifying products against
7 those requirements.

8
9 **Keywords:** genomics, next generation sequencing, high throughput sequencing, massively
10 parallel sequencing, NGS, HTS, MPS, workflow, pipeline, bioinformatics, analysis, regulatory
11
12

13 **OPEN SOURCE NOTICE:** This IEEE project incorporates open source software [File download:
14 <https://gitlab.com/IEEE-SA/2791/schema/-/archive/v1.3.0-alpha/schema-v1.3.0-alpha.zip>; Web
15 link: <https://gitlab.com/IEEE-SA/2791/schema/tree/v1.3.0-alpha>] under a BSD 3-Clause license
16 either within the draft or as a separate file. Any person contributing material to IEEE open source
17 software during standards development, SA ballot, or Public Review is required to provide the
18 appropriate license to IEEE (IEEE Contributor License Agreement or CLA). Please note that any
19 Contributions to IEEE open source software that is submitted without first providing the
20 appropriate CLA to IEEE will not be eligible for inclusion either in the draft standard or the open
21 source software, and may not be considered by the comment resolution group. The appropriate
22 CLA form (<https://app.box.com/s/2nlzxsru3jewdy9ello4rfkokwt9k841>) must be submitted to os-
23 fb@ieee.org, if applicable. Note that Federal or Crown employees must submit the CLA that
24 includes Appendix A. *

The Institute of Electrical and Electronics Engineers, Inc.
3 Park Avenue, New York, NY 10016-5997, USA

Copyright © 2018 by The Institute of Electrical and Electronics Engineers, Inc.
All rights reserved. Published <Date Published>. Printed in the United States of America.

IEEE is a registered trademark in the U.S. Patent & Trademark Office, owned by The Institute of Electrical and Electronics
Engineers, Incorporated.

PDF: ISBN 978-0-XXXX-XXXX-X STDXXXXX
Print: ISBN 978-0-XXXX-XXXX-X STDPDXXXXX

IEEE prohibits discrimination, harassment, and bullying.

For more information, visit <http://www.ieee.org/web/aboutus/whatis/policies/p9-26.html>.

No part of this publication may be reproduced in any form, in an electronic retrieval system or otherwise, without the prior written permission
of the publisher.

1 **Important Notices and Disclaimers Concerning IEEE Standards Documents**

2 IEEE documents are made available for use subject to important notices and legal disclaimers. These
3 notices and disclaimers, or a reference to this page, appear in all standards and may be found under the
4 heading “Important Notices and Disclaimers Concerning IEEE Standards Documents.” They can also be
5 obtained on request from IEEE or viewed at <http://standards.ieee.org/ipr/disclaimers.html>.

6 **Notice and Disclaimer of Liability Concerning the Use of IEEE Standards** 7 **Documents**

8 IEEE Standards documents (standards, recommended practices, and guides), both full-use and trial-use, are
9 developed within IEEE Societies and the Standards Coordinating Committees of the IEEE Standards
10 Association (“IEEE-SA”) Standards Board. IEEE (“the Institute”) develops its standards through a
11 consensus development process, approved by the American National Standards Institute (“ANSI”), which
12 brings together volunteers representing varied viewpoints and interests to achieve the final product. IEEE
13 Standards are documents developed through scientific, academic, and industry-based technical working
14 groups. Volunteers in IEEE working groups are not necessarily members of the Institute and participate
15 without compensation from IEEE. While IEEE administers the process and establishes rules to promote
16 fairness in the consensus development process, IEEE does not independently evaluate, test, or verify the
17 accuracy of any of the information or the soundness of any judgments contained in its standards.

18 IEEE Standards do not guarantee or ensure safety, security, health, or environmental protection, or ensure
19 against interference with or from other devices or networks. Implementers and users of IEEE Standards
20 documents are responsible for determining and complying with all appropriate safety, security,
21 environmental, health, and interference protection practices and all applicable laws and regulations.

22 IEEE does not warrant or represent the accuracy or content of the material contained in its standards, and
23 expressly disclaims all warranties (express, implied and statutory) not included in this or any other
24 document relating to the standard, including, but not limited to, the warranties of: merchantability; fitness
25 for a particular purpose; non-infringement; and quality, accuracy, effectiveness, currency, or completeness
26 of material. In addition, IEEE disclaims any and all conditions relating to: results; and workmanlike effort.
27 IEEE standards documents are supplied “AS IS” and “WITH ALL FAULTS.”

28 Use of an IEEE standard is wholly voluntary. The existence of an IEEE standard does not imply that there
29 are no other ways to produce, test, measure, purchase, market, or provide other goods and services related
30 to the scope of the IEEE standard. Furthermore, the viewpoint expressed at the time a standard is approved
31 and issued is subject to change brought about through developments in the state of the art and comments
32 received from users of the standard.

33 In publishing and making its standards available, IEEE is not suggesting or rendering professional or other
34 services for, or on behalf of, any person or entity nor is IEEE undertaking to perform any duty owed by any
35 other person or entity to another. Any person utilizing any IEEE Standards document, should rely upon his
36 or her own independent judgment in the exercise of reasonable care in any given circumstances or, as
37 appropriate, seek the advice of a competent professional in determining the appropriateness of a given
38 IEEE standard.

39 IN NO EVENT SHALL IEEE BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL,
40 EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO:
41 PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS;
42 OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY,
43 WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR
44 OTHERWISE) ARISING IN ANY WAY OUT OF THE PUBLICATION, USE OF, OR RELIANCE
45 UPON ANY STANDARD, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE AND
46 REGARDLESS OF WHETHER SUCH DAMAGE WAS FORESEEABLE.

1 **Translations**

2 The IEEE consensus development process involves the review of documents in English only. In the event
3 that an IEEE standard is translated, only the English version published by IEEE should be considered the
4 approved IEEE standard.

5 **Official statements**

6 A statement, written or oral, that is not processed in accordance with the IEEE-SA Standards Board
7 Operations Manual shall not be considered or inferred to be the official position of IEEE or any of its
8 committees and shall not be considered to be, or be relied upon as, a formal position of IEEE. At lectures,
9 symposia, seminars, or educational courses, an individual presenting information on IEEE standards shall
10 make it clear that his or her views should be considered the personal views of that individual rather than the
11 formal position of IEEE.

12 **Comments on standards**

13 Comments for revision of IEEE Standards documents are welcome from any interested party, regardless of
14 membership affiliation with IEEE. However, IEEE does not provide consulting information or advice
15 pertaining to IEEE Standards documents. Suggestions for changes in documents should be in the form of a
16 proposed change of text, together with appropriate supporting comments. Since IEEE standards represent a
17 consensus of concerned interests, it is important that any responses to comments and questions also receive
18 the concurrence of a balance of interests. For this reason, IEEE and the members of its societies and
19 Standards Coordinating Committees are not able to provide an instant response to comments or questions
20 except in those cases where the matter has previously been addressed. For the same reason, IEEE does not
21 respond to interpretation requests. Any person who would like to participate in revisions to an IEEE
22 standard is welcome to join the relevant IEEE working group.

23 Comments on standards should be submitted to the following address:

24 Secretary, IEEE-SA Standards Board
25 445 Hoes Lane
26 Piscataway, NJ 08854 USA

27 **Laws and regulations**

28 Users of IEEE Standards documents should consult all applicable laws and regulations. Compliance with
29 the provisions of any IEEE Standards document does not imply compliance to any applicable regulatory
30 requirements. Implementers of the standard are responsible for observing or referring to the applicable
31 regulatory requirements. IEEE does not, by the publication of its standards, intend to urge action that is not
32 in compliance with applicable laws, and these documents may not be construed as doing so.

33 **Copyrights**

34 IEEE draft and approved standards are copyrighted by IEEE under U.S. and international copyright laws.
35 They are made available by IEEE and are adopted for a wide variety of both public and private uses. These
36 include both use, by reference, in laws and regulations, and use in private self-regulation, standardization,
37 and the promotion of engineering practices and methods. By making these documents available for use and
38 adoption by public authorities and private users, IEEE does not waive any rights in copyright to the
39 documents.

40

41

1 Photocopies

2 Subject to payment of the appropriate fee, IEEE will grant users a limited, non-exclusive license to
3 photocopy portions of any individual standard for company or organizational internal use or individual,
4 non-commercial use only. To arrange for payment of licensing fees, please contact Copyright Clearance
5 Center, Customer Service, 222 Rosewood Drive, Danvers, MA 01923 USA; +1 978 750 8400. Permission
6 to photocopy portions of any individual standard for educational classroom use can also be obtained
7 through the Copyright Clearance Center.

8 Updating of IEEE Standards documents

9 Users of IEEE Standards documents should be aware that these documents may be superseded at any time
10 by the issuance of new editions or may be amended from time to time through the issuance of amendments,
11 corrigenda, or errata. A current IEEE document at any point in time consists of the current edition of the
12 document together with any amendments, corrigenda, or errata then in effect.

13 Every IEEE standard is subjected to review at least every ten years. When a document is more than ten
14 years old and has not undergone a revision process, it is reasonable to conclude that its contents, although
15 still of some value, do not wholly reflect the present state of the art. Users are cautioned to check to
16 determine that they have the latest edition of any IEEE standard.

17 In order to determine whether a given document is the current edition and whether it has been amended
18 through the issuance of amendments, corrigenda, or errata, visit IEEE Xplore at <http://ieeexplore.ieee.org/>
19 or contact IEEE at the address listed previously. For more information about the IEEE-SA or IEEE's
20 standards development process, visit the IEEE-SA Website at <http://standards.ieee.org>.

21 Errata

22 Errata, if any, for all IEEE standards can be accessed on the IEEE-SA Website at the following URL:
23 <http://standards.ieee.org/findstds/errata/index.html>. Users are encouraged to check this URL for errata
24 periodically.

25 Patents

26 Attention is called to the possibility that implementation of this standard may require use of subject matter
27 covered by patent rights. By publication of this standard, no position is taken by the IEEE with respect to
28 the existence or validity of any patent rights in connection therewith. If a patent holder or patent applicant
29 has filed a statement of assurance via an Accepted Letter of Assurance, then the statement is listed on the
30 IEEE-SA Website at <http://standards.ieee.org/about/sasb/patcom/patents.html>. Letters of Assurance may
31 indicate whether the Submitter is willing or unwilling to grant licenses under patent rights without
32 compensation or under reasonable rates, with reasonable terms and conditions that are demonstrably free of
33 any unfair discrimination to applicants desiring to obtain such licenses.

34 Essential Patent Claims may exist for which a Letter of Assurance has not been received. The IEEE is not
35 responsible for identifying Essential Patent Claims for which a license may be required, for conducting
36 inquiries into the legal validity or scope of Patents Claims, or determining whether any licensing terms or
37 conditions provided in connection with submission of a Letter of Assurance, if any, or in any licensing
38 agreements are reasonable or non-discriminatory. Users of this standard are expressly advised that
39 determination of the validity of any patent rights, and the risk of infringement of such rights, is entirely
40 their own responsibility. Further information may be obtained from the IEEE Standards Association.

1 Participants

2 At the time this draft standard was completed, the P2791 Working Group had the following membership:

3 **Raja Mazumder, Chair**
4 **Vahan Simonyan, Vice Chair**

5
6 Ogan Abaan 18 Paul Duncan 30 Rahi Navelkar
7 Jonas Almeida 19 Josep Gelpi 31 Asa Oudes
8 Gil Alterovitz 20 Carole Goble 32 Janisha Patel
9 Payal Banerjee 21 Jeremy Goecks 33 John Penn
10 Amanda Bell 22 Jonathan Jacobs 34 Megan Pottersbusch
11 Surajit Bhattacharya 23 Robel Kahsay 35 Jonathan Pryke
12 Lee Black 24 Jonathon Keeney 36 Stian Soiland-Reyes
13 Ben Busby 25 Charles Hadley King 37 Dan Taylor
14 Kristy Cloyd-Warwick 26 Jonathan LoTempio 38 Jason Travis
15 Ryan Connor 27 Xeandong Meng 39 Paul Walsh
16 Michael Crusoe 28 David Michaels 40 Jianchao Yao
17 Dennis Dean 29 Hiroki Morizono

41

42 The following members of the individual/entity balloting committee voted on this standard. Balloters may
43 have voted for approval, disapproval, or abstention.

44 *[To be supplied by IEEE]*

45 Balloter1 48 Balloter4 51 Balloter7
46 Balloter2 49 Balloter5 52 Balloter8
47 Balloter3 50 Balloter6 53 Balloter9

54

55 When the IEEE-SA Standards Board approved this standard on <Date Approved>, it had the following
56 membership:

57 *[To be supplied by IEEE]*

58 **<Name>, Chair**
59 **<Name>, Vice Chair**
60 **<Name>, Past Chair**
61 **Konstantinos Karachalios, Secretary**

62 SBMember1 65 SBMember4 68 SBMember7
63 SBMember2 66 SBMember5 69 SBMember8
64 SBMember3 67 SBMember6 70 SBMember9

71 *Member Emeritus

72

1 Introduction

2 This introduction is not part of P2791/D4, Draft Standard for Bioinformatics Analyses Generated by High-Throughput
3 Sequencing (HTS) to Facilitate Communication.

4 The P2791 specification enables the description of bioinformatic genome analysis workflows in a
5 standardized way. P2791 addresses the tremendous variability and uncertainty in communicating
6 bioinformatics workflows and data related to analysis as a result of high throughput sequencing (HTS). The
7 need to resolve issues in communication was felt particularly strongly between the United States Food and
8 Drug Administration (FDA) and the entities that submit any work to the FDA for regulatory analysis that
9 includes an HTS component¹² (<https://doi.org/10.5731/pdajpst.2016.006734> and [PMC5510742](https://doi.org/10.1093/bioinformatics/btu1074)). A plan for
10 what would become P2791 and initial goals of the project were drafted in a collaboration between the
11 George Washington University and the FDA in 2014. The project has grown since then to include
12 publications, workshops, applied use cases, and a large community of participants and collaborators. P2791
13 Objects created according to this standard are intended 1) to be both human and machine readable, 2) to be
14 applied to genomic analysis workflows, and 3) to be able to capture details related to a workflow in such a
15 way as to facilitate efficient communication and improve reproducibility and interoperability. Efforts were
16 made to accommodate as many tools, platforms or scripts as possible, and to be adaptable to future
17 developments in this field under a unified set of descriptions to standardize and streamline the
18 representations of such complex bioinformatics processes.

19 P2791 is a standard and a P2791 Object is an instance of that standard. High throughput sequencing (HTS),
20 also referred to as next-generation sequencing (NGS) or massively parallel sequencing (MPS), has
21 increased the pace at which we generate, compute and share genomic data in biomedical sciences. As a
22 result, scientists, clinicians and regulators are now faced with a new data paradigm that is less portable,
23 more complex and most of all poorly standardized. The P2791 Objects are written in JSON format to
24 encode important information on the execution of computational pipelines, or for the creation of knowledge
25 bases. P2791 can be considered to be process oriented (for software pipelines) and/or product oriented (for
26 knowledge bases). The goal of using a P2791 Object is to streamline communication of these otherwise
27 difficult to elucidate details between stakeholders in academia, industry and regulatory agencies.

28 Standardized HTS data processing descriptions and data formats will promote interoperability and simplify
29 the verification of the bioinformatics protocols applied against data. To do this, a schema has been
30 developed to represent instances of computational analysis as a P2791 Object. A P2791 Object includes:

- 31 — Information about parameters and versions of the executable programs in a pipeline
- 32 — Reference to input and output test data for verification of the pipeline
- 33 — A usability domain
- 34 — Keywords
- 35 — A list of agents involved along with other important metadata, such as their specific contribution

36 Knowledge of input data is intended to be captured according to existing efforts, such as Minimum
37 Information Required about a Glycomics Experiment (MIRAGE)³, Minimum Information about a
38

¹ Alterovitz G et al. Enabling Precision Medicine via standard communication of NGS provenance, analysis, and results. *PLoS Biol.* 2018 Dec; 16(12):e3000099 DOI: <https://doi.org/10.1371/journal.pbio.3000099>.

² Simonyan V, Goecks J and Mazumder R. Biocompute Objects—A Step towards Evaluation and Validation of Biomedical Scientific Computations. *PDA J Pharm Sci Technol.* 2017 Mar-Apr;71(2):136-146

³ Kolarich, Daniel; Rapp, Erdmann; Struwe, Weston B.; Haslam, Stuart M.; Zaia, Joseph; McBride, Ryan; Agravat, Sanjay; Campbell, Matthew P.; Kato, Masaki; Ranzinger, Rene; Kettner, Carsten; York, William S. (1 April 2013). "The Minimum Information Required for a Glycomics Experiment (MIRAGE) Project: Improving the Standards for Reporting Mass-spectrometry-based Glycoanalytic Data". *Molecular & Cellular Proteomics.* 12 (4): 991–995. doi:10.1074/mcp.O112.026492. ISSN 1535-9476. PMC 3617344. PMID 23378518

1 Proteomics Experiment (MIAPE)⁴, Standards for Reporting Enzymology Data (STRENDA)⁵ and to be
2 in accordance with Minimum Information Standards⁶. In addition to all the information captured in the
3 P2791 Object, the P2791 Object itself is intended to be independent of the execution environment,
4 whether it is a local or a cloud-based infrastructure.

5

⁴ Taylor, C. F.; Paton, N. W.; Lilley, K. S.; Binz, P. A.; Julian Jr, R. K.; Jones, A. R.; Zhu, W.; Apweiler, R.; Aebersold, R.; Deutsch, E. W.; Dunn, M. J.; Heck, A. J. R.; Leitner, A.; Macht, M.; Mann, M.; Martens, L.; Neubert, T. A.; Patterson, S. D.; Ping, P.; Seymour, S. L.; Souda, P.; Tsugita, A.; Vandekerckhove, J.; Vondriska, T. M.; Whitelegge, J. P.; Wilkins, M. R.; Xenarios, I.; Yates Jr, J. R.; Hermjakob, H. (2007). "The minimum information about a proteomics experiment (MIAPE)". *Nature Biotechnology*. 25 (8): 887–893. doi:10.1038/nbt1329. PMID 17687369

⁵ Tipton, K.F., Armstrong, R.N., Bakker, B.M., Bairoch, A., Cornish-Bowden, A., Halling, P.J., Hofmeyr, J.-H., Leyh, T.S., Kettner, C., Raushel, F.M., Rohwer, J., Schomburg, D., Steinbeck, C. (2014) Standards for Reporting Enzyme Data: The STRENDA Consortium: What it aims to do and why it should be helpful. *Perspect. Sci.* 1(1.6):131-137. DOI: 10.1016/j.pisc.2014.02.012

⁶ Taylor, Chris F (2008). "Promoting coherent minimum reporting guidelines for biological and biomedical investigations: the MIBBI project". *Nature Biotechnology*. 26 (8): 889–896. doi:10.1038/nbt.1411. PMC 2771753. PMID 18688244

1 Contents

2	1. Overview	1
3	1.1 General	1
4	1.2 Scope	2
5	1.3 Purpose	2
6	2. Normative references.....	2
7	3. Definitions, acronyms, and abbreviations	2
8	3.1 Acronyms and abbreviations	2
9	4. P2791 Standard.....	3
10	4.1 General	3
11	Annex A (informative) Bibliography	5
12		
13		

1 Draft Standard for Bioinformatics 2 Analyses Generated by High- 3 Throughput Sequencing (HTS) to 4 Facilitate Communication

5 1. Overview

6 1.1 General

7 The P2791 standard captures relevant information from a high throughput sequencing workflow as a P2791
8 Object in order to enable a user to understand and interpret the workflow efficiently and with high
9 confidence. P2791 is a standard that was initially created with a goal of improving efficiency in regulatory
10 review. Pursuant to this, workflow steps and prerequisites to execute workflow steps are recorded in detail
11 in a P2791 Object. Information is recorded using key/value pairs in JavaScript Object Notation (JSON),
12 adhering to the P2791 JSON Schema.

13 Information in P2791 Objects is organized by domains;

- 14 • The Provenance Domain - tracks metadata about the P2791 Object
- 15 • The Usability Domain - tracks what was done
- 16 • The Extension Domain - provide user-defined fields
- 17 • The Description Domain - captures a description of external resources, pipeline steps, and the
18 relationships of I/O objects
- 19 • The Execution Domain - describes information needed for deployment, software configuration and
20 running applications in a dependent environment
- 21 • The Parametric Domain - captures all parameters that customize a computational flow
- 22 • The Input and Output Domain - contains a list of global input and output files
- 23 • The Error Domain - describes errors, including the limits of detectability, false positives, false
24 negatives, statistics confidence of outcomes, and description of errors (i.e. empirical or
25 algorithmic).

1

2 This document should be read in conjunction with the open source P2791 JSON Schema files
3 (<https://w3id.org/2791/>) which are referred to from the text, for instance “*provenance_domain.json*” refers
4 to (https://w3id.org/2791/provenance_domain.json). Files are kept separate for organization. References in
5 the P2791 Object schema (`$ref`) to these files should be replaced with the proper domain from the
6 appropriate file. For example, line 142 of “*p2791object.json*” (“`$ref`”:
7 “*provenance_domain.json*”) is a reference to the structure specified in the
8 *provenance_domain.json* file. The P2791 Object Schema builds on the JSON Schema by adding domains
9 in a way that facilitates the communication of bioinformatics workflows. A description of the domain files
10 follows.

11 1.2 Scope

12 This standard establishes detailed and structured communication of bioinformatics protocols in order to
13 facilitate bioinformatics workflow related exchange and communication between regulatory agencies,
14 pharmaceutical companies, bioinformatics platform providers and researchers. Detailed communication
15 helps ensure responsibility, verify bioinformatics protocol, track provenance information and promote
16 interoperability.

17 1.3 Purpose

18 The standard allows for the cross platform communications of complex computation from inception to
19 manufacturing of medical products and services. Another goal of this standard is to improve efficiency and
20 speed of communication.

21 2. Normative references

22 The following referenced documents are indispensable for the application of this document (i.e., they must
23 be understood and used, so each referenced document is cited in text and its relationship to this document is
24 explained). For dated references, only the edition cited applies. For undated references, the latest edition of
25 the referenced document (including any amendments or corrigenda) applies.

26
27 JSON (RFC8259): <https://tools.ietf.org/html/rfc8259>
28 JSON Schema: draft-handrews-json-schema-01
29 JSON Schema Validation: draft-handrews-json-schema-validation-01
30

31 3. Definitions, acronyms, and abbreviations

32 For the purposes of this document, the following terms and definitions apply. The *IEEE Standards*
33 *Dictionary Online* should be consulted for terms not defined in this clause. ¹

34 3.1 Acronyms and abbreviations

35 JSON JavaScript Object Notation

¹*IEEE Standards Dictionary Online* is available at: <http://dictionary.ieee.org>.

1 SCM Source Control Management

2 4. P2791 Standard

3 4.1 General

4 This document describes the P2791 standard for describing bioinformatic workflows. A P2791 “Object” is
5 an instance of the P2791 standard, and is a text file written in JSON data structure that shall consist of all
6 domains required by the P2791 Schema (<https://w3id.org/2791/>). The P2791 Schema is the formal
7 definition of the standard against which instances of the standard can be validated. JavaScript Object
8 Notation (JSON) is a textual format used by both instances of Objects and the formal P2791 Schema, and
9 the JSON Schema is the language used to express the P2791 Schema.

10
11 A valid Object must conform to the P2791 JSON Schema (see section 4.3), and therefore invokes all of the
12 requirements of the JSON Schema (while a valid Object file must conform to the schema, the schema file is
13 not technically required to create the Object file). Later versions of P2791 may be updated for conformance
14 with future JSON Schema versions. The minimum requirement to execute the standard is the fully
15 organized P2791 Object containing all domains in JSON Schema format. Pursuant to JSON schema, the
16 fields required for a valid Object are listed at the top of the *2791object.json* file.

17
18 All the files in the repository are linked together (using JSON pointers as described by the JSON Schema),
19 being referenced by the ‘*2791object.json*’ file. The *error_domain.json* is an optional domain to further
20 describe empirical and algorithmic sources and measures of error for a bioinformatics workflow
21 (https://w3id.org/2791/error_domain.json), and the *extension_domain.json* is an optional domain that
22 contains user-defined fields.

23
24 At its top-level, Objects have the following three required metadata fields: "spec_version",
25 "object_id", and "etag". These lines are external to all domains. Everything except for the etag,
26 object_id, and spec_version shall be included in the generation of an ETag (see
27 <https://tools.ietf.org/html/rfc7232#section-2.3>) - which can be "strong" or "weak" (see
28 <https://tools.ietf.org/html/rfc7232#section-2.1>). It is recommended that the ETag be deleted or updated if
29 the object file is changed (except in cases using weak ETags in which the entirety of the change comprises
30 a simple re-writing of the JSON).

31
32 object_id is a string that follows the JSON Schema format of namespace/ref shall be a unique
33 identifier. Users are free to number Object files in the manner of their choosing, however, in order to avoid
34 naming conflicts, it is recommended that a domain be registered with a registration authority, such as the
35 one at <https://www.biocomputeobject.org/registry.html>. For example,
36 http://www.example.com/exampleproject/1.3.0/schemas/ABC_object0001.json, where “ABC” is a registered
37 domain, and “_object0001.json” is an arbitrary identifier, chosen by the owner of that domain.

38
39 The remaining top level fields are domains that partition workflow into meaningful subunits. These are the
40 Description domain, Error domain, Execution domain, IO domain, Parametric domain, Provenance domain,
41 and Usability domain.

42
43 The Description Domain of an Object contains a description of external resources, pipeline steps, and the
44 relationship of I/O objects (https://w3id.org/2791/description_domain.json).

45
46 The Error Domain contains information related to the bounds of detection (such as the minimum sequence
47 depth and minimum sequence coverage), and statistical analyses of the pipeline (such as the false negative
48 and false positive rates). It is recommended that the keys directly under *empirical_error* and

1 algorithmic_error use a full URI. Resolving the URI should give a JSON Schema or textual
2 definition of the field. Other keys are not allowed in error_domain.
3

4 The Execution Domain of an Object describes details of deployment, software configuration, and running
5 applications in a dependent environment (https://w3id.org/2791/execution_domain.json). This may include
6 scripts, drivers, environment variables, and other software prerequisites.
7

8 The IO Domain of an Object is a list of global input and output files that may exist on local machine or on
9 another machine (https://w3id.org/2791/io_domain.json). It does not include references to intermediate
10 files.
11

12 The Parametric Domain of a P2791 Object includes any parameters used in a workflow
13 (https://w3id.org/2791/parametric_domain.json). This is intended for use when parameters are changed
14 from default settings.
15

16 The Provenance Domain contains metadata related to the Object
17 (https://w3id.org/2791/provenance_domain.json). It is used to track the flow of data from original source to
18 final computation, and includes contributors, reviewers, and versioning. In the event that a P2791 Object
19 retrospectively references an existing Object (such as an example Object), the derived_from field
20 within the Provenance Domain shall reference the specific Object by object_id field. In the event that
21 the Object is an example Object or is created de novo without reference to existing work, this field is not
22 included. In the event that the Object is an example or template P2791 Object, best practice is to state this
23 in the Usability Domain, along with relevant details (such as completeness of data, whether data is real or
24 artificial, etc.).

25 The Usability Domain of an Object is a plain language description of what was done in the workflow
26 (https://w3id.org/2791/usability_domain.json). This should align with the actual steps described elsewhere
27 in the Object. The Usability Domain conveys the purpose of the Object.
28

29 The Extension Domain allows a user to define additional fields and is optional. The Extension Domain is
30 for the inclusion of any additional structured information. A valid JSON schema for each extension used in
31 this domain is expected to be specified. The schema should be name spaced, and it is recommended that
32 resolving the namespaced URI will provide the extension's JSON Schema. The URL should be provided in
33 the required "extension_schema" field. If execution portability is desired, then the included script
34 should be in the Common Workflow Language v1.0 (<https://w3id.org/cwl/v1.0/>) or later format. In order to
35 avoid potential naming conflicts, it is recommended that users register their domain with
36
37

1 **Annex A**

2 (informative)

3 **Bibliography**

4 Bibliographical references are resources that provide additional or helpful material but do not need to be
5 understood or used to implement this standard. Reference to these resources is made for informational use
6 only.

7 [B1] Community User Guide for Best Practices. <https://w3id.org/biocompute/1.3.1>

8 [B2] JSON Schema: A Media Type for Describing JSON Documents. [https://tools.ietf.org/html/draft-](https://tools.ietf.org/html/draft-handrews-json-schema-01)
9 [handrews-json-schema-01](https://tools.ietf.org/html/draft-handrews-json-schema-01)

10