

New Tools for Building BioCompute Objects on DNAnexus, Seven Bridges, and Galaxy for deposit into an access-controlled BioComputeDB

September 13th, 2022

Jonathon Keeney, Ph.D.

George Washington University



Introduction Agenda

- Brief introduction to BioCompute
- Previous work
- Updates
- Today's workshop
 - BioComputeDB
- Introduction of other speakers

Need for Guidelines

- Lack of standards or guidelines for documentation and reporting of workflows
- Rich metadata not captured in workflow languages

Purpose

- What kind of data needs to be present in order to understand a computational analysis?
- How does that data need to be represented?

- Dozens of working groups
- 3 workshops
- Multiple draft specifications and schemas

IEEE SA
STANDARDS
ASSOCIATION

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

IEEE Engineering in Medicine and Biology Society

Developed by the
IEEE Standards Committee

IEEE Std 2791™-2020

 **IEEE**

STANDARDS



Electronic Submissions; Data Standards; Support for the International Institute of Electrical and Electronics Engineers Bioinformatics Computations and Analyses Standard for Bioinformatic Workflows

A Notice by the [Food and Drug Administration](#) on [07/22/2020](#)



Accepted at CBER, CDER, and CFSAN for most drug applications

PUBLISHED DOCUMENT



AGENCY:

Food and Drug Administration, Health and Human Services (HHS).



ACTION:

Notice.



1



SUMMARY:

The Food and Drug Administration (FDA or Agency) is announcing support for use in regulatory submissions the current version of the International Institute of Electrical and Electronics Engineers (IEEE) bioinformatics computations and analyses standard for bioinformatic workflows (BioCompute) and an update to



DOCUMENT DETAILS

Printed version:

[PDF](#)

Publication Date:

[07/22/2020](#)

Agencies:

[Food and Drug Administration](#)

Dates:

Submit either electronic or written comments on the notice by August 21, 2020.

Comments Close:

[08/21/2020](#)

Document Type:

Notice

Domain-based organization

- Who contributed to the work?
- What was their role?

Provenance Domain

User attribution and role in the work

Execution Domain

Description Domain

Parametric Domain

Usability Domain

IO Domain

Error Domain

Extension Domain

Object Information

Domain-based organization

- What was the execution environment?
 - e.g. HIVE, Galaxy, Seven Bridges, DNAnexus, command line, etc.
- Environmental variables?
- Software prerequisites?

Provenance Domain

Execution Domain

Execution environment needed to run the analysis

Description Domain

Parametric Domain

Usability Domain

IO Domain

Error Domain

Extension Domain

Object Information

Domain-based organization

- Keywords
- Description space for each step
- Input/output to describe what each step is doing

Provenance Domain

Execution Domain

Description Domain

Description of each step, including dependencies and IO

Parametric Domain

Usability Domain

IO Domain

Error Domain

Extension Domain

Object Information

Domain-based organization

- List of all parameters
 - Each parameter associated with a step

Provenance Domain

Execution Domain

Description Domain

Parametric Domain

List of parameters

Usability Domain

IO Domain

Error Domain

Extension Domain

Object Information

Domain-based organization

- Overview of pipeline
 - Purpose, goals, outcomes, context
 - Any other relevant information or author comments

Provenance Domain

Execution Domain

Description Domain

Parametric Domain

Usability Domain

Free text description

IO Domain

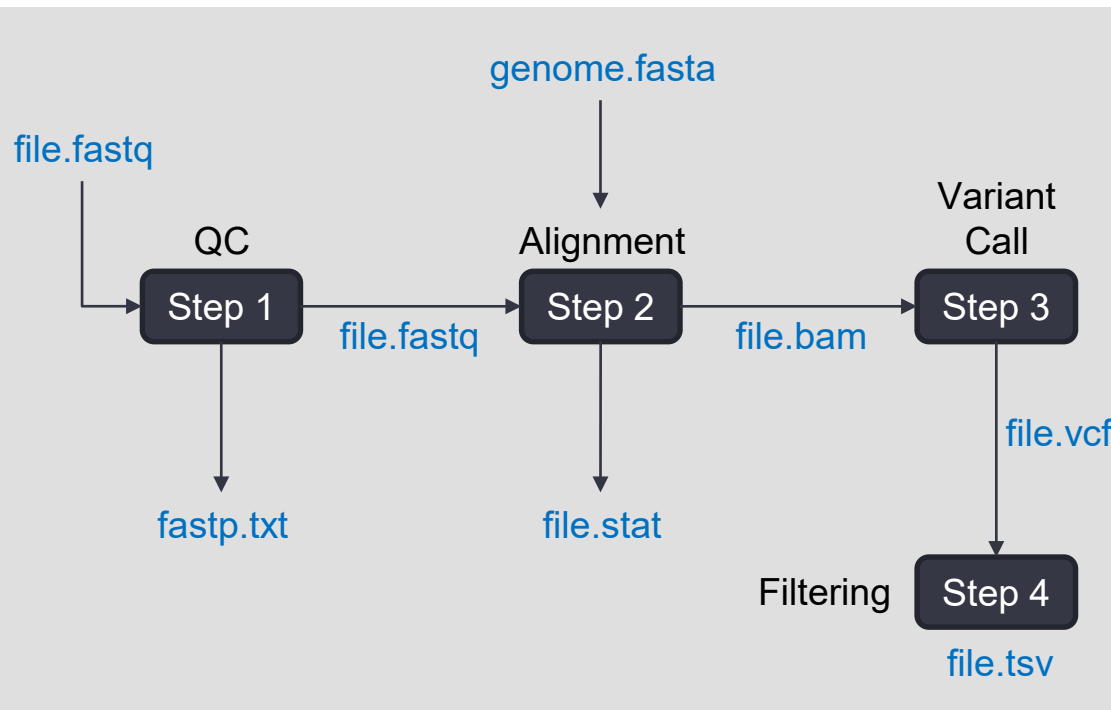
Error Domain

Extension Domain

Object Information

Domain-based organization

- Global input/output
 - Overview of what the pipeline needs and what it produces



Provenance Domain

Execution Domain

Description Domain

Parametric Domain

Usability Domain

IO Domain

Global inputs/outputs to pipeline

Error Domain

Extension Domain

Object Information

Domain-based organization

- Limits of detection
 - Empirical
 - Algorithmic

Provenance Domain

Execution Domain

Description Domain

Parametric Domain

Usability Domain

IO Domain

Error Domain

Error in the pipeline

Extension Domain

Object Information

Domain-based organization

- User-defined Domain
- Can extend beyond the base schema
- Requires a schema to validate

Provenance Domain

Execution Domain

Description Domain

Parametric Domain

Usability Domain

IO Domain

Error Domain

Extension Domain

User-defined

Object Information

Domain-based organization

- Object ID
- Specification Version
- eTag

Provenance Domain

Execution Domain

Description Domain

Parametric Domain

Usability Domain

IO Domain

Error Domain

Extension Domain

Object Information

Metadata about the pipeline

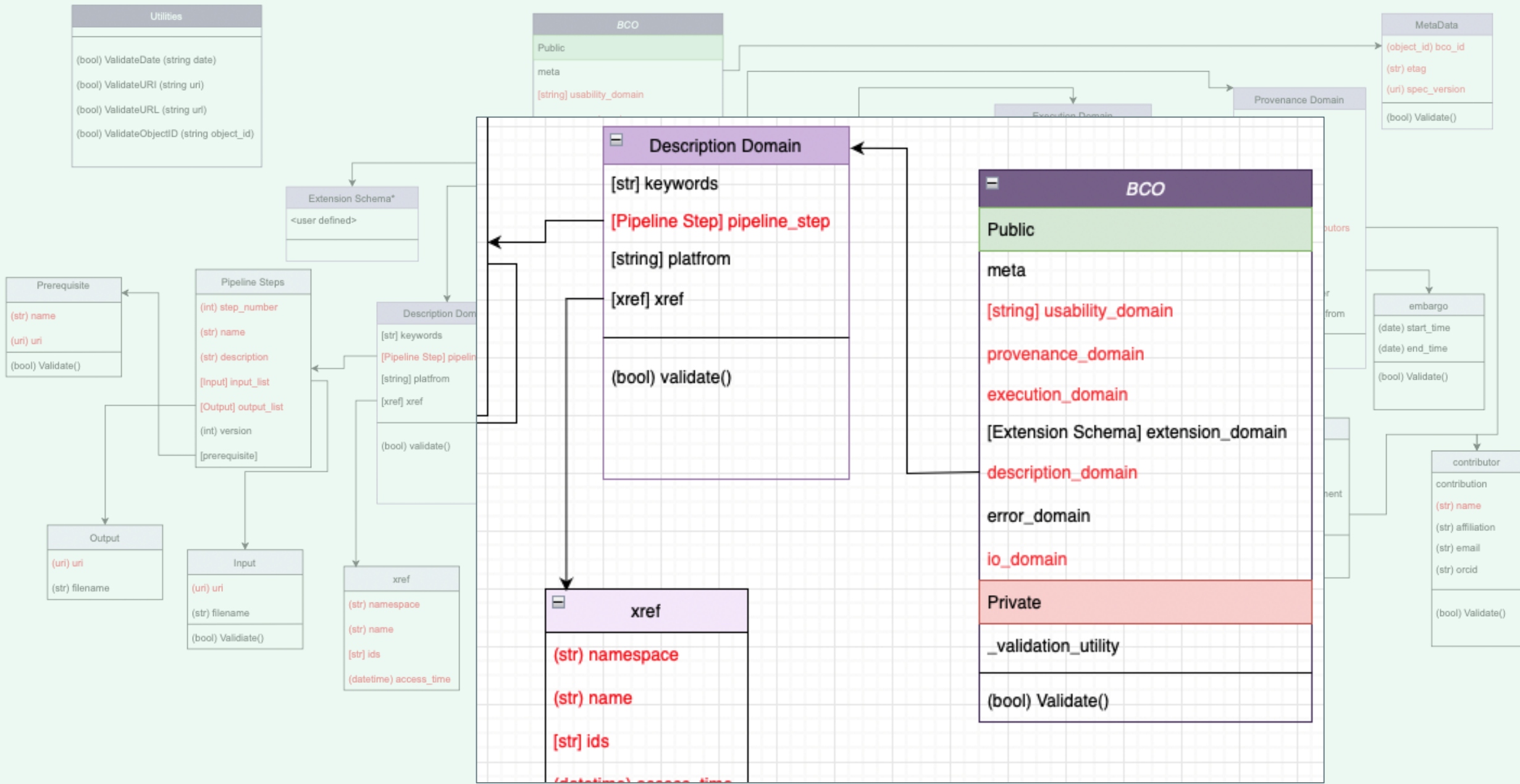
Previously...

- BCO exporter on HIVE
- Form-based editor for making BCOs
- Improved Documentation
- Combination BCOs with other projects
 - Research Objects/Nextflow
 - Common Workflow Language
- Reviewer's perspective on using BCOs
 - King et al.

Updates: BCO Python Library

- Importable package for Python
- Programmatically work with BCOs
 - Creates a Python Class Object from a BCO
 - Includes tools for working with the Objects

```
>>> import bco_
```



Example usage

- Retrieve publications for specific date range from all authors in Provenance Domain
- Sort studies by number of samples
- Check your environment for compatibility with a particular workflow
- Workflow validator: do file types line up correctly?
- Compare two workflows

BCO_001

1: Alignment

Hexagon

2: Variant Calling

Heptagon v1

AA profile: yes
min coverage: 12

3: Annotation

Annotation
mapper v1

dbSNP v136

4: Enrichment

GSEA

BCO_002

1: Alignment

BLAST

2: Variant Calling

Heptagon v2

AA profile: yes
min coverage: 10

3: Filtering

tblquery

4: Annotation

Annotation
mapper v1

dbSNP v137

Workflow Overlap

1: Alignment

2: Variant Calling

3: Annotation



Sort by

Score

Updated

Added

Name

Citation Count

Publication Date



Display as

Compact

Detailed

This site uses cookies. By continuing to browse the site you are agreeing to our use of cookies. Find out more here.

BLAST (EBI) |

Find regions of sequence similarity and alignments between a query sequence and database sequences.

Sequence analysis

Genomics

Sequence similarity search

Sequence alignment

Web API

Web application

Web service

BLAST

EBI Tools

Job Dispatcher Tools

BLAST |

A tool that finds regions of similarity between biological sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance.

Sequence analysis

Bioinformatics

Sequence similarity search

Sequence alignment

Command-line tool

Web application

Web API

BLAST

UCSC Genome Browser |

Large database of publicly available sequence and annotation data along with an integrated tool set for examining and comparing the genomes of organisms, aligning sequence to genomes, and displaying and sharing users own annotation data.

Model organisms

Genotype and phenotype

Data visualisation

Data integration and warehousing

Genomics

Rare diseases

Database search

Genetic variation analysis

Genome visualisation

Sequence alignment

PCR primer design

Web application

Database portal

Apache 2.0

Rare Diseases

Source for Tool Categories

This site uses cookies. By continuing to browse the site you are agreeing to

BLAST (EBI) |

Find regions of sequence similarity and alignments between a query sequence a

Sequence analysis Genomics

Sequence similarity search Sequence alignment

Web API Web application Web service BLAST EBI Tools Job I

BLAST |

A tool that finds regions of similarity between biological sequences. The pro
statistical significance.

Sequence analysis Bioinformatics

Sequence similarity search Sequence alignment

Command-line tool Web application Web API BLAST

UCSC Genome Browser |

Large database of publicly available sequence and annotation data along with
sequence to genomes, and displaying and sharing users own annotation data.

Model organisms Genotype and phenotype Data visualisation Data i

Database search Genetic variation analysis Genome visualisation Sequence alignment PCR primer design

Web application Database portal Apache 2.0 Deep Disease

EDAM - Bioscientific data analysis ontology

Last uploaded: February 26, 2021



Summary Classes Properties Notes Mappings Widgets

Jump to:

- ⊕ Data
- ⊕ DeprecatedClass
- ⊕ Format
- ⊖ Operation
 - ⊕ Alignment
 - ⊖ Analysis
 - ... Disease transmission analysis
 - ⊕ Enrichment analysis
 - ⊕ Expression analysis
 - ⊕ Genetic variation analysis
 - ⊕ Image analysis
 - ⊕ Network analysis
 - ⊕ Pathway analysis
 - ⊕ Phylogenetic analysis
 - ⊖ Protein function prediction
 - ⊕ Binding site prediction
 - ⊕ Molecular docking
 - ⊕ Peptide immunogenicity prediction
 - ⊕ Protein function comparison
 - Protein signal peptide detection**
 - ⊕ Protein-nucleic acid interaction analysis
 - ⊕ Protein-protein interaction analysis
 - ... Subcellular localisation prediction
 - ⊕ Sequence analysis
 - ⊕ Spectral analysis
 - ⊕ Structure analysis
 - ⊕ Text mining
 - ⊕ Transmembrane protein analysis
 - ⊕ Annotation
 - ⊕ Calculation
 - ⊕ Classification
 - ⊕ Clustering
 - ⊕ Comparison
 - ⊕ Conversion
 - ⊕ Correlation
 - ⊕ Data handling

Details Visualization Notes (0) Class Mappings (0)

Preferred Name	Protein signal peptide detection
Definitions	Detect or predict signal peptides and signal peptide cleavage sites in protein sequences. Methods might use sequence motifs and features, amino acid composition, profiles, machine-learned classifiers, etc.
ID	http://edamontology.org/operation_0418
comment	Methods might use sequence motifs and features, amino acid composition, profiles, machine-learned classifiers, etc.
Created in	beta12orEarlier
hasDefinition	Detect or predict signal peptides and signal peptide cleavage sites in protein sequences.
inSubset	http://purl.obolibrary.org/obo/edam#operations http://purl.obolibrary.org/obo/edam#edam
label	Protein signal peptide detection
prefixIRI	operation_0418
prefLabel	Protein signal peptide detection
subClassOf	Protein function prediction Protein feature detection



Search or jump to...



Pull requests Issues Marketplace Explore



biocompute-objects / bcotool Public

Edit Pins Unwatch 4 Fork 4 Star 2

Code Issues 5 Pull requests 1 Discussions Actions Projects Wiki Security Insights Settings

dev 4 branches 3 tags

Go to file Add file Code

About

This branch is 11 commits ahead, 3 commits behind main.

Contribute

skeene01 Merge pull request #20 from biocompute-objects/patch-2 4 days ago 57 commits

bcotool	Merge branch 'dev' into patch-2	4 days ago
data_tests	Fix links to cwl scripts	5 months ago
.gitignore	Update for conversion function	6 months ago
LICENSE	Update license	2 years ago
README.md	Update for conversion function	6 months ago
requirements.txt	adding requirements.txt	15 months ago
tst.json	Updates to API function	2 years ago

README.md

BCO-TOOL

<https://github.com/biocompute-objects/bcotool/tree/dev>

This is a Command Line Tool that allows for the manipulation of BioCompute Objects. Several functionalities are provided (detailed more below in supported modes).

To install:

Run the Git Clone command in the location you would like the repository:

No description, website, or topics provided.

- Readme
- MIT license
- 2 stars
- 4 watching
- 4 forks

Releases 3

1.2.0 Latest on Mar 3

+ 2 releases

Packages

No packages published Publish your first package

Contributors 4

- HadleyKing Hadley King
- acoleman2000 Alex Coleman
- skeene01
- rajamazumder



Provenance Domain

Usability Domain

Extension Domain (Optional)

Description Domain

Execution Domain

Parametric Domain (Optional)

I/O Domain

Error Domain (Optional)

Object Information:

Spec Version:

Prefix:

[CHECK PREFIX PERMISSION](#)

eTag:

PROVENANCE DOMAIN ?

*Name

*Version

*License

Created

Modified

Embargo

Start Time

End Time

*Review

Date	Status	*Reviewer Name	Reviewer Contribution	Reviewer Affiliation	Reviewer email	Reviewer ORCID	Comment	
<input type="text" value=""/> <input type="button" value="Calendar"/>	<input type="text" value=""/> <input type="button" value="Dropdown"/>	<input type="text"/>	<input type="text"/> <input type="button" value="Dropdown"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="button" value="Add"/>
								<input type="button" value="REMOVE"/>

*Contributors

Name	Contribution	Affiliation	eMail	ORCID	
<input type="text"/>	<input type="text"/> <input type="button" value="Dropdown"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="button" value="Add"/>
					<input type="button" value="REMOVE"/>

Portal Update



BioCompute Documentation

External site

- User Guide
- Best Practices
- SOP
- Tutorials

IEEE 2791-2020

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



BCO TSC

The Technical Steering Committee of the BioCompute Partnership (TSC) is a body of experienced professionals with BioCompute standard subject matter expertise. See here for the Meeting notes and agenda for all past and the upcoming meetings.

News and Events

FDA Notice on BioCompute

Electronic Submissions; Data Standards; Support for the International Institute of Electrical and Electronics Engineers Bioinformatics Computations and Analyses Standard for Bioinformatic Workflows.

<https://biocomputeobject.org>

Cloud-based tools for BioCompute

See our resources page for additional tools and services.



AWS instance of HIVE is temporarily down. Check back later for access.

Access AWS HIVE, the High-Performance Integrated Virtual Environment on AWS

BioCompute Builder



Use the BioCompute Builder or view objects in the database.

The BioCompute Builder is a platform-



BioCompute has been merged into the main Galaxy repository. This BioCompute enabled instance of Galaxy on AWS is therefore no longer operational. Thank

Tweets from @BioComputeObj

BioCompute Retweeted

Hadley King @HadleyKingIV · Aug 18

Tools from @SevenBridges @dnanexus and @galaxyproject will be featured! #BioCompute #BioComputeObjects

3 replies

BioCompute @BioComputeObj · Aug 18

A new BioCompute workshop will be held on September 13th via WebEx! Learn how to create BCO records from your #workflows directly on popular #bioinformatics platforms, and how to deposit them into an access-controlled database. Please register here: [workshop.org/#building-bio](#)





BioCompute Documentation

External site

- User Guide
- Best Practices
- SOP
- Tutorials

IEEE 2791-2020

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



BCO TSC

The Technical Steering Committee of the BioCompute Partnership (TSC) is a body of experienced professionals with BioCompute standard subject matter expertise. See here for the Meeting notes and agenda for all past and the upcoming meetings.

News and Events

FDA Notice on BioCompute

Electronic Submissions; Data Standards; Support for the International Institute of Electrical and Electronics Engineers Bioinformatics Computations and Analyses Standard for Bioinformatic Workflows.

object.biocompute@gmail.com

Cloud-based tools for BioCompute

See our resources page for additional tools and services.



AWS instance of HIVE is temporarily down. Check back later for access.

Access AWS HIVE, the High-Performance Integrated Virtual Environment on AWS

BioCompute Builder



Use the BioCompute Builder or view objects in the database.

The BioCompute Builder is a platform-



BioCompute has been merged into the main Galaxy repository. This BioCompute enabled instance of Galaxy on AWS is therefore no longer operational. Thank

Tweets from @BioComputeObj

BioCompute Retweeted



Hadley King

@HadleyKingIV · Aug 18

Tools from @SevenBridges @dnanexus and @galaxyproject will be featured! #BioCompute #BioComputeObjects

3 replies



BioCompute

@BioComputeObj · Aug 18

A new BioCompute workshop will be held on September 13th via WebEx! Learn how to create BCO records from your #workflows directly on popular #bioinformatics platforms, and how to deposit them into an access-controlled database. Please register here:



Updates: ARGOSDB

- Database of regulatory-grade genomes of infectious diseases
- Knowledgebase BCOs
- Document curation process

Updates: FDA-ARGOS

- Embleema and George Washington University's bioinformatics research collaboration to expand and develop database for high-quality sequences that can help researchers fight against infectious disease outbreaks
- This project will expand datasets publicly available in FDA-ARGOS, improve quality control by developing quality matrix tools and scoring approaches that will allow the mining of public sequence databases, and identify high-quality sequences for upload to the FDA-ARGOS database as regulatory-grade sequences.
- Building on expansions during the COVID-19 pandemic, this project aims to further improve the utility of the FDA-ARGOS database as a key tool for medical countermeasure development and validation.

FDA-ARGOS Project Outcomes

- Identify genomes of microbial species of high clinical relevance qualified as regulatory-grade sequences from public resources and generate annotation data model
- Develop more comprehensive and reliable quality control (QC) assessments of sequence representations
- Prepare NCBI submission packages and deposit regulatory-grade sequences to the FDA-ARGOS database and provide documentation, outreach, and training to FDA personnel

Search by: BCROID, dataset file name, title, description or categories



27 results found.



reviewed protein dataset ARGOS_000033 in FASTA format. [Salmonella typhimurium]

Salmonella typhimurium reference proteome sequences

```
>sp|Q9DB34-1|CHM2A_<ORG>
sp|Q9DB34-1|CHM2A_<ORG>
Vacuolar protein sorting-associated protein 2
OS= <Organism Name>
GN=Ccdc15 PE=1 SV=122
MDLLFGRRKTP
```

Salmonella typhimurium (strain LT2 / SGSC1412 / ATCC 700720) reference proteome fasta sequences.

reviewed protein dataset ARGOS_000002 in FASTA format. [Influenza A]

Influenza A reference proteome sequences

```
>sp|Q9DB34-1|CHM2A_<ORG>
sp|Q9DB34-1|CHM2A_<ORG>
Vacuolar protein sorting-associated protein 2
OS= <Organism Name>
GN=Ccdc15 PE=1 SV=122
MDLLFGRRKTP
```

Influenza A (A/Puerto Rico/8/1934 H1N1) reference proteome fasta sequences.

reviewed protein dataset ARGOS_000006 in CSV format. [SARS-CoV-2]

SARS-CoV-2 reference proteins list

uniprotkb_ac...	entry_name...
PODTC4...	VEMP_SARS2...
PODTDI...	R1AB_SARS2...

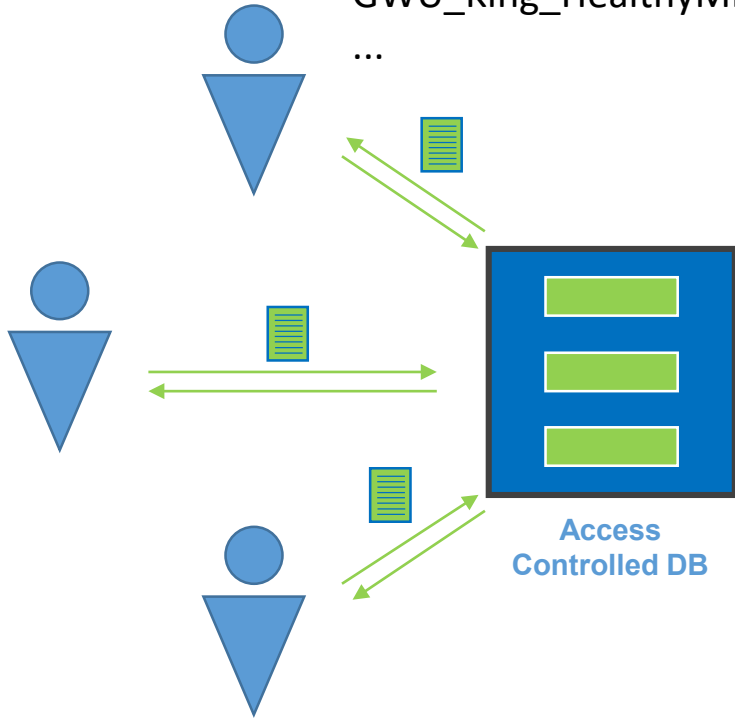
SARS-CoV-2 (Wuhan-Hu-1) reference protein accessions and summary annotations.

[View Details](#)

BioComputeDB

- Centralized registry of workflow Objects
- User-controlled based on Prefix
- Existing DBs at GW and FDA

BCOs:
CBER_HIVE_AAJ001.1
CBER_HIVE_ABC021.4
CBER_HIVE_FFX017.7
...

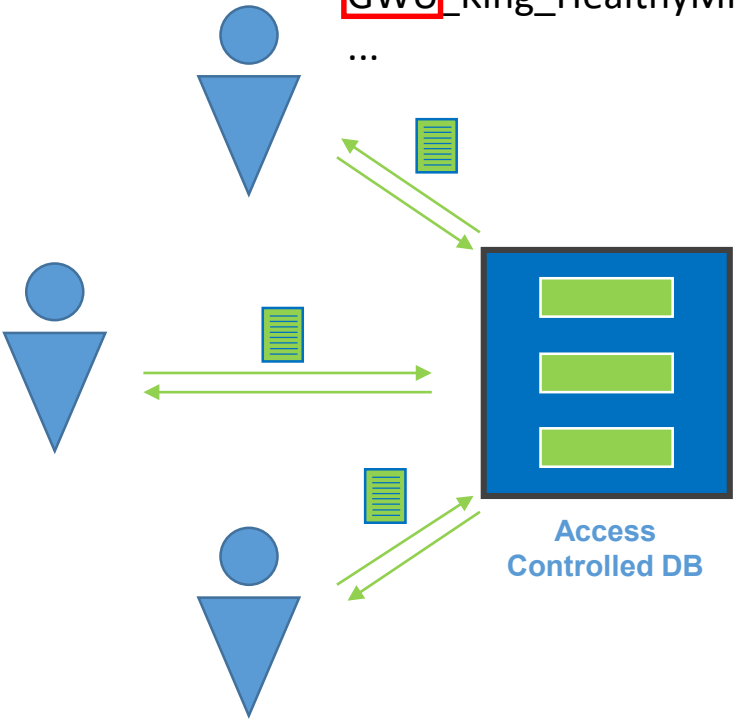


BCOs:
GWU_Vaccine_Keeney01
GWU_Vaccine_Keeney02
GWU_King_HealthyMicrobiomePipeline
...



BCOs:
NCI_Miller_Schuez_EtAl
NCI_Astling_And Smith
NCI_Erikson_EtAl
...

BCOs:
CBER_HIVE_AAJ001.1
CBER_HIVE_ABC021.4
CBER_HIVE_FFX017.7
...



BCOs:
GWU_Vaccine_Keeney01
GWU_Vaccine_Keeney02
GWU_King_HealthyMicrobiomePipeline
...



BCOs:
NCI_Miller_Schuez_EtAl
NCI_Astling_And Smith
NCI_Erikson_EtAl
...

Acknowledgements



BCO Founding Members:

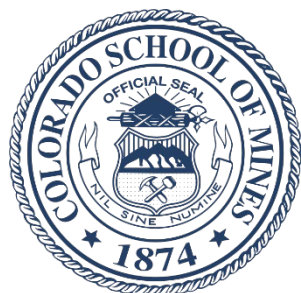
Raja Mazumder

Vahan Simonyan

- Konstantinos Karagiannis
- Mark Walderhaug
- Anton Golikov
- Hadley King
- Tianyi Wang
- Rohan Panigrahi
- Sean Keeney
- Lam Phuc
- John Torcivia-Rodriguez
- Alhanouf Altuwayjiri
- Michael Crusoe
- Stian Soiland-Reyes



The University of Manchester



BAA: 75F40119C10136



Workshop Schedule

10:30 - 12:00PM	BioCompute and Galaxy
	Introduction Updates Questions BioComputeDB Concept, demo functionality with Portal, API, prefixes Questions Galaxy Demo BCO builds on the platform Questions
12:00 - 12:30PM	Lunch Break
12:30 - 2:00PM	DNAexus
	Demo BCO nexus interface to present overview, in-memory editor, schema/standard driven, WDL, DNAexus integration Questions Demo use cases Build from scratch, start with a JSON file BCO metadata template and layer in workflow metadata from WDL and DNAexus; export BCO to DNAexus and run the workflow Questions
2:00 - 2:15PM	Break
2:15 - 3:45PM	Seven Bridges
	Demo Seven Bridges tool interface Questions Demo use cases Generate BCOs from platform output Questions
3:45 - 4:00PM	Open Discussion Future plans and feedback