



BioCompute Database and Transfer Mechanism Development Workshop

Raja Mazumder, Ph.D.

Principal Investigator

Professor, GW

Chair, BioCompute Executive Steering Committee

mazumder@gwu.edu

Hadley King, M.S.

Operational Lead

Chair, BioCompute Technical Steering Committee

hadley_king@gwu.edu

Jonathon Keeney, Ph.D.

Co-Investigator

Assistant Research Professor, GW

Managing Director, BioCompute Executive Steering Committee

keeneyjg@gwu.edu

Janisha Patel, M.S.

Training & Outreach Lead

Technical Writer

janishapatel@gwu.edu

Workshop Guidelines

- If you have questions during the talk, please type them into the chat.
- Q&A
 - Moderator: Charles Hadley King
 - Use “raise hand” feature during Q&A session to ask a question.
- This workshop will be recorded, and the recording will be sent to all attendees shortly after.

Thank you!

Goals of this Workshop

1. Introduce BioCompute Objects (BCO) for computational analysis
2. Provide BioCompute resources for future reference
3. Demonstrate tools available via BioCompute Portal and BCODB

Agenda

- Introduction: Use cases and **BioCompute**
- BioCompute Portal Walkthrough
- Demo of user account and DB access
- Description of DB and schema
- Transfer from Galaxy, HIVE, & local machine
- Q&A



Use-Case Examples

Test
Submission

Tuberculosis
Detection

Vaccine Safety

Provenance Domain

Name R Safety Assessment Algorithm for Aluminum in Infant Vaccines.

Version 1.0

License <https://creativecommons.org/licenses/by/4.0/> 

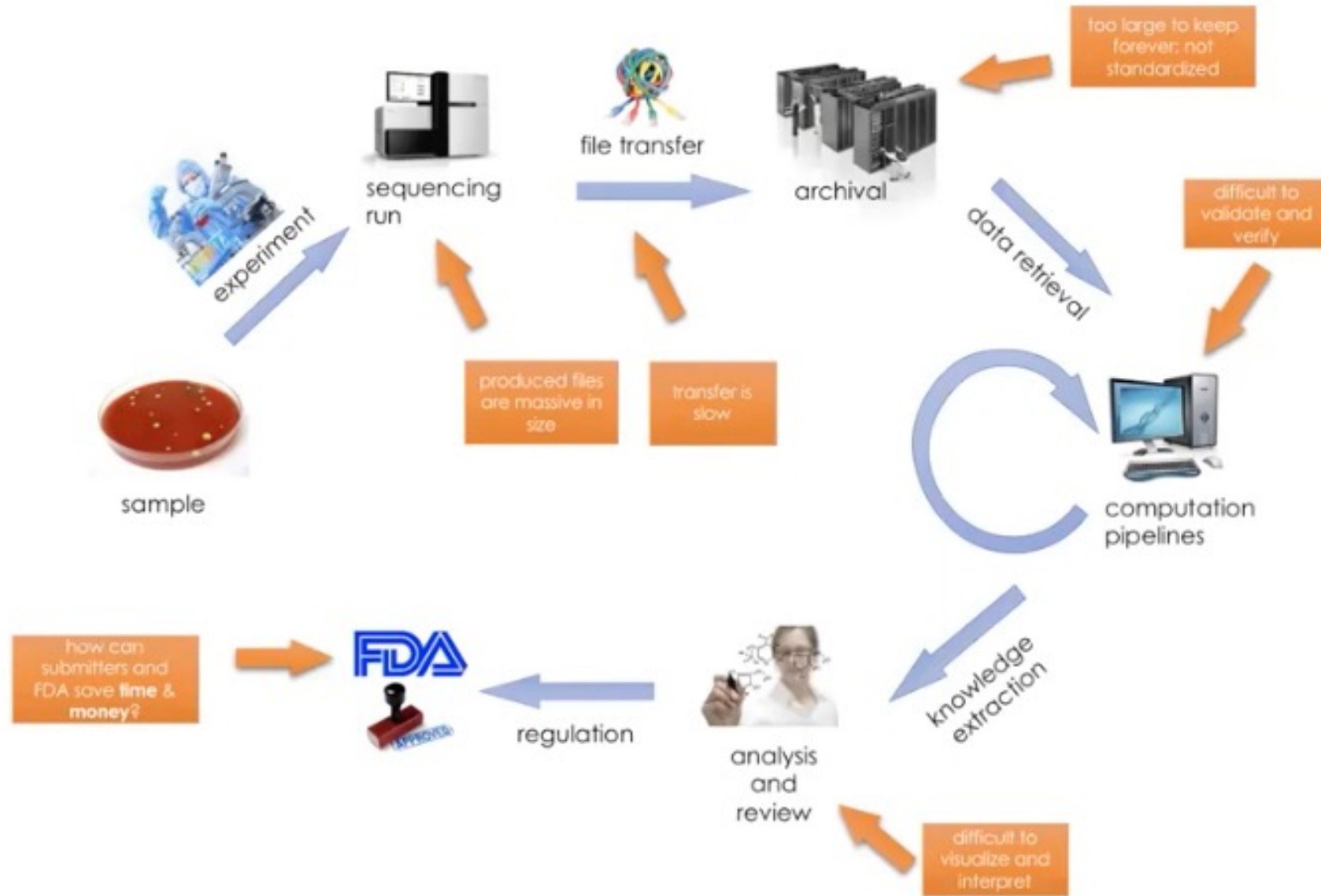
Created 2020-04-30T18:03:25.679Z

Modified 2021-12-03T19:23:45.074Z

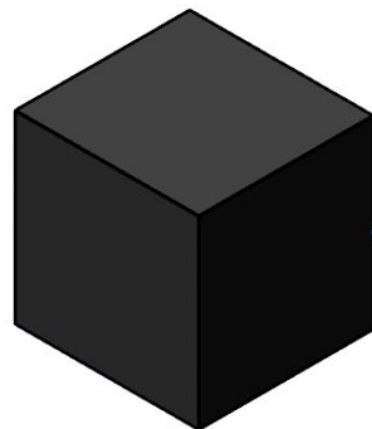
Contributors

Name	Contribution	Affiliation	eMail
Charles Hadley King	createdBy, curatedBy	George Washington University	hadley_king@gwu.edu
Mark Walderhaug	curatedBy, curatedBy, authoredBy	U.S. Food and Drug Administration	mark.walderhaug@fda.hhs.gov
RJ Mitkus	authoredBy	U.S. Food and Drug Administration	rj.mitkus@fda.hhs.gov

NGS life cycle



NGS Data Flows



```
$ fastq-dump -X 2 SRR001666 --split-3
W: $ fastq-dump -X 2 SRR001666 --split-3
R: $ fastq-dump -X 2 SRR001666 --split-3
N: $ fastq-dump -X 2 SRR001666 --split-3
+
@: $ fastq-dump -X 2 SRR001666 --split-3
-
: Read 2 spots for SRR001666
G: Written 2 spots for SRR001666
+
I: $ head SRR001666_1.fastq SRR001666_2.fastq
@: ==> SRR001666_1.fastq <==
+
I: @SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
+
: GGCTGATGCCCGCTGCCGATGGCGTCAATCCACC
+
I: +SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
+
: IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII9IG9IC
+
I: @SRR001666.2 071112_SLXA-EAS1_s_7:5:1:801:338 length=36
+
: GTTCAGGGATACGACCTTTGTATTTAAGAAATCTGA
+
I: +SRR001666.2 071112_SLXA-EAS1_s_7:5:1:801:338 length=36
+
: IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII6IBI
+
I: ==> SRR001666_2.fastq <==
+
I: @SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
+
: AAGTACCCTTAACAACCTAAGGGTTTTCAAATAGA
+
I: +SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
+
: IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII/
+
I: @SRR001666.2 071112_SLXA-EAS1_s_7:5:1:801:338 length=36
+
: AGCAGAAGTCGATGATAATACGCGCTTTTATCAT
+
I: +SRR001666.2 071112_SLXA-EAS1_s_7:5:1:801:338 length=36
+
: IIIIIIIIIIIIIIIIIIIIIIIIIIGI>IIIII-I)8I
```

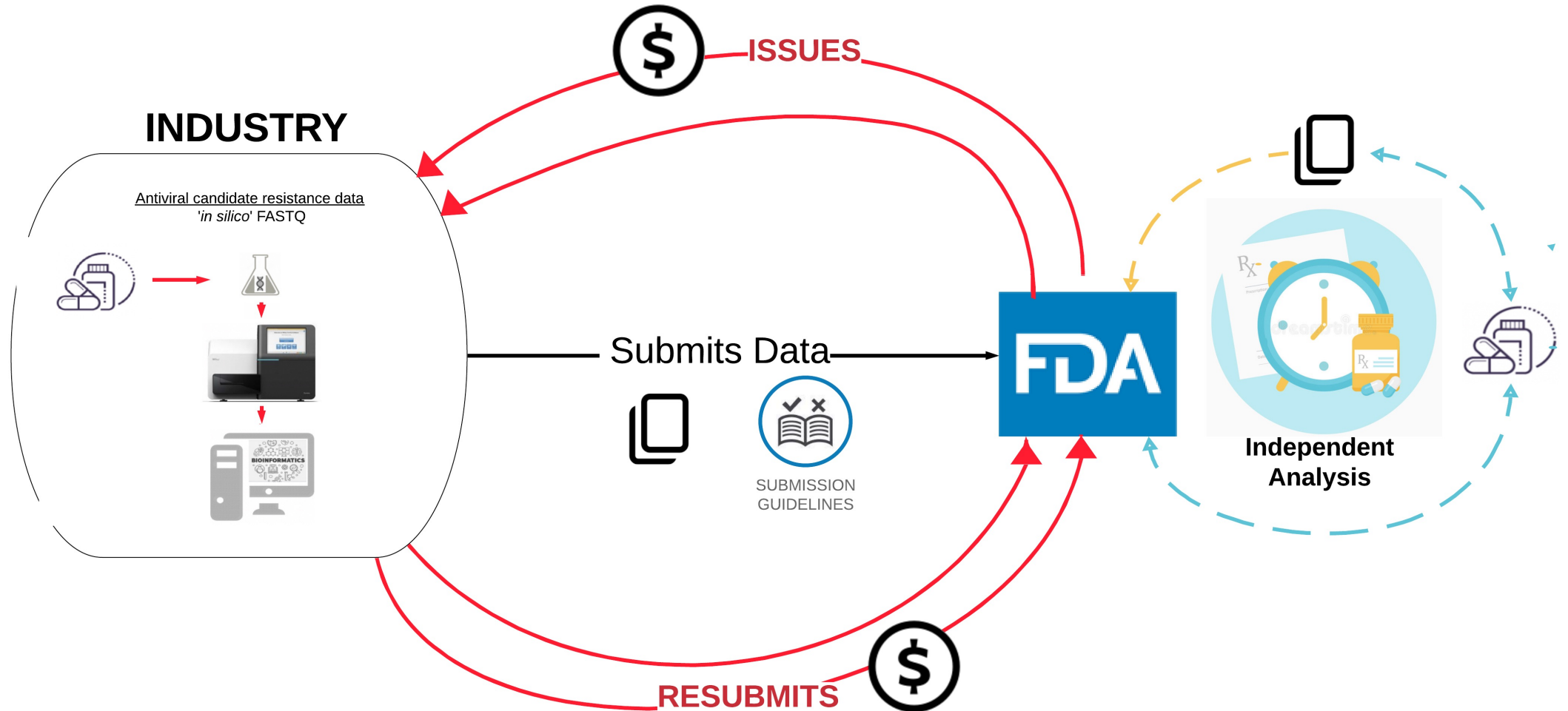
- Ancestry
- Cancer
- Microbiome
- Disease correlation
- Agriculture
- Synthetic biology
- Livestock
- Metagenomics
- Personalized medicine

Introduction to BioCompute

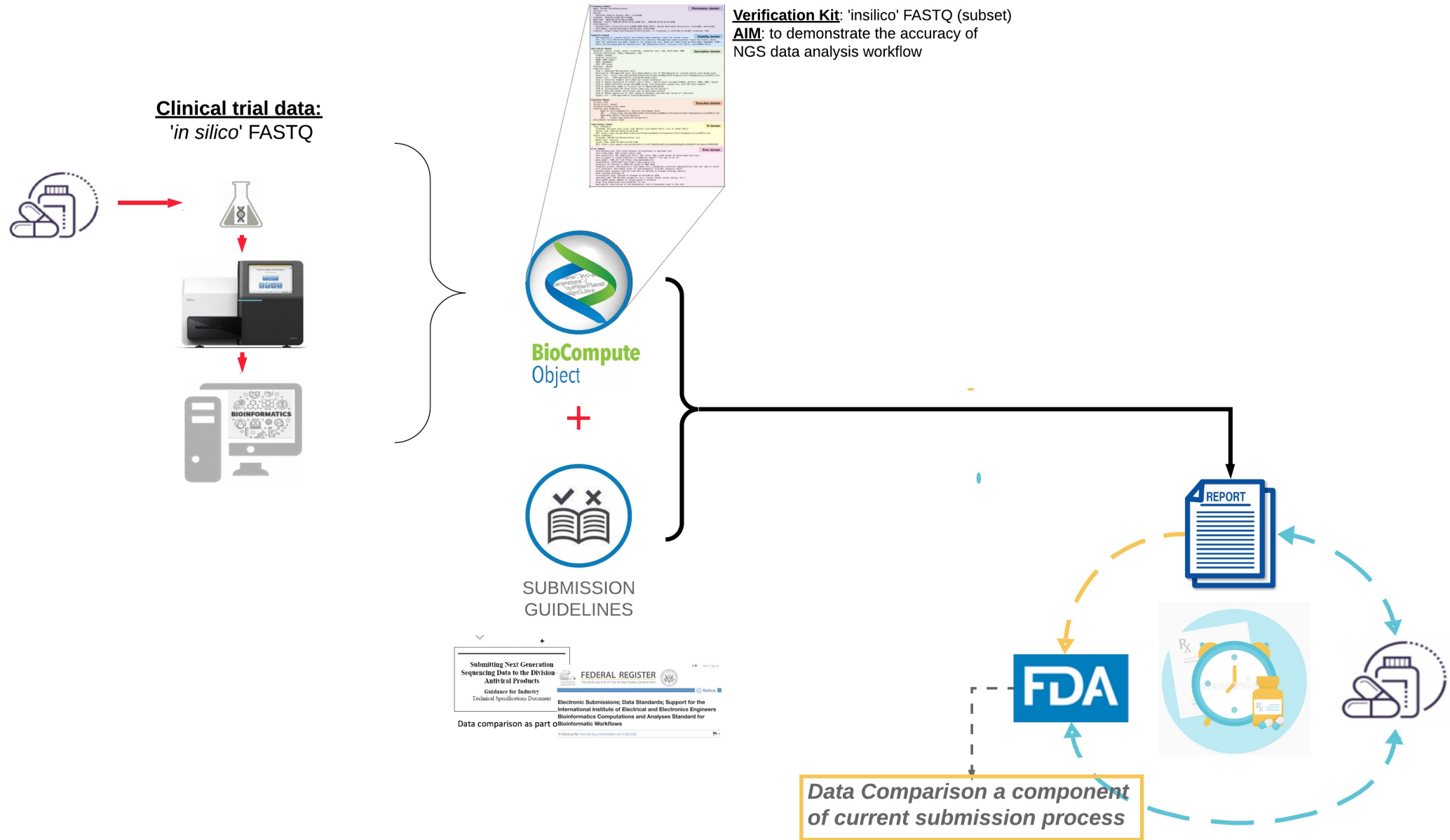


BCO Value in Regulatory Submission

Parameters? Tool versions?
Data provenance?
Reproducibility issues?
Missing pipeline steps?



BCO Value in Regulatory Submission



A solution should...

In 2014, the Genomics Working Group convened a special session to discuss the problem, and came up with four key focus areas

1. Be **human readable**: like a GenBank sequence record
2. Be **machine readable**: structured information with predefined fields and associated meanings of values
3. Contain enough information to understand the computational pipelines, interpret information, maintain records, and reproduce experiments
4. Be **immutable**: ensure information has not been altered

Written in JSON

Categorized by domains

- Adaptable to other schemas

Adheres to and encourages F.A.I.R. principles

- Fully open source
- Preserves data provenance (unique IDs for versioning)

<p>Top Level</p> <p>BCD ID: https://w3id.org/biocompute/1.3.0/examp1ns/FDA-NA-TestsBreastCancer Checksum: 06DACE78679F35BA87A3D06FFED4ED24A4F588C2571264C37E5F1B3ADE8A431 Specification: https://w3id.org/biocompute/1.3.0/</p>	Metadata
<p>Provenance Domain</p> <p>Name: FDA-NA-TestsBreastCancer Version: 1.0 Review: approved: Natalie Abrams, NIH ; createdBy Created: 2018-05-24T09:48:17-0500 Modified: 2018-06-22T14:06:14-0400 Embargo: Start: 2008-09-26T14:43:43-0400 End: 2008-09-26T14:43:45-0400 Contributors: Janisha Patel (http://orcid.org/0000-0002-8824-4637), George Washington University; createdBy, modifiedBy Dara Baker, George Washington University; authoredBy License: https://spdx.org/licenses/CC-BY-4.0.html -> Licensing is inferred by OncoMX Licensing. Pub=</p>	Parametric domain
<p>Usability Domain</p> <p>FDA-approved or cleared nucleic acid-based human biomarker tests for breast cancer The .xlsx file FDA-NA-TestsBreastCancer.xlsx contains FDA-approved human biomarker tests for breast cancer. Each row represents one gene linked to its respective test. Genes are identified by UniProtKB, HGNC name Tests are distinguished by manufacturer, FDA submission ID(s), clinical trial ID(s) and PubMed ID(s).</p>	Usability domain
<p>Extension Domain</p> <p>Dataset Extension: Comment: Unique column headers for the dataset test_disease_use: FDA-listed disease corresponding to approved test test_trade_name: FDA-listed product name test_manufacturer: FDA-listed patent company for the approved test test_submission: FDA submission ID(s), web links; FDA-listed patent ID associated with test test_is_panel: A single biomarker or biomarker panel? Y for yes, N for no gene_symbol: HGNC ID from https://www.genenames.org uniprotKB_ac: UniProtKB from https://www.uniprot.org biomarker_id: Matched to EDN IDs based on HGNC Name biomarker_origin: Characteristic that makes this a biomarker; molecular abnormalities that can lead to cancer ncit_biomarker: Searchable terms for gene/Biomarker from NCI Thesaurus (NCIT)</p>	Extension domain
<p>Description Domain</p> <p>Keywords: cancer, breast cancer, biomarker, biomarker test, FDA, UniProtKB, EDN External References: (Name, Namespace, Ids) PubMed: pamed; UniProt: accession; EDN: EDN number; HGNC: hgncName; GTR: GTR terms; Platform: Manual Pipeline Steps: Step 1: Download FDA-approved tests Description: FDA approved tests were downloaded a list of FDA approved or cleared nucleic acid based tests Input List: https://www.fda.gov/MedicalDevices/ProductsandMedicalProcedures/InVitroDiagnostics/ucn330711.htm Output List: ~/FDA-approved-or-cleared-NAC-based-tests</p>	Description domain
<p>Execution Domain</p> <p>Scripts: none Script Driver: manual Software Prerequisites: None External Data Endpoints: Name In Vitro Diagnostics > Nucleic Acid Based Tests URL: https://www.fda.gov/MedicalDevices/ProductsandMedicalProcedures/InVitroDiagnostics/ucn330711.htm Name NCI Genetic Testing Registry URL: https://www.ncbi.nlm.nih.gov/gtr/ Environment: Variables: None</p>	Execution domain
<p>Parametric Domain</p> <p>N/A</p>	Parametric domain
<p>Input/Output Domain</p> <p>Input Subdomain: Filename: Multiple test files from "Nucleic Acid Based Tests: List of Human Tests" Access Time: 2018-10-10T11:34:02-5:00 URL: https://www.fda.gov/MedicalDevices/ProductsandMedicalProcedures/InVitroDiagnostics/ucn330711.htm Output Subdomain: Filename: FDA-NA-TestsBreastCancer.xlsx Media Type: xlsx/csv Access Time: 2018-10-10T11:37:02-5:00 URL: https://docs.google.com/spreadsheets/d/1x1UY7WJNE7HyCgH5sYpxEuq4btgVIRhwgF2oc81WFH28Y/edit#gid=1492826383</p>	IO domain
<p>Error Domain</p>	Error domain

Provenance Domain

Usability Domain

Extension Domain

Description Domain

Execution Domain

Parametric Domain

IO Domain

Error Domain

IEEE Standard



Institute of Electrical and Electronics
Engineers Standard

IEEE 2791-2020 approved January 2020

<https://standards.ieee.org/content/ieee-standards/en/standard/2791-2020.html>



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

PUBLISHED DOCUMENT

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

ACTION:

DOCUMENT

Printed ve
PDF

Publicatio
07/22/20

Agencies:

PUBLISHED DOCUMENT

AGENCY:

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

ACTION:

Notice.

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 include this standard in the FDA Data Standards Catalog for the submission of high-throughput sequencing (HTS) data in new drug applications (NDAs), abbreviated new drug applications (ANDAs), biologics license applications (BLAs), and investigational new drug applications (INDs) to the Center for Biologics Evaluation and Research (CBER), Center for Drug Evaluation and Research (CDER), and Center for Food Safety and Applied Nutrition (CFSAN).

DATES:

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

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

Document Citation:

85 FR 44304

Page:

44304-44305 (2 pages)

Agency/Docket Number:

Docket No. FDA-2020-N-1450

Document Number:

2020-15771

DOCUMENT DETAILS

ACCESS: Private | NAME: test-workflow | ORG: dnanexus.science | ADDED BY: sam.westreich | ID: workflow-FQ7P7Vj05922F6k6J3b87yQ6

CREATED: 2018-12-10 23:16:23

Edit tags

Revision: 1 | Latest | Edit | Fork | Export | Run Workflow rev1

SPEC | WORKFLOW DIAGRAM

INPUTS

file	Input 1	REQUIRED	workflow-app-1
file	Input 2	REQUIRED	workflow-app-2

OUTPUTS

file	Output 1	REQUIRED	workflow-app-1
file	Output 2	REQUIRED	workflow-app-2



Projects | Data | Apps

Identifiers and File name(s) | Search | P Queries | Save Query | Copy files to project

Start Query From:

- Case
- File
- Sample
- Portion
- Slide
- Analyte
- Aliquot
- Drug therapy
- Radiation therapy
- Follow up
- New Tumor Event

File (ADD FILTER)

Data Format (Remove filters)

Experimental Strategy (Remove filters)

Disease Type (Remove filters)



Man Home HIVE Portal Links

CensusScope

Loading	Status
Building histogram	Done 100%
Preparing alignments	Done 100%
Visualizing alignments in track	Done 100%
Fetching alignments	Done 100%
Creating mutation heat diagram	Done 100%

Taxonomy Help | Taxonomy Details

Parameters

Progress

Results

Taxonomy Details

Convergence

Phylogenetic Tree

Text Tree

Table

Subtree

What's Next?

Alignment



Galaxy Administration

Galaxy Administration | Analyze Data | Workflow | Shared Data | Admin | Help | User | Using 35.7

Administration

- Security
 - Manage users
 - Manage groups
 - Manage roles
- Data
 - Manage quotas
 - Manage data libraries
- Server
 - Reload a tool's configuration
 - Profile memory usage
 - Manage jobs
 - Manage installed tool shed repositories
- Tool sheds
 - Search and browse tool sheds
- Form Definitions
 - Manage form definitions
- Sample Tracking
 - Manage sequencers and external services
 - Manage request types
 - Sequencing requests
 - Find samples

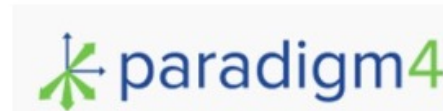
Repository Actions | Tool Shed Actions

Genome/Exome paired analysis (SNVMix1)

Boxes are red when tools are not available in this repository (this page displays SVG graphics)



BioCompute participants



Agenda

- Introduction to BioCompute
- **BioCompute Portal Walkthrough**
- Demo of user account and DB access
- Description of DB and schema
- Transfer from Galaxy, HIVE, & local machine
- Q&A



BioCompute Portal Walkthrough



Demo

How to navigate the Portal?

How to create a BioCompute Object?

How do you edit a BioCompute Object?



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.


Cloud-based tools for BioCompute

See our resources page for additional tools and services.



Access AWS HIVE, the High-Performance Integrated Virtual Environment, on AWS. HIVE is a cloud-based environment optimized for the storage and analysis of extra-large data, such as biomedical data, clinical data, next-generation sequencing (NGS) data, mass spectrometry files, confocal microscopy images, post-market

BioCompute Builder



Use the BioCompute Builder or view objects in the database. The BioCompute Builder is a platform-free, form-based editor. The builder walks a user through building a BCO through text boxes, indicating which entries are required to adhere to the IEEE standard.



Use Galaxy on AWS, the open source, web-based platform for data intensive biomedical research. Assemble your [pipeline] in the workspace, designate the outputs in the module boxes, and record the entire pipeline as a BCO.

Tweets by @BioComputeObj

BioCompute Retweeted

 **GW SMHS** @GWSMHS
A consortium led by @GW_HIVE_Lab, @rmazumde, and @Embleema, with @TempleUniv, has received @FDA funding to advance the work against #infectiousdiseases. Read more: bit.ly/FaDCvY via @TheVDT


Embleema and George Washington University-led Co... METUCHEN, N.J.--(BUSINESS WIRE)--Oct 28, 2021-- valdetadailytimes.com

Nov 15, 2021

Portal home

Demo



User account and DB access





Demo

- Register a new account
- Build an Object
- Create a draft
- Edit a draft
- Save edits
- Validate Object
- Publish



Description of DB and Schema



BioCompute Schema Files

IEEE.org | IEEE Xplore Digital Library | IEEE Standards | IEEE Spectrum | More Sites

<https://opensource.ieee.org/2791-object/ieee-2791-schema/>

IEEE SA OPEN Menu

ieee-2791-schema

2791 object > ieee-2791-schema

Project information

Repository

Issues 1

Merge requests 1

CI/CD

Deployments

Monitor

Packages & Registries

Analytics

Wiki

Snippets

ieee-2791-schema
Project ID: 116

24 Commits 2 Branches 3 Tags 276 KB Files 276 KB Storage 1 Release

master ieee-2791-schema

History Find file Download Clone



Update README.md

Joshua Gay authored 1 year ago

45683af9



README Other Auto DevOps enabled

Name	Last commit	Last update
.gitignore	Creates initial release of BioCompute Obje...	2 years ago
2791object.json	replaces https://w3id.org/2791/ with https://w3id.org/2791-object/	1 year ago
AUTHORS	Update AUTHORS	1 year ago
CONTRIBUTORS	Update CONTRIBUTORS	1 year ago
LICENSE	Update LICENSE	1 year ago
README.md	Update README.md	1 year ago
description_domain.json	replaces https://w3id.org/2791/ with https://w3id.org/2791-object/	1 year ago
error_domain.json	replaces https://w3id.org/2791/ with https://w3id.org/2791-object/	1 year ago

Demo



Repositories

Transfer from external sources using Swagger site for API access





<https://biocomputeobject.org/api/docs/?format=openapi>

Explore

BioCompute Object Data Base API (BCODB API) 1.3.0

[Base URL: biocomputeobject.org/]
<https://biocomputeobject.org/api/docs/?format=openapi>

A web application that can be used to create, store and edit BioCompute objects based on BioCompute schema described in the BCO specification document.

- [Terms of service](#)
- [Contact the developer](#)
- MIT License

Schemes
HTTPS ▾

[Django Login](#) [Authorize](#) 🔒

Filter by tag

Account Management

GET	<code>/api/accounts/activate/{username}/{temp_identifier}</code> Activate an account	<code>api_accounts_activate_read</code> 🔒
POST	<code>/api/accounts/describe/</code> Account details	<code>api_accounts_describe_create</code> 🔒
POST	<code>/api/accounts/new/</code> Account creation request	<code>api_accounts_new_create</code> 🔒

Group Management

POST	<code>/api/groups/create/</code> Create group	<code>api_groups_create_create</code> 🔒
------	---	---

Swagger site for
API access ✓

Demo



Transfer from Galaxy



Thank you!

Your time and feedback are greatly appreciated.



Acknowledgments



Konstantinos Karagiannis
Eric Donaldson
Mark Walderhaug
Carolyn Wilson
Anton Golikov



Chris Armstrong

The logo for SevenBridges is a dark blue rectangle with the text 'SevenBridges' in white.

SevenBridges

Dennis A. Dean
Jeffrey Grover
Soner Koc

The logo for The George Washington University features the text 'THE GEORGE WASHINGTON UNIVERSITY' in blue, with horizontal lines above and below the text.

THE GEORGE
WASHINGTON
UNIVERSITY

WASHINGTON, DC

Raja Mazumder
Jonny Torcivia

The logo for DNAexus features the text 'DNAexus' in black, with 'DNA' in a larger font and 'exus' in a smaller font, and a blue 'x' in the middle.

DNAexus®

Omar Serang
John Didion



Vahan Simonyan
Jeremy Goecks
Gil Alterovitz
Carole Goble
Jonas Almeida
Dan Taylor
Ntino Krampis
Michael Crusoe
Stian Soiland-Reyes
Konstantinos Krampis
Elaine Thompson
Nicola Soranzo
Jason Travis
Nan Xiao

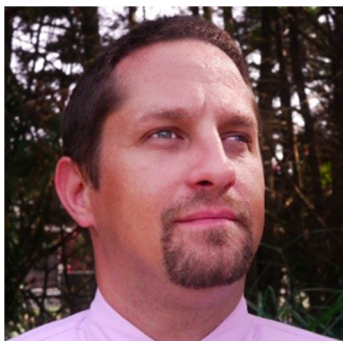
Acknowledgements and Contact



Raja Mazumder, Ph.D., PI
Professor
The George Washington University
mazumder@gwu.edu



Jonathon Keeney, Ph.D., Co-I
Assistant Research Professor
The George Washington University
keeneyjg@gwu.edu



Charles Hadley King
Technical Lead
The George Washington University
hadley_king@gwu.edu



Janisha Patel
Outreach & Training Lead
The George Washington University
janishapatel@gwu.edu

Recording on workshop: <youtube link>

- Introduction: Use cases and BioCompute
- BioCompute Portal Walkthrough
- Demo of user account and DB access
- Description of DB and schema
- Transfer from Galaxy, HIVE, & local machine
- Q&A



Q&A

