# Starting Cloud-Based Research on NHLBI BioData Catalyst® (BDC)

BDC for NDSA Workshop

**Day 1: July 13  |  11 AM ET**

**NIH** National Heart, Lung, and Blood Institute | **BioData** **CATALYST**®

# Statement of Conduct

The BDC Consortium is dedicated to **providing a harassment-free experience for everyone**, regardless of gender, gender identity and expression, age, sexual orientation, disability, physical appearance, body size, race, or religion (or lack thereof). We do not tolerate harassment of community members in any form. Sexual language and imagery is generally not appropriate for any venue, including meetings, presentations, or discussions.

**Web Resource**: Statement of Conduct

BioData CATALYST

# Agenda: Day 1

| Topic | Time |
| --- | --- |
| **Welcome! Introductions, Housekeeping, and Icebreaker** | 30 min |
| **BDC's Role in Data Science + Discussion**<br>→ Present Day Data Science Challenges<br>→ BDC's Approach<br>→ Future of Data Science and discussion | 1 hour |
| **30 MINUTE LUNCH BREAK @ 12:30 ET** | |
| **Using BDC in a Research Project**<br>→ Data Discovery and Exploration<br>→ Analysis on BDC | > 2 hours |
| **Closing & Recap + Working Session** | 40 min |

# Agenda: Day 2

| Topic | Time |
|---|---|
| **Reproducible Research** | 1.5 hour |
| **30 MINUTE LUNCH BREAK @ 12:40 ET** | |
| **Using BDC as a teaching platform**<br>→ Use Cases and Guest Speaker with Q&A<br>→ Benefits and Challenges of Cloud Teaching with discussion | > 1.5 hour |
| **Workflows and Cost Estimation**<br>→ Sign up for $500 in Pilot Credits | > 1 hour |
| **Closing & Recap** | 40 min |

# Pre-Work Reminder

**Step 1:
Join the BDC Community**

Then, introduce yourself on the forum and download today's slides!
bit.ly/BDC-NDSA-Forum

**Step 2:
Sign up for Seven Bridges**

**I _have_ eRA Commons**

Create an account using eRA Commons

**Create Seven Bridges Account**

**I _don't have_ eRA Commons**

Follow the instructions posted in the forum

**Alternative Login Instructions**

NIH National Heart, Lung, and Blood Institute | BioData CATALYST

# Introductions

# Meet Your Instructors

**Emily Hughes**
*BDC Powered by PIC-SURE*
Harvard Medical School

**Kat Thayer**
*BDC Powered by Terra*
Broad Institute

**David Roberson**
*BDC Powered by Seven Bridges*
Seven Bridges/Velsera

**Justin Dorsheimer**
*BDC Powered by Gen3*
University of Chicago

# TAs and Live Support

**Cera Fisher**
*BDC Powered by Seven Bridges*
Seven Bridges/Velsera

**Aarthi Krishnan**
*BDC Powered by Seven Bridges*
Seven Bridges/Velsera

**Michael Corace**
*BDC Powered by PIC-SURE*
Harvard Medical School

**Amber Voght**
User Engagement Specialist,
BDC Coordinating Center

**Kaleena Narwani**
User Engagement Specialist,
BDC Coordinating Center

NIH National Heart, Lung, and Blood Institute | BioData CATALYST®

# BDC's Role
# in Data Science

# Curricula of:

## BDC's Role in Data Science

**Learning goals**

- Acknowledge the state of the data science landscape
- Assess BDC's impact on field of data science
- Consider the trajectory of where data science is going
- Discover ways to support your and your students' science data interests using BDC

# Present Day
# Data Science Challenges

Perspective

# Industry Representation

**Citizenship**

73.3% US Citizens

**Top 3 areas of Employment**
1. Genetic Counseling (45.7%)
2. Research (30.4%)
3. Academic (23.4%)

**Employment Status**

78% Permanent Positions

**Age**

45.1% 25-40 years old

**Gender Identity and Sexual Orientation**

74.7% Women

23.3% Men

0.5% Nonbinary/Transgender

6.9% LGBTQIA

**Disability Status**

3.4% Reported having a disability

ASHG: Human Genetics & Genomics Workforce Survey

# Industry Representation

Race, Ethnicity, & Ancestry

| Category | Percentage |
|---|---|
| American Indian or Alaskan Native | <1% |
| Native Hawaiian or other Pacific Islander | <1% |
| Middle Eastern or North African | 1.1% |
| Black, African American, or African | 1.5% |
| Hispanic, Latino, or Spanish | 2.0% |
| Multiracial | 4.8% |
| Asian | 7.4% |
| White | 67.0% |
| Not identified | 15.4% |

ASHG: Human Genetics & Genomics Workforce Survey

# Data Representation



- 0.1% genomic differences come from variations among ~3 billion bases in our DNA
- Most permitted DNA samples used in research come from people with European Ancestry (78%)
- Active Initiatives: Polygenic Risk Methods in Diverse Populations Consortium, Electronic Medical Records and Genomics Network, and the Human Pangenome Reference Consortium

Genome.gov: Diversity in Genomics Research Fact Sheet

# The rate of data generation is accelerating rapidly

**Doubling Time of Health Knowledge**



- More biomedical data will be generated this year than all previous years **combined**

- Diverse data modalities including EHR data, Survey, Sequencing, Transcriptomics, Metabolomics, Proteomics, Imaging, Sensor, E-Phys, Flow Cytometry, and so on

# Scalability

- Capturing data
    - Size of raw data
    - Funding considerations
    - Saving information

- Analyzing data
    - Computing time and power
    - Comparing data sets and harmonization

- Sharing results
    - Can be slow to impact other fields

Sources - Genome.gov:
    Cost of Sequencing a Human Genome
    Genomic Data Science Fact Sheet

# Reproducibility

- Data harmonization and documentation

- Use diverse, yet comparable datasets

- Working in silos
  - Hyper local
  - Breaking through the noise

- Study errors
  - Type I Errors - false positives and generalizability

# Access

- IRB and data sharing

- Data privacy laws
  - California Privacy Law

- Special omics considerations
  - Patient De-identification
  - Secondary data scope allowance

# Case Study: Giraffe

# Brief Discussion

Of the challenges discussed, what resonated with you most?
What other challenges come to mind?

# BDC's Approach

# BDC

# Mission   Vision



The *mission* is to develop and integrate advanced cyberinfrastructure, leading edge tools, and FAIR data to support the NHLBI research community.

The *vision* is to be a community-driven ecosystem implementing data science solutions to democratize data and computational access to advance Heart, Lung, Blood, and Sleep science.

# Using the Cloud to store and analyze growing health data

- Immediate scaling -- no need to wait to purchase and install hardware.

- Levels the playing field -- even researchers at institutions without large compute infrastructure investments can access powerful data and compute resources.

- Many researchers can access data without needing to physically copy it.

- Data and methods in a single place streamlines reproducibility.



Old model: send data to compute

New model: send compute to data

# What BDC offers

**Managing the Computing Environment**

Elastic Computing

**Easier Access to many High Value Datasets**

**Tooling**

Data Discovery

Statistical Analysis Tools (R, SAS)

Other Specialized Workflows

**Community and Peer Interactions**

# The Computing Environment

No need to **download** and **manage** (multiple) large datasets

No **computer system** to **manage**

Pay **only** for what you **use**

**Help desk** and **documentation**

# Platforms and Services

## Explore Data
- PIC-SURE
- Gen3

## Analyze Data
- Seven Bridges
- Terra

→ **View BDC Services**

**What Do You Want to Do Today?**

**Explore Available Data**

### BDC-Gen3

Gen3 is a software platform that allows partner organizations and grant approved researchers to search and access harmonized datasets. Users can search over project and study-specific genomic and phenotypic data and export selected cohorts to analytical workspaces in a scalable, reproducible, and secure manner.

Launch | Documentation | Learn

### BDC-PIC-SURE

Explore available data through *BDC-PIC-SURE* with interactive search and visualizations for feasibility assessment. Use query results to create a cohort, with the ability to choose specific variables of interest to export into an analysis environment.

Launch | Documentation | Learn

**Analyze Data in Cloud-based Shared Workspaces**

### BDC-Seven Bridges

Utilize collaborative workspaces for analyzing genomics data at scale. Access hosted datasets along with

### BDC-Terra

Share and compute across large genomic and genomic-related datasets. Terra offers a stand-alone computational

# Community engagement and support

*Though the primary goal of BDC is to build a data science platform, at its core, this is a people-centric endeavor. BDC is also building a **community of practice** working to collaboratively solve technical and scientific challenges.*



- User-driven, vibrant community
- Peer-to-peer mentoring
- Support available via platforms
- Community Forum
- Community Hours & Showcases

Join the community: https://biodatacatalyst.nhlbi.nih.gov/contact/ecosystem

# Questions?

# The Future of Data Science

- How do you see the field of data science **evolving** in the **next five years**?

- What are **topics** we in the profession should be paying attention to right now? Why?

- What can be done to help better **prepare** future data scientists?

**→ Answer now on Jamboard**

# Data Science Community Engagement

# Fellowship programs

## 52 BDC Fellows over three cohorts

- Hands-on, modularized training:
  - Introduction to the platform
  - Creating Your Own Tool/Workflow
  - Performing GWAS
  - Estimating Cloud Costs

- Scientific projects were executed which otherwise might not have been possible

- Retrospective interviews for user-defined development

## Four Bench2Bassinet Fellows

- Hands-on, modularized training:
  - Introduction to the platform
  - Creating Your Own Tool/Workflow
  - Interoperability
  - Estimating Cloud Costs

**Read about the BDC Fellows →**

NIH National Heart, Lung, and Blood Institute | BioData CATALYST®

# Workshops & Courses

- **UW Summer Institute for Statistical Genetics**
  - >300 students over four years
- **American Thoracic Society**
  - Asynchronous "flipped" approach for 50 attendees
- **CHARGE Consortium Annual Meeting, 2022**
  - Hands-on workshop for 75+ attendees, with interactive demo
- **Howard University Data Science Collaboration**
- **AIM-AHEAD PRIME**
- **PRIDE Programs**
  - Held six sessions with 50+ students in 2022



**image credit:** https://www.biostat.washington.edu/suminst/sisg

# Questions?

# Break

While you wait…

If you **have an eRA Commons account** and would like to follow along with the upcoming demonstration, visit https://picsure.biodatacatalyst.nhlbi.nih.gov/ and log into the platform.

BREAK

# Using BDC in a Research Project

NIH — National Heart, Lung, and Blood Institute

**BioData CATALYST** ®

# Curricula of:

# **Using BDC in a Research Project**

**Learning goals:**

- Search and select data relevant to your research question
- Create and use a project in a cloud-computing analysis workspace
- Discover some available tools and workflows

# Platforms and Services

## Explore Data
- PIC-SURE
- Gen3

## Analyze Data
- Seven Bridges
- Terra

→ **View BDC Services**

---

### What Do You Want to Do Today?

**Explore Available Data**

| BDC-Gen3 | BDC-PIC-SURE |
|---|---|
| Gen3 is a software platform that allows partner organizations and grant approved researchers to search and access harmonized datasets. Users can search over project and study-specific genomic and phenotypic data and export selected cohorts to analytical workspaces in a scalable, reproducible, and secure manner. | Explore available data through *BDC-PIC-SURE* with interactive search and visualizations for feasibility assessment. Use query results to create a cohort, with the ability to choose specific variables of interest to export into an analysis environment. |
| Launch   \|   Documentation   \|   Learn | Launch   \|   Documentation   \|   Learn |

**Analyze Data in Cloud-based Shared Workspaces**

| BDC-Seven Bridges | BDC-Terra |
|---|---|
| Utilize collaborative workspaces for analyzing genomics data at scale. Access hosted datasets along with | Share and compute across large genomic and genomic-related datasets. Terra offers a stand-alone computational |

# Finding Data on BDC

# Introduction to dbGaP

BDC ingests various datasets from the **Database of Genotypes and Phenotypes, or dbGaP** (https://www.ncbi.nlm.nih.gov/gap/)

What is dbGaP?

- Public repository for individual phenotype, exposure, genotype, and sequence data
- Main purpose is to archive and distribute the results of studies investigating the association between genotype and phenotype
- Researchers submit a Data Access Request (DAR) and are able to download the study files when authorized for research

# How is data organized in dbGaP?

- Data is organized into **studies**
  - Each study has a specific **accession number** or unique identifier (e.g., phs000007)

- Studies have multiple **subjects**, or study participants

- Data organized by **consent groups**, based on consents given by subjects (research purposes their data can be used for)

- Studies consist of **phenotypic** and/or **genotypic** data
  - Phenotypic data is generally referred to as **variables**
  - Genotypic data is generally referred to as **samples**

# Data Available in BDC

**3.42**
**Petabytes of data**

**280,000+**
**Participants**

**490,000+**
**Data files**

**150,000+**
**Whole genomes**

# Data Available in BDC

**BDC is always ingesting new data**

Check BDC website for a full list of studies available on the ecosystem

**Resources → Data**

Click "Explore Studies"

EXPLORE STUDIES



## BioData Catalyst Studies

The filterable data table below provides metadata on all non-COVID studies available in BioData Catalyst. Note that some parts of the ecosystem may lag in showing some datasets. View COVID-19 studies here.

Filter by Study Name

| | | Accession | Study Name |
|---|---|---|---|
| ☐ | › | phs001607.v2.p2 | NHLBI TOPMed: Pulmonary Fibrosis Whole Genome Sequencing |
| ☐ | › | phs001607.v2.p2 | NHLBI TOPMed: Pulmonary Fibrosis Whole Genome Sequencing |
| ☐ | › | phs001601.v1.p1 | NHLBI TOPMed - NHGRI CCDG: Penn Medicine BioBank Early Onset Atrial Fibrillation Study |
| ☐ | › | phs001927.v1.p1 | NHLBI TOPMed: SubPopulations and InteRmediate Outcome Measures In COPD Study (SPIROMICS) |
| ☐ | › | phs000179.v6.p2 | Genetic Epidemiology of COPD (COPDGene) Funded by the National Heart, Lung, and Blood Institute |
| ☐ | › | phs000954.v3.p2 | NHLBI TOPMed: The Cleveland Family Study (CFS) |

https://biodatacatalyst.nhlbi.nih.gov/resources/data/studies

# Bring-Your-Own Data

- To support **flexibility and analysis**, we allow researchers to bring their own data and workflows into the ecosystem.

- Users can upload data for which they have the appropriate approval, provided that they do not violate the terms of their Data Use Agreements, Limitations, or IRB policies and guidelines.

**Web resource**: Bring Your Own Data

# Gen3 - Key Features

1. Source of Truth - File Object Persistence & Dataset Metadata
2. Interoperability - Standards-based integration points with other systems
3. Data Access - eRA Commons / dbGaP Authorization Inheritance
4. Data Ingestion - Robust data ingestion pipeline

# Gen3 - Discovery Page

A tool for discovery of released datasets (fully open, no required approval to discovery available data).

# Gen3 - Exploration Page

A dynamic summary statistics display and cohort builder for export:

- Search facets leveraging harmonized variables.

Standardized Cohort Handoff support to move cohort to analysis workspaces (e.g. Broad's Terra System, Velsera's Seven Bridges system).

# Empowering researchers to access data

*BDC Powered by PIC-SURE* facilitates approachable research for all skill levels.

Search at the variable value and genomic variant level
+
Apply filters to create a cohort
=
Dataframe ready for research without opening any files or mapping to data dictionaries

# Submitting a Data Access Request (DAR)

BDC uses dbGaP infrastructure for managing access to controlled-access data

Requirements:
1. An NIH eRA Commons account (or other valid NIH login). To learn more about this, visit <u>Understanding eRA Commons</u>.
2. User must have Principal Investigator status. Those who are not PIs can ask their PI to add them as a data downloader.

# Submitting a Data Access Request (DAR)

Components of a DAR:

- Research Use Statement (2200 characters)
- Non-technical Summary (1100 characters)
- BioData Catalyst-specific Cloud Use Statement [Template language available]

For more information, step-by-step instructions, and template language, visit the "Submitting a dbGaP Data Access Request" page of the BDC documentation.

# Live Demo: PIC-SURE Authorized Access

# Breakout room

In your breakout rooms:

1.  **Try:** exploring data in PIC-SURE. If you are not authorized to access data, explore available data in Open Access.
2.  **Discuss:** How could PIC-SURE be useful in your research process?

Feel free to ask TAs any questions.

# Breakout rooms - converge

Would someone like to discuss their PIC-SURE use case?
How could PIC-SURE be helpful in your research process?

# Analysis on BDC

# Curricula of:

# Using BDC in a Research Project

**Learning goals:**

✓ Search and select data relevant to your research question

● Create and use a project in a cloud-computing environment

● Discover some available tools and workflows

# Analysis Platforms

# Workflows

*Workflows (aka pipelines) are a series of steps performed by an external compute engine that are often used for automated, bulk analysis (such as aligning genomic reads)*

# BDC-Terra

- Can write your own in WDL
- Can access 1,500+ public workflows in our methods repository

# *BDC-Seven Bridges*

A curated collection of **800⁺** bioinformatics tools & workflows:

- Optimized for speed & cost in the cloud

- Fully <u>parameterized</u> & customizable

- Accessible via the user interface & API

- Tool descriptions and helpful hints



Open to the public @ platform.sb.biodatacatalyst.nhlbi.nih.gov/public/apps

# Workspaces and Projects

*Workspaces (BioData Catalyst powered by Terra) and Projects (BioData Catalyst powered by Seven Bridges) are dedicated space where you and your collaborators can access and organize the same data and tools and run analyses together.*

# *BDC-Terra*

- Dashboard
  - General overview of the workspace that includes documentation on the workspace itself, cloud information, owners, and tags
  - Good documentation makes your analysis easy to share (with others, as well as with your future self) and reproduce.

# *BDC-Terra*

- Data
  - Import your own data or access data that is stored in Terra
  - Convenient spreadsheet formatted data tables help keep track of all project data, no matter where files are stored in the cloud.

# *BDC-Terra*

- Analyses
  - Interrogate and visualize your data in real time using Galaxy, Jupyter Notebooks, or RStudio
  - All three apps run on virtual machines or clusters of machines in a workspace Cloud Environment.

# *BDC-Terra*

- Workflows
  - Collect, configure (set up) and run workflows for bulk analyses

# BDC-Terra

- Workflows
  - Can edit WDL script directly or use UI to assign inputs and outputs

# *BDC-Terra*

- Useful Workspaces
  - [Working with GnomAD data](#)
  - [Whole Genome Analysis Pipeline](#)
  - [Workflows Tutorial](#)

- Useful Workflows
  - [Processing for Variant Discovery](#)
  - [CRAM to BAM](#)
  - [Generate Sample Map](#)
  - [HaploType Caller](#)

# Seven Bridges - Projects organize files, methods, and results



Project owner

Project owners can add collaborators to the project and define permissions

Project

Files    Apps    Tasks

File 1    App 1    Task 1

File 2    App 2

...    ...

File *n*

New file 1

...

A user selects the files, apps, and parameters to create tasks

Completed tasks generate new files that are stored in the project

Also known as *workspaces* or *sandboxes*

Easily manage collaborators and permissions

# User friendly workflow editor enables reproducibility by default
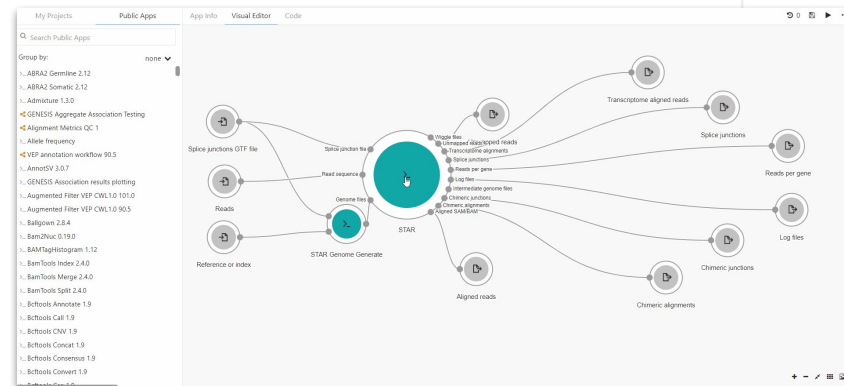
Common Workflow Language enables **portability**, **reproducibility**, and **scalability**

Use or combine 800+ optimized tools and workflows to construct your analysis

Seamlessly import workflows from external public repos

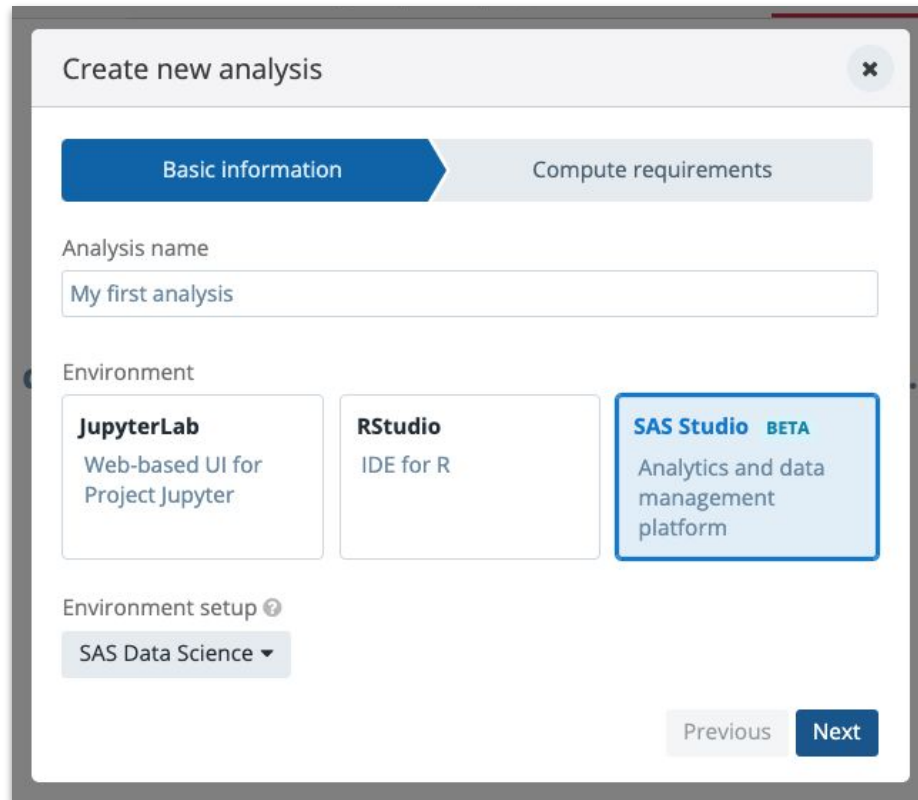Create your own tools with our CWL Tool Editor

Expose or lock parameters appropriately

# Interactive analysis

**Fast prototyping** and implementation of custom tertiary analysis tools using interactive Java, Python and R in the JupyterLab environment as well as RStudio.

All project files available within JupyterLab, RStudio, and SAS. Over 50 instances to select from.



Create new analysis

Basic information | Compute requirements

Analysis name

My first analysis

Environment

**JupyterLab**
Web-based UI for Project Jupyter

**RStudio**
IDE for R

**SAS Studio** BETA
Analytics and data management platform

Environment setup ⓘ

SAS Data Science ▾

Previous | **Next**

# Scale to 100's and 1000's of tasks in parallel using batching

Only one input per task can be selected for batching.

- Turn on the batching option on the draft task page, and select batch criteria: by File, or File metadata (e.g. Sample ID, Library ID).
- For each batch criteria match, a task will be created.

# Detailed documentation and tutorials

## Comprehensive tips for reliable and efficient analysis set-up

BIODATA CATALYST POWERED BY SEVEN BRIDGES

### OBJECTIVE

We have prepared this guide to help you with your first set of projects on BioData Catalyst powered by Seven Bridges. Each section has specific examples and instructions to demonstrate how to accomplish each step. We also highlight potential stumbling blocks so you can avoid them as you get set up. If you need more information on a particular subject, our Knowledge Center has additional information on all of the platform features. Additionally, our support team is available 24/7 to help!

### HELPFUL TERMS TO KNOW

**Tool** refers to a stand-alone bioinformatics tool or its Common Workflow Language (CWL) wrapper that is created or already available on the platform.

**Workflow / Pipeline** (interchangeably used) – denotes a number of tools connected together in order to perform multiple analysis steps in one run.

**App** stands for a CWL wrapper of a tool or a workflow that is created or already available on the platform.

**Task** – represents an execution of a particular tool or workflow on the platform. Depending on what is being executed (tool or workflow), a single task can consist of only one tool execution (tool case) or multiple executions (one or more per each tool in the workflow).

**Job** – this refers to the "execution" part from the "Task" definition (see above). It represents a single run of a single tool found within

## Troubleshooting Failed Tasks

BIODATA CATALYST POWERED BY SEVEN BRIDGES

Tasks and examples described in this guide are available as a public project on the Platform.

Often the first step to a user becoming comfortable using BioData Catalyst powered by Seven Bridges is their gaining confidence in resolving issues they encounter on their own. This confidence usually comes with experience – the experience with bioinformatics tools and Linux environment in general, but also the experience with the platform features.

However, one of the reasons for developing the platform in the first place is to enable an additional level of abstraction between the users and low-level command line work in the terminal. Even though there are a number of platform features that help with tracking down the issues, the less-experienced users can still face challenges with troubleshooting because the whole process might assume familiarity digging through the tool and system messages.

Fortunately, there is a set of steps that most often brings us to the solution. Based on internal knowledge and experience, the Seven Bridges team has come up with the *Troubleshooting Cheat Sheet* (Figure 1) which should help you navigate through the process of resolving the failed tasks.

SevenBridges

**Troubleshooting CHEAT SHEET**

Error Message

A Informative

JS expression evaluation error

Other

Running Running

## Visit the Knowledge Center

NIH National Heart, Lung, and Blood Institute

BioData CATALYST®

# Getting Help - Contacting Support from the platform

24/7 Help Desk can help you with failed analyses, login issues, or any other platform issue.

# Whole person: integrating social, environmental, and genetic factors

BioData CATALYST

NIH National Heart, Lung, and Blood Institute

# GWAS is a great place to start

**Genome wide association study**

- Method that helps scientists identify genes associated with a particular disease or trait

- Used frequently on BDC

- Identifies signals of significance but further experimental follow up required to understand the functional biology

# Manhattan plot is a main GWAS output



P-value gets smaller as "height" increases

—

SNP Significance gets higher as "height" increases

# Interactions between Genetics and SDoH

- **Social determinants of health**
  - Income and social protection
  - Education
  - Unemployment and job insecurity
  - Working life conditions
  - Food insecurity
  - Housing, basic amenities and the environment
  - Early childhood development
  - Social inclusion and non-discrimination
  - Structural conflict
  - Access to affordable health services of decent quality

[Link to pubmed entry](#)

SDoH can be used as covariates in GWAS?

# Run association pipelines out of the box

- **GENESIS**
- Plink
- EPACTS
- STAAR

No login required!

Publish your apps to share with the world!



**Public apps for your data analysis**

We offer publicly available Common Workflow Language workflows and tools to enable reproducible bioinformatics.

Browse 848 apps

BUILD A PIPELINE

**RNA-seq alignment - STAR 2.5.4b**

Toolkit version: STAR 2.5.4b

This workflow performs the first step of RNA-seq analysis - alignment to a reference genome and transcriptome. STAR (Spliced Transcripts Alignment to a Reference), an ultrafast RNA-seq aligner, is used in this workflow. STAR is capable of mapping full length RNA sequences and detecting de novo cano...

Alignment   RNA-Seq

Copy   ▶ Run

**Whole Exome Sequencing - BWA +...**
Toolkit version: GATK 4.1.0.0

**Fusion Transcript Detection -...**
Toolkit version: Fusion Transcript Detection - ChimeraScan 1.0

**Whole Genome Sequencing - BWA +...**
Toolkit version: GATK 4.1.0.0

https://platform.sb.biodatacatalyst.nhlbi.nih.gov/public/apps

# GWAS Additional Features

- Annotation Explorer
  - Use for prep for Aggregate test
  - Post GWAS annotations

- Aggregate and Sliding window association test apps

- Data Overview feature
  - Explore variant frequency across TOPMed Freeze8 studies on GRCh38

- Study Variable Explorer
  - Search for annotated variables across studies and manually annotate and compare

**GENESIS Model Explorer App**

The GENESIS Model Explorer App is a Shiny application developed by the Genetic Analysis Center at the University of Washington in collaboration with Seven Bridges Genomics.

Learn more

Open

**LocusZoom Shiny App**

LocusZoom Shiny App allows users to visualize and interactively explore the results of a single variant association test.

Open

**OmicCircos App**

OmicCircos App is R Shiny application created around OmicCircos R package for more effective generation of high-quality circular plots for visualizing variations in

Open

# Researcher Spotlight: Dr. Jamie Murkey

Primary question: *What role do social factors play in influencing the pathophysiology of cardiovascular disease (CVD) and racial/ethnic disparities in CVD events?*

Aim 1: Perform a comparative assessment between two bioinformatics tools to develop and test a cloud-based telomere length estimation workflow.
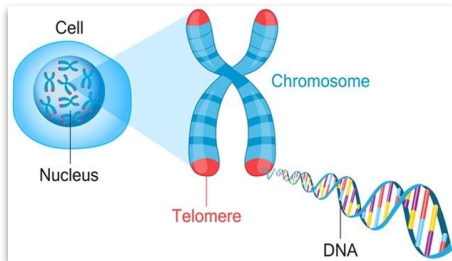
He developed **TeleGen**, a bioinformatics workflow that generated telomere length data for 100% of TOPMed MESA Participants



## Researcher Spotlight

Investigating the impact of social and environmental drivers of cardiometabolic health outcome disparities

Dr. Jamie Murkey is an Intramural Social and Environmental Epidemiology Postdoctoral Fellow at the National Institute of Environmental Health Sciences (NIEHS).

### Why BDC?

*"The BDC ecosystem enabled me to develop a cloud-based telomere estimation workflow, which was essential for completing my dissertation research. Members of the BDC community were supportive and available throughout that process. I believe that the BDC ecosystem serves as an important resource for genomic and non-genomic researchers alike, which can remove historical barriers for advancing science."*

Linkedin Profile

# Live Demo: Seven Bridges

# Breakout room

In your breakout rooms:

1.  **Try:** Create your own project/sandbox on BDC-Seven Bridges
    a.  Name project in this format - **NDSA_First_Last**
    b.  Add your TA to the project
    c.  Copy files, apps and JupyterLab from the main NDSA Workshop project

2.  **Discuss:** What apps or software environments would you like to run on Seven Bridges for your **research** or to help with **teaching**

**Feel free to ask TAs any questions.**

# Breakout rooms - converge

- What was the experience like getting your sandbox project setup?

- What software would you need for your particular use cases?