NASA GL4U: On-Demand Bioinformatics Training Using Space Biology Omics Data

A RASA

2024 ASGSR Educator Session 3

Amanda M. Saravia-Butler, Ph.D. Science and GL4U Lead, NASA GeneLab Contractor: Amentum

Biological & Physical Sciences

National Aeronautics and Space Administration



GL4U Background / Objectives

Background:

- NASA's GeneLab project empowers researchers with open access to space-relevant multi-omics data through the <u>Open Science Data</u> <u>Repository (OSDR)</u>
- GeneLab for Colleges and Universities (GL4U) offers bioinformatics training for space biology
- GL4U was conceptualized in 2020 through a collaboration with SJSU and USRA, and the first GL4U bootcamp was held virtually in June 2021



Objectives:

- Educate and train the next generation of scientists to process, analyze, and interpret space-relevant 'omics data using
 publicly available data and bioinformatics tools
- Maximize the number of scientists who understand Space Biology data and utilize NASA's open-source 'omics data and tools
- **Provide educators with** the **knowledge and resources** required to train and inspire their students using GeneLab bioinformatic analyses as an entree into Space Biology

GL4U Content Design

GL4U materials are organized into introduction and omics-specific modules

Introduction Module

- Overview Lecture (NASA, SMD, BPS, OSDR, GeneLab)
- Jupyter Lab Tutorial
- Unix Jupyter Notebook (hands-on training)
- R Jupyter Notebook (hands-on training)

The Intro Module serves as a pre-requisite for any omicsspecific module

Unix Intro JN

5. Running commands

Using the foundational rules described above, we will begin running some commands.

date

date is a command that prints out the date and time. This particular command doesn't require any argument

[2]: date Sat Jul 8 22:40:55 PDT 2023

When we run date with no arguments, it uses some default settings, like assuming we want to know the time in our computer's currently set time zone. But we can provide optional arguments to date.

Optional arguments most often require putting a dash in front of them in order for the program to interpret them properly.

Here, we are adding the -u argument to tell the date program to report UTC time instead of the local time which will be the same if the computer we're using happens to be set to UTC time:

n [3]: date -u

Sun Jul 9 05:41:28 UTC 2023

Note that if we try to run the command above without the dash, we get an error (ignore the message that print out highlighted in red, we wouldn't normally see that outside of a notebook):

n [4]: date u

date: invalid date 'u'

Note

date-u date-u: command not found

Notice that the error above comes from the program date. So the program we wanted to use is actually responding to us, but it doesn't seem to know what to do with the letter \mathbf{u} we gave it. And this is because wasn't prefixed with a dash, like $-\mathbf{u}$.

Let's see what happens if we try to enter this without the "space" separating date and the optional argument -u, the computer won't know how to break apart the command and we get a different error (again, ignoring the red output):

R Intro JN

1d. Data frame manipulations

Much of the data we work with in bioinformatics is in the data frame or matrix format. For example, gene expression data is usually held in matrix format, with samples as columns and genes as rows, where each entry (or cell) in the matrix contains the expression of a particular gene in a particular semple.

When analyzing numerical data in table format, it can be useful to be able to perform mathematical functions on all cells in a data frame, such as adding a value to all cells or taking the log of all cells. Fortunately, R makes that easy for us to do. Below are some examples of common mathematical manipulations we often perform on data frames in

bioinformatics.

Add a value to all cells

In R, you can add, subtract, multiply, or divide the number in every cell of a data frame by a specific value very easily. Run the command in the next cell to add 1 to every value in your myDF data frame.

	column1	column2	column3	
	column	conunna	columno	
	<dpl></dpl>	<dpl></dpl>	<dbl></dbl>	
row1	2	3	4	
row2	5	6	7	
row3	8	9	10	
row4	11	12	13	
row5	14	15	16	
row6	17	18	19	
row7	20	21	22	
row8	23	24	25	
row9	26	27	28	
row10	1	1	1	

Omics-specific Modules (RNAseq, Amplicon Seq, etc.)

- Omics Data Lectures
 - Experimental design
 - Sample preparation and quality control
 - Data processing tools and visualizations
 - Results analysis and interpretation
- Hands-on data processing and analysis of an OSDR dataset via Jupyter Notebooks (JNs)

AmpSeq JN

It's worth noting again that these are not interpretable as "real" numbers of anything (due to the nature of sequencing data), but they can still be useful as relative metrics of comparison within a study.

As a reminder:

 the left, "Chaot", is an estimate of total richness (total number of unique "things")
 the right, Shanon, is a metric of diversity – which incorporates "richness" and "evenness" (the relativ proportions of all our unique things to each other)

Looking at the plots above

1. Do you notice any immediate differences between the flight and ground-control groups?

Some thoughts

We can also modify the parameters of the plot_richness() phyloseq function to group samples based our treatment groups, which can be helpful sometimes:

In [32]: plot_richness(ASV_pbyseq, x = "treatment", color = "treatment", measures = c("Chaol", "Shannom calculated color paneously (values = unique) (apb(calculate)) (abs(color)) + theme_logend, title = element_blank(), text = element_text(size = 18), axis.text.x = element_text(angle = 96, upst = 0.5, h) (ust = 1))



RNAseq JN

4c. Volcano Plot

Finally, let's make a volcano plot to identify a few interesting genes. A volcano plot is a scatterplot which shows
the relationship of the adjusted p-value to the log2 fold change. Genes with large fold changes that are also
statistically significant by adjusted p-value are labeled.
First we'll use the default settings from the EnhancedVol cano() function: lon2 Fold Change cutoff > 12L and

First, we'll use the default settings from the EnhancedVolcano() function: log2 Fold Change cutoff > [2], and the adjusted p-value cutoff is < 10e-6.

Note: You can read more about the EnhancedVolcano() function and see some examples clicking here

- - pointSize = 3.0 labSize = 6.0,
 - colAlpha=0.5)

ggsave(file.path(DGE_plots,'GLDS-104_volcano_DGE.png'), width = 6.5, height = 8.5, dpi = 300)







GL4U: Previous Bootcamps





2021 – Space biology & RNAseq

- Virtual 1-week long bootcamp with 14 **SJSU** students
- Collaboration with USRA/SJSU; SJSU compute system
- Space biology- and RNAseq-specific lectures and hands-on instruction using Jupyter Notebooks (JNs)

2023 - Space biology & Amplicon Seq

- In-person 4-day bootcamp with 24 CSULA students
- Collaboration with JPL/CSULA; NSF ACCESS compute system
- Space biology- and Amplicon Seq-specific lectures and JNs







2022 – Space biology & RNAseq

- Virtual ~1.5-week long bootcamp
- Collaboration with JPL; SMCE compute system
- 6 professors and 4 graduate students from 4 Institutions
- Space biology- and RNAseq-specific lectures and hands-on instruction using JNs
- Resources to teach at home institution

On-demand courses now available!

GL4U On-Demand

5

GL4U: Intro On-Demand Content

Recorded lectures on YouTube

Hands-on activities via GitPod



GL4U:RNAseq On-Demand Content

Recorded lectures on YouTube

Hands-on activities via GitPod

🔲 🕒 YouTube		Search			Q .		⊕ ¢ ♠		Second Table Control C
Ser reds Start reds Set reds S	RNA Sc ered App to Gename, trans ysis	Re-depictor or File	Inte RMAs	Cinerate CMA, fragr size saleet, cold in Library QC Sequence and Doub of millions Dos of millions Dos of millions Dos of sillions urteleavirue-sec-basice-applicatione-	ett, sikais hitter, Akuada M	GLU: RNA Sequencing Max Genetal-4/4	× ; 22224 22224 2 22224 2 22224 2 22224 2 22224 2 22224 2 22224 2 22224 2 22224 2 22224 2 22224 2 22224 2 22224 2 22224 2 22224 2 22224 2 22224 2 22224 2 22224 2 2 22224 2 2 2 2 2 2 2 2 2 2 2 2 2	Filter files by name Q / GL4U_RNAseq_JNs / Modified Mme Modified © [0 01-RNAseq_pro 20m ago > [0 03-RNAseq_ana 20m ago	Bash B
ANASeq Lecture for AN ANASeq Lecture for AN ANASeq Lecture for AN Softwarehers NASA GeneLab Softwarehers Softwarehers Normalizatio C. For every gene i calculated. This of genes are not of genes are not	On with in a sample, is performed t differential	Search Search Median (the ratios (sar d for each sam y expressed, th	of Ratio		ntoad & Clip C Save	GL4U: RNA Sequencing NASA Genetab- 3/4 C S S GL4U: RNAseq Lecture 1of RNAseq Overview part3of3 GL4U: RNAseq Lecture 1of RNAseq Coverview part2of3 GL4U: RNAseq Lecture 1of RNAseq Coverview part2of3 NASSA Genetab	E [□] Q A × : 2024 2024		Raw Reads Trim/Filter Raw Data (TrimGalorel) Raw Data QC (FastQC -> multiQC) Trimmed Data QC (FastQC -> multiQC) Trimmed Data QC (FastQC -> multiQC) Trimmed Data QC (multiQC) Trimmed
should have sim	mple1 Sample2	pseudo-reference	ratio of	ratio of		GL4U: RNAseq Lecture 2of Statistics Overview 2024 NASA GeneLab	2		1. RNAseq processing: fastq to counts
FF2A 1	1489 906	1161.5	1489/1161.5 =	906/1161.5 =		4 GL4U: RNAseq Lecture 1of RNAseq Overview part1of3	2 2024		
APCD4	22 12	16.0	1.28	0.78		NASA GeneLab			Here we are going to setup a directory structure to store the output data we'll generate, then we will follow the steps in the Genetab standard RNAseq processing pipeline to generate raw count data. All of these steps are most easily done at a Unix-like command line, so this notebook uses a "Bash" kernel, a common language used in a Unix-like
ABCD1	22 13	16.9	793/570.2 =	410/570.2 =					environment.
MEFV 7	/93 410	570.2	1.39	0.72		Statistic	s		This is notebook 1 of 2 of GL4U's RNAseg Module Set. It is expected that GL4U's Introduction Module Set has been completed already.
BAG1	76 42	56.5	76/56.5 = 1.35 521/883.7 =	42/56.5 = 0.74 1196/883.7 =					
MOV10 5	521 1196	883.7	0.590	1.35		Overvie	W		Next: 2. RNAseq analysis
					anatti 012 yauti menulimba keni				
 ◀ ► ► ◀) 49:14 / 1:05:0	04								Table of Contents
GL4U: RNAseq Lecture 2of2 Sta NASA GeneLab 504 subscribers Subscriber	atistics Overviev	v 2024	凸 Like	🖓 🧀 Share 🛓 Dov	nload 💥 Clip 🗔 Save …				0. Setup 1. Raw Data Quality Control (QC) 1. Raw Data QC with FastQC

GL4U RNAseq Bootcamp Using On-Demand Content

GL4U: RNAseq Virtual Bootcamp, 9/2024

- Virtual 5-day bootcamp, 9/27 9/29; 10/4 10/5
- 34 students and 2 educators from 4 institutions
 - 14 students, 2 professors from AAMU
 - 10 students from CSULA
 - 1 student from UCSC
 - 9 students from UNM
- Compute resources: GitPod free compute resources
- Bootcamp covered GL4U Introduction RNA Sequencing modules

Sophomore

B.S. Completed

Master's Student

Junior

Senior

Educator





8

Participant Feedback

"I learned more about bioinformatics than I have in any other higher education class. This will greatly improve my comprehension when reading genetics papers. I cannot emphasize enough how valuable this workshop was, it needs to continue for other students, and I will absolutely be looking to attend more workshops such as this in the future."

"This was an amazing experience. Amanda was terrific. It really opened up my eyes to space biology and really showed me this is something I would love to be able to pursue." "When I first took the survey, I was intimidated by some of the topics because I was unfamiliar with them, but **after this bootcamp I have a better understanding of RNA sequencing**. Dr. Saravia has **inspired me to look into space biology careers**, a career field I was not considering prior to this bootcamp."



Certificates of Completion

GL4U: Introduction





Certificate of Completion

This certifies that on the 5th day of October, 2024

FirstName LastName

has successfully completed the OSDR GL4U: Introduction Module Set

Samrawit Gebre

Samrawit Gebre Open Science Data Repository, Project Manager NASA Ames Research Center



GL4U: RNA Sequencing



Certificate of Completion

This certifies that on the 5th day of October, 2024

FirstName LastName

has successfully completed the OSDR GL4U: RNA Sequencing Module Set

Samrawit Gebre

Samrawit Gebre Open Science Data Repository, Project Manager NASA Ames Research Center



Collaborate on Data Mining/Publications



Cell Press Package 2020: The Biology of Spaceflight <u>https://www.cell.com/c/the-biology-of-spaceflight</u>



Nature Portfolio Collection 2024: Space Omics and Medical Atlas across orbits (SOMA) https://www.nature.com/collections/ebdbcahdgc



GL4U: Intro Course Available On Canvas

Enroll on Canvas: https://canvas.instructure.com/enroll/63FHH6

GL4U: Intro > Syllabus			GL4U: Intro > N	Iodules ····································		
Home Syllabus Modules	Course Syllabus Part I: Pre-Course Survey		Home Syllabus Modules Assignments	MODULE 1: Pre-Course Survey Complete All Items		
Assignments Files Discussions Grades Part II People Lucid (Whiteboard) Part II	Part II: Introduction to NASA Open Life Sciences 1. Intro to NASA, Science Mission Directorate, and Space Biology		Discussions Grades People	GL4U Introduction On-Demand Pre-Course Survey O pts Submit	0	
	2. Intro to NASA's Open Science Data Repository and GeneLab		Lucid (Whiteboard)	MODULE 2: Introduction to NASA Sciences Prerequisites: MODULE 1: Pre-C Complete All Ite Complete All Ite View	ourse Survey	
	 Setting up GitPod Intro to Jupyter Lab environment Intro to coding in a Unix-based command-line interface Intro to coding in R 			LECTURE 2: Introduction to NASA Open Life Sciences and GeneLab View	0	
	Part IV: Short Read Sequencing Overview			MODULE 3: Coding, The Basics Prerequisites: MODULE 1: Pre-Course Survey, MODUL Introduction to NASA Sciences Complete All Iter	ILE 2:	
	 Library preparation overview Sequencing by synthesis Sequencing data quality assessment 			LECTURE 3: Introduction to the Command Line, R, and Jupyter View HANDS-ON ACTIVITY 1: Getting Started With GitPod opts Submit		
	Part V: Post-Course Survey			 HANDS-ON ACTIVITY 2: Jupyter Introduction o pts Submit HANDS-ON ACTIVITY 3: Unix Introduction 	0	

GL4U: Intro Course Available On Canvas

Enroll on Canvas: https://canvas.instructure.com/enroll/63FHH6



GL4U: RNAseq Course Available On Canvas

Enroll on Canvas: https://canvas.instructure.com/enroll/KYXWN6

GL4U: RNAseq	> Syllabus		GL4U: RNAseq > Modules				
Home Syllabus Modules Assignments Files Discussions Grades People Lucid (Whiteboard)	Course Syllabus Part I: Pre-Course Activities 1. Verify Successful Completion of GL4U: Introduction Course 2. Pre-Course Survey Part II: RNA Sequencing Overview and Data Pre-processing 1. What came before RNAseq 2. RNAseq experimental design 3. Raw data QC 4. Trim/filter data and QC	Home Syllabus Modules Assignments Files Discussions Grades People Lucid (Whiteboard)	Home Syllabus Modules Assignments Files Discussions Grades People Lucid (Whiteboard)	• MODULE 1: Pre-Course Activities Complete All Items Image: Submit Complete All Items			
	Part III: RNAseq Data Processing 1. Alignment and QC 2. Assess strandedness 3. Quantitation and QC			 MODULE 2: RNA Sequencing Overview and Preprocessing Prerequisites: MODULE 1: Pre-Course Ad Complete All Items LECTURE 1: RNA Sequencing Overview Part 1 of 3 View 			
	Part IV: RNAseq Data Analysis 1. Statistics overview 2. Count data normalization 3. Differential gene expression analysis 4. Data visualization			 HANDS-ON ACTIVITY 1: RNAseq Pre-processing <i>o</i> pts Submit MODULE 3: RNAseq Data Processing Prerequisites: MODULE 1: Pre-Course Activities, MODULE 2: RNA Sequencing Overview and Preprocessing 	0 2:		
	Part V (Optional): Apply Your Knowledge 1. Data analysis on any OSD dataset Part VI: Post-Course Survey			Complete All Items Complete All Items LECTURE 2: RNA Sequencing Overview Part 2 of 3 View HANDS-ON ACTIVITY 2: RNAseq Processing 0 pts Submit	0 0		

GL4U: RNAseq - Apply Your Knowledge

Enroll on Canvas: https://canvas.instructure.com/enroll/KYXWN6



1a. Set up Directory Path Variable

https://nasa.edgebioinformatics.org/home



My Projects My uploads Job Queue

AMS

EDGE bioinformatics is an open-source bioinformatics platform with a user-friendly interface that allows scientists to perform a number of bioinformatics analyses using state-of-the-art tools and algorithms. NASA EDGE takes an updated EDGE Bioinformatics framework and has only the NASA GeneLab Illumina amplicon sequencing data processing pipeline (AmpIllumina) integrated.





LOS Alamos NATIONAL LABORATORY Terms of Use, Privacy

- Awarded ROSES F.7: Support for Open-Source Tools, Frameworks, and Libraries (3 years)
 - Enhancing Analysis Capabilities of Biological Data With the NASA EDGE Bioinformatics Platform

15

What's In Store For GL4U?

Pending funding...

- Video-recorded walk-throughs of hands-on activities
- Automatic assignment grading
- Automatic GL4U certification upon course completion
- Set-up stable Jupyter Lab environment
 - GitPod has a 30min time-out, requiring users to re-run all previous commands upon GitPod restart
- Host live (weekly / biweekly) office hours
- Add additional omics datatypes (e.g. Amplicon seq, Metagenomics, scRNAseq, etc.)
- Expand internationally
- Incorporate NASA EDGE Bioinformatics training
- Continue to host live annual bootcamps

Acknowledgements

Open Science for Life in Space Teams



Open Science Analysis Working Group Members



Support

NASA Space Biology Program NASA Science Mission Directorate NASA Human Research Program NASA Biological and Physical Sciences

Questions?

For more info about GL4U:



Sign up for the GL4U mailing list

 Stay up-to-date on future GL4U events / bootcamps: Send and e-mail to <u>GL4U-join@lists.nasa.gov</u> with the Subject: **subscribe**

Enroll in GL4U: Intro On Canvas



Enroll in GL4U: RNAseq On Canvas

