



TISSUE REPOSITORY



SCIENTIFIC COMMUNITY



OPEN SCIENCE

DATA REPOSITORY



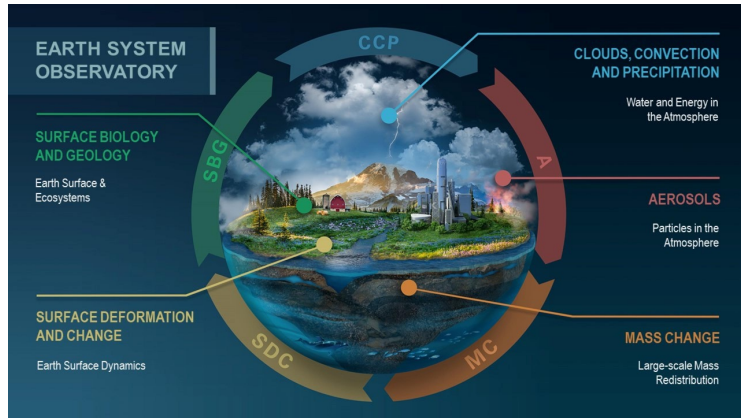
GENELAB
OMICS

+

ALSDA
PHENOTYPIC

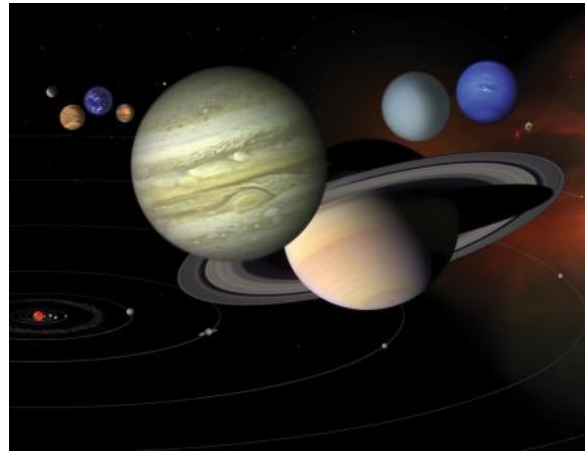


NASA Science Mission Directorate (SMD)



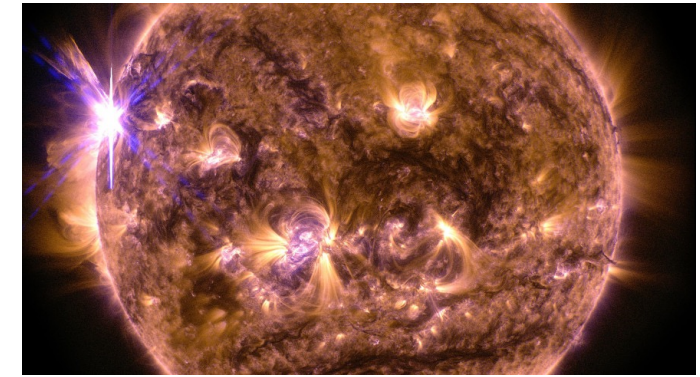
Earth Science

The study of planet Earth



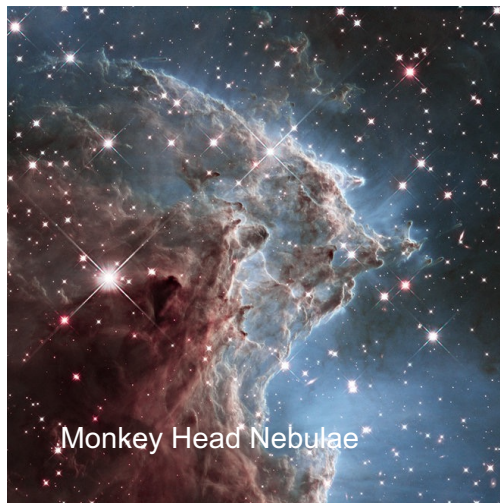
Planetary Science

The study of the origin and history of the solar system and the potential for extraterrestrial life



Heliophysics

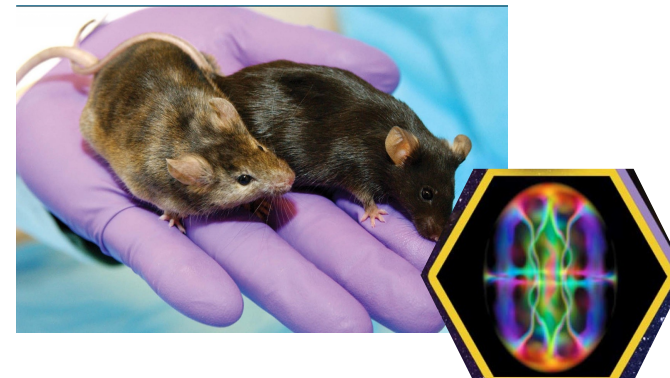
The study of the sun and its effects on the solar system



Monkey Head Nebulae

Astrophysics

The study of the origin, structure, evolution and destiny of the universe



Biological and Physical Sciences

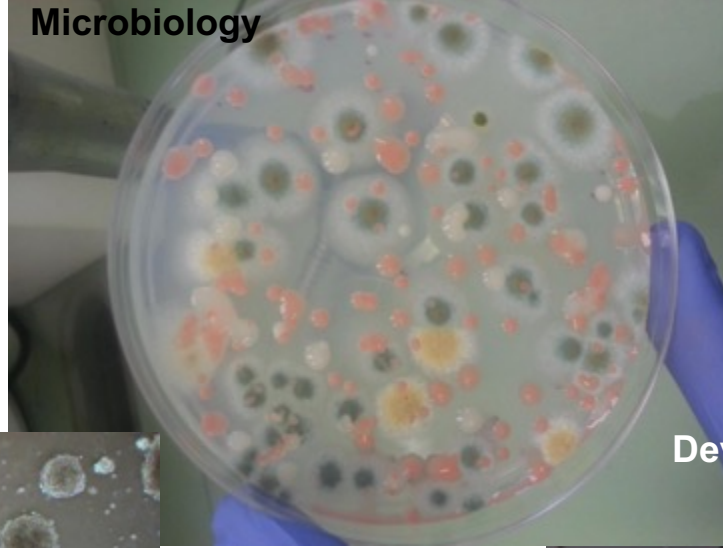
The study of how spaceflight environment effects biological and physical systems

BPS: Space Biology Program

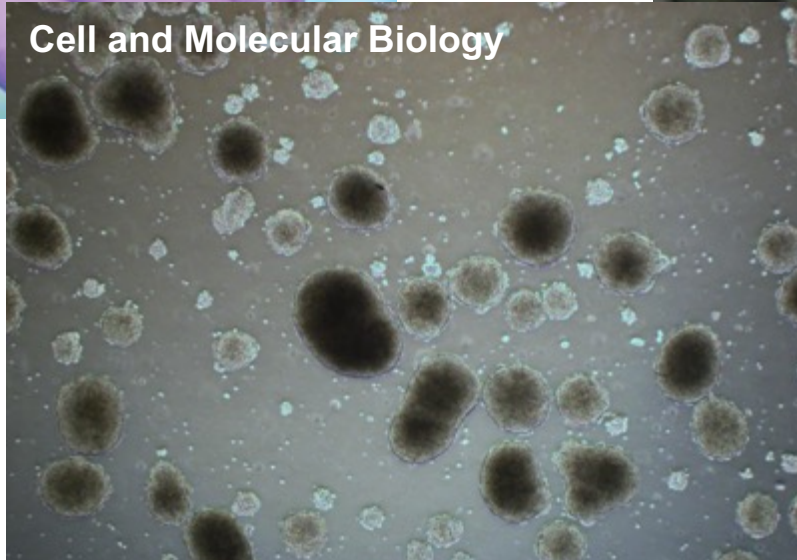
Animal Biology



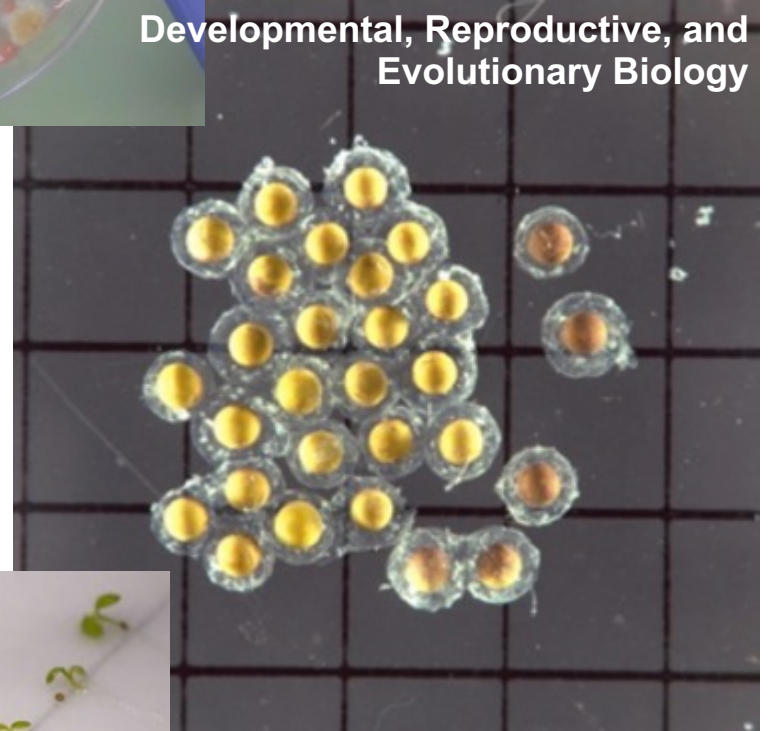
Microbiology



Cell and Molecular Biology



Developmental, Reproductive, and Evolutionary Biology

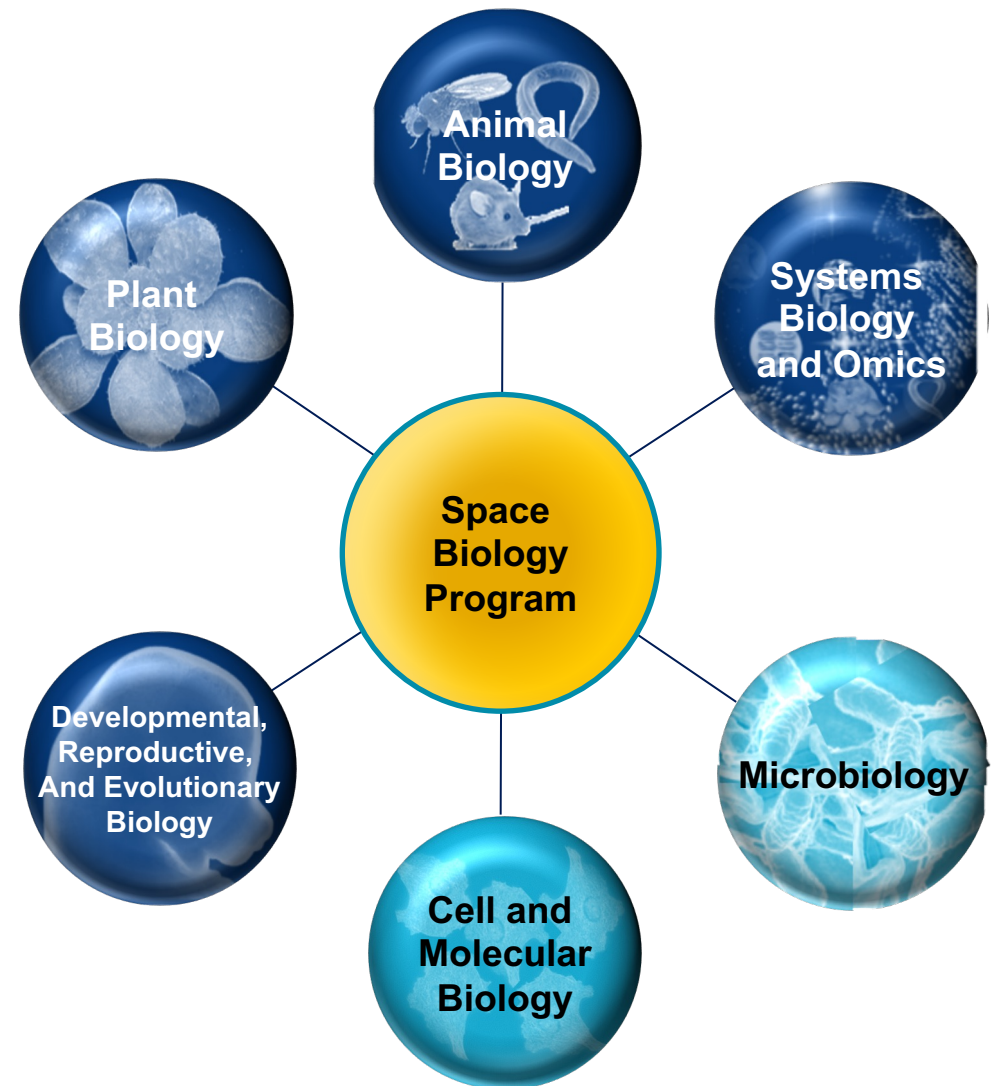


Plant Biology



Space Biology Objectives

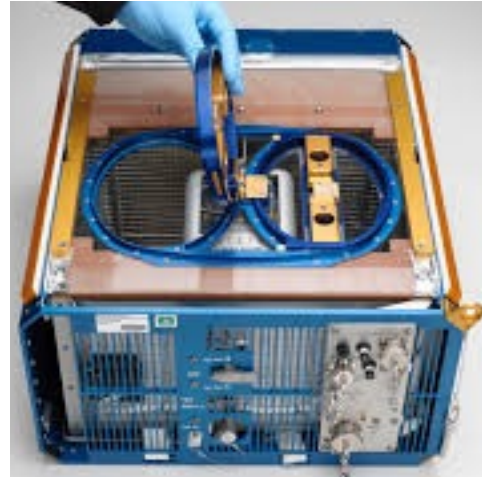
- Discover how biological systems respond to the space environment
- Identify the underlying mechanisms and develop models for biological systems in space
- Develop technologies to enable spaceflight research
- Promote open science through the GeneLab Data System and Ames Life Sciences Data Archive
- Provide mechanistic understanding to support human health in space
- Transfer the knowledge and technology of space-based research to benefit life on Earth



Flight Hardware To Support Space Biology Research



NASA: Bioculture System



NASA: Rodent Habitat



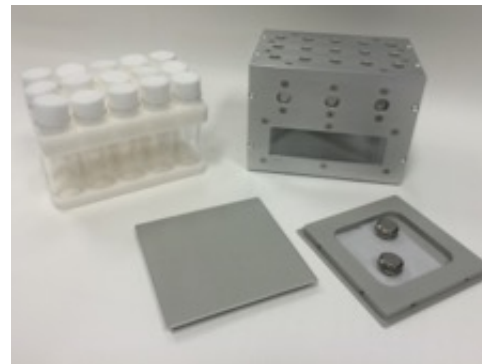
NASA-JAXA: Life Sciences Glovebox



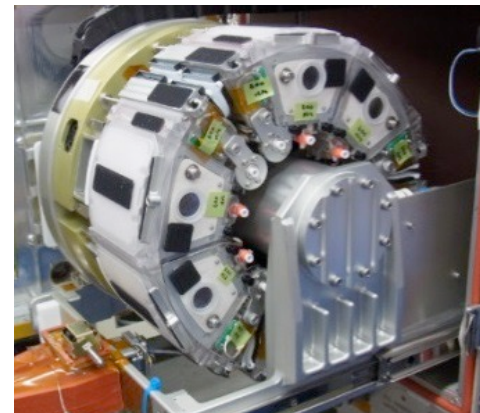
BioServe: BioCell and Fluid Processing Apparatus



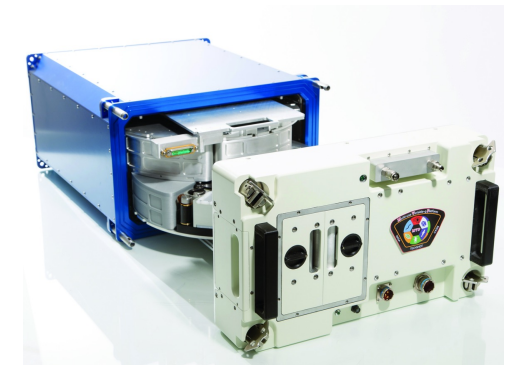
NASA: WetLab-2 (qPCR)



NASA: Vented Fly Box



JAXA: Mouse Habitat Unit

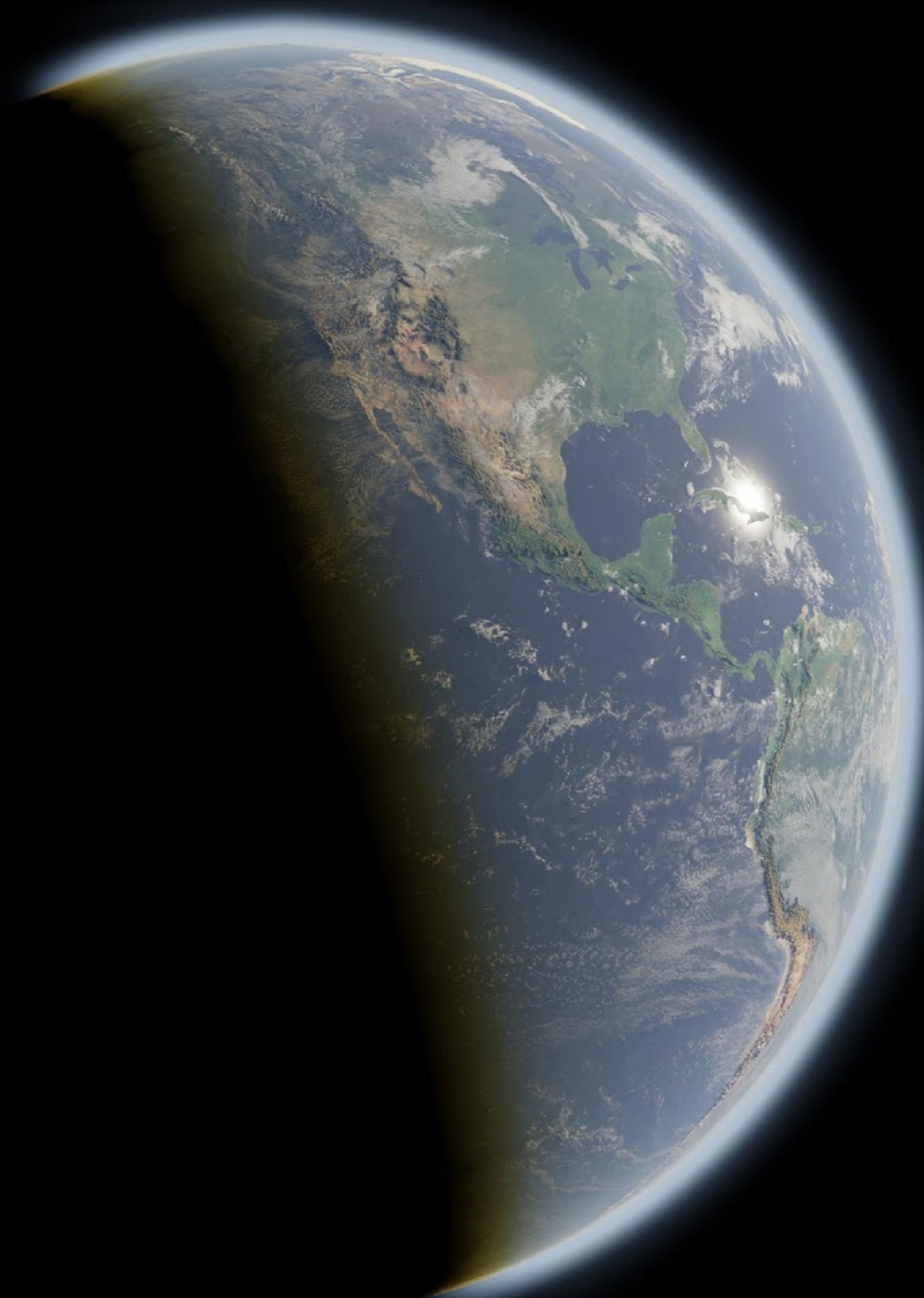


Redwire: Multi-Use Variable-G Platform

Open Science

“We define open science as a collaborative culture enabled by technology that empowers the **open sharing of data, information, and knowledge** within the scientific community and the wider public to accelerate scientific research and understanding.”

Ramachandran, R., Bugbee, K., & Murphy, K. J. **Moving from Open Data to Open Science.** Earth and Space Science, Wiley Publication. <https://doi.org/10.1029/2020EA001562>



OSDR Mission and Vision

Mission

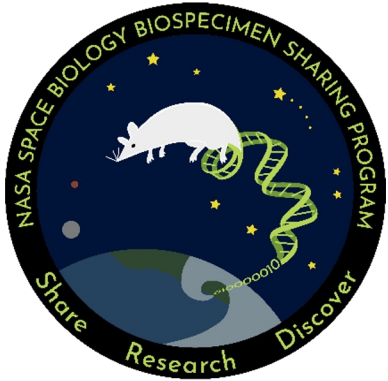
- To enable biological discovery and health resilience for space exploration through open science.

Vision

- Design and deploy a unique repository, the Open Science Data Repository (OSDR), housing standardized metadata and data from spaceflight or spaceflight-relevant samples, following FAIR (Findable, Accessible, Interoperable, and Reusable) principles
- Engage with the scientific community by prioritizing the needs of the Open Science Analysis Working Group (AWG) members, and supporting open source, open science, citizen science, and education initiatives such as GeneLab for High Schools (GL4HS) and GeneLab for Colleges and Universities (GL4U)
- Partner with spaceflight-relevant projects through sample sharing or augmentation of experimental samples to expand analyses
- Process spaceflight-relevant samples answering knowledge gaps, establishing best-practices for sample processing and providing common omics processing platforms for the space biology community, leading to a more cohesive sets of independent datasets
- Curate, beyond best-practice, spaceflight-relevant datasets and make processed data publicly available as expediently as possible
- Provide open-access to computing resources and visualization tools for raw and processed data, democratizing the access to spaceflight-relevant data and disseminating knowledge of how life responds to the space environment. This platform will be an essential tool to discover new countermeasures for human exploration of space, and will inevitably benefit life on Earth as well.

NASA Biological Open Science Resources

Biospecimen Sharing Program (BSP)



Dissection and preservation of rodent tissues from Flight and Ground investigations. Coordination of internal tissue sharing



NASA Biological Institutional Scientific Collection (NBISC)



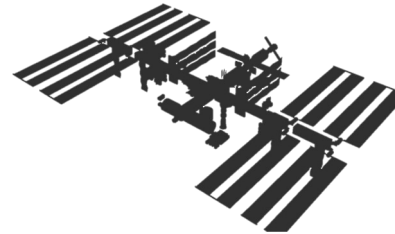
Collection of non-human specimens and space microbial culture



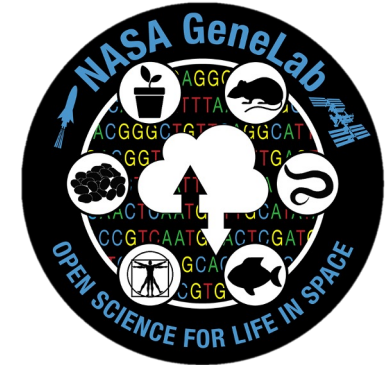
Ames Life Sciences Data Archive (ALSDA)



Collection and curation of mission, project, and imaging data



NASA GeneLab (GL)



Collection and curation of omics data

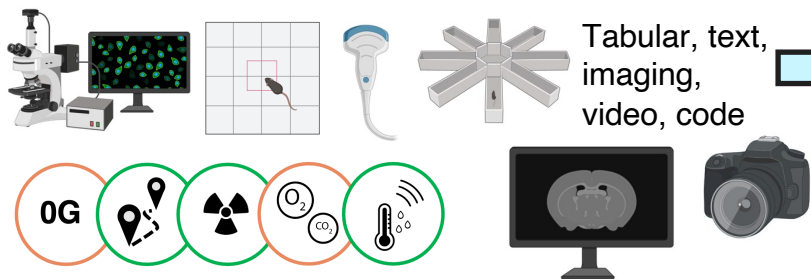


NASA Internal Program

Open-Source Science Programs – Available Globally

Integrating Biological Data Repositories

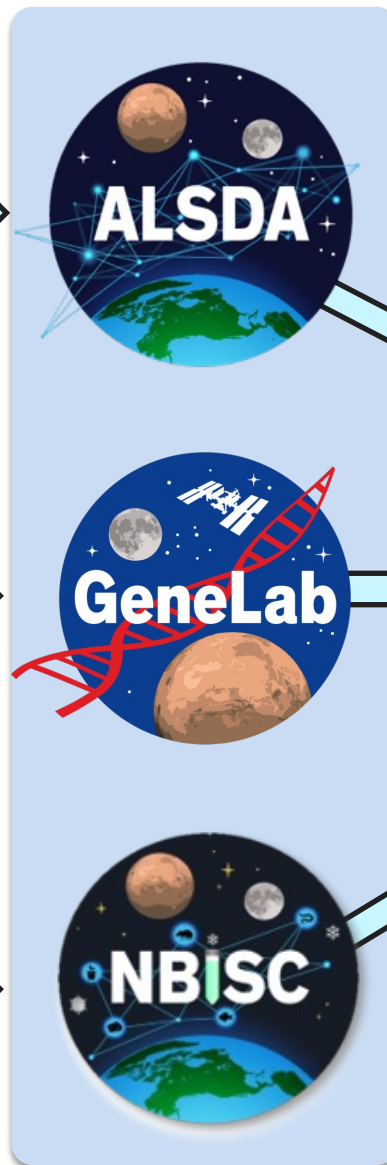
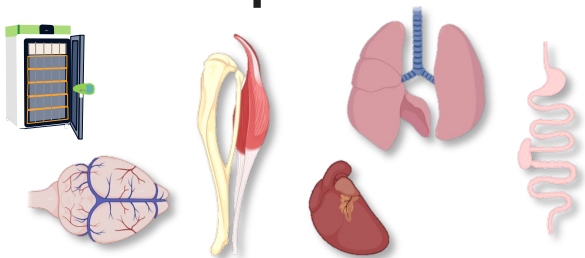
Physiological/Phenotypic/Imaging/ Environmental Telemetry Data



Molecular/Omics Data

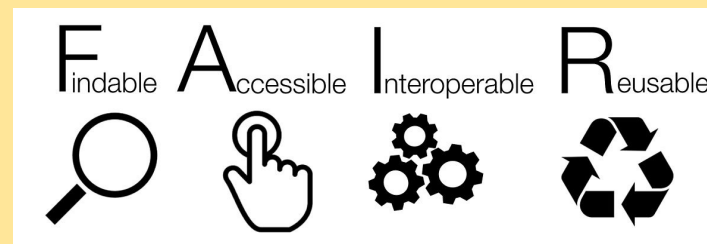


Biospecimens



**NASA Open Science
Data Repository (OSDR)**
osdr.nasa.gov/bio

- Single Submission Portal (BDME)
- User Interface/Website Tool for RDSAs (Research Data Submission Agreements)
- Maximally Open Access with Necessary Controls for Sensitive Data
- Data Maximally FAIR



Open Science Data Repository (osdr.nasa.gov)



- User-submitted data
- GeneLab-generated data
- Ingested data

Open Science Projects

Open Science Projects primary goals aim to increase collaborative scientific data sharing, analysis and more rapid scientific advancement.

GeneLab

GeneLab, an open science multi-omics repository, covering transcriptomics, metagenomics, epigenomics, proteomics, and metabolomics. Studies comprise of data from model organisms including microbes, plants, fruit flies, rodents and humans.

[Learn more GeneLab](#)



BSP

The NASA Space Biology Biospecimen Sharing Program (BSP) collects biospecimens to maximize the scientific return from biological spaceflight and associated ground investigations and to encourage and broaden participation from the scientific community in space biology-related research.

[Learn more about BSP](#)



ALSDA

Ames Life Sciences Data Archive (ALSDA) collects, curates, and makes available space-relevant higher-order phenotypic datasets. Datasets that enable scientists to perform retrospective analysis across missions, experiments, life science disciplines, research subjects, and species.

[Learn more about ALSDA](#)



NBISC

NASA Biological Institutional Scientific Collection (NBISC) is a biorepository of non-human samples collected from NASA-funded spaceflight investigations and correlative ground studies. The purpose of NBISC is to receive, store, document, preserve, and make the collection available to the scientific community.

[Learn more about NBISC](#)



Open Science Data Repository (osdr.nasa.gov)



Open Science Projects

Open Science Projects primary goals aim to increase collaborative scientific data sharing, analysis and more rapid scientific advancement.

GeneLab

GeneLab, an open science multi-omics repository, covering transcriptomics, metagenomics, epigenomics, proteomics, and metabolomics. Studies comprise of data from model organisms including microbes, plants, fruit flies, rodents and humans.

[Learn more GeneLab](#)



BSP

The NASA Space Biology Biospecimen Sharing Program (BSP) collects biospecimens to maximize the scientific return from biological spaceflight and associated ground investigations and to encourage and broaden participation from the scientific community in space biology-related research.

[Learn more about BSP](#)



ALSDA

Ames Life Sciences Data Archive (ALSDA) collects, curates, and makes available space-relevant higher-order phenotypic datasets. Datasets that enable scientists to perform retrospective analysis across missions, experiments, life science disciplines, research subjects, and species.

[Learn more about ALSDA](#)



NBISC

NASA Biological Institutional Scientific Collection (NBISC) is a biorepository of non-human samples collected from NASA-funded spaceflight investigations and correlative ground studies. The purpose of NBISC is to receive, store, document, preserve, and make the collection available to the scientific community.

[Learn more about NBISC](#)



Open Science Data Repository (https://osdr.nasa.gov/bio/repo/)



Open Science Data Repository Search

Search Datasets

Sort By: Release Date

Items per page: 25 1 - 25 of 435



Persistence of Escherichia coli in the microbiomes of red Romaine lettuce (Lactuca sativa cv. 'Outregeous')- does seed sanitization matter?

Study
OSD-385

Organisms	Factors	Assay Types	Release Date	Description
Microbiota	Treatment Seed Sanitization Tissue	Amplicon Sequencing	19-Apr-2024	Seed sanitization via chemical processes removes/reduces microbes from the external surfaces of the seed and thereby could have an impact on the plants,health or productivity. To determine the impact ...

Highlights: *cgene*



Transcriptional profiling of heart tissue from mice flown on the RRRM-2 mission

Study
OSD-580

Organisms	Factors	Assay Types	Release Date	Description
Mus musculus	Spaceflight Age Euthanasia Location	transcription profiling	03-Jan-2024	In the Rodent Research Reference Mission (RRRM-2), forty female C57BL/6NTac mice were flown on the International Space Station. To assess differences in outcomes due to age, twenty 12 week-old and twe...

Highlights: *cgene*



Transcriptional profiling of tibialis anterior muscle from mice flown on the RR-23 mission

Study
OSD-576

Organisms	Factors	Assay Types	Release Date	Description
Mus musculus	Spaceflight	transcription profiling	12-Dec-2023	The objective of the Rodent Research-23 mission (RR-23) was to better understand the effects of spaceflight on the eyes, specifically on the structure and function of the arteries, veins, and lymphati...

Highlights: *cgene*



Ionizing radiation induces transgenerational effects of DNA methylation in zebrafish

Study
OSD-524

Organisms	Factors	Assay Types	Release Date	Description
Danio rerio	Ionizing Radiation Generation	DNA methylation profiling	31-Aug-2023	Ionizing radiation is known to cause DNA damage, yet the mechanisms underlying potential transgenerational effects of exposure have been scarcely studied. Previously, we observed effects in offspring ...

General Search Filters

Data Source

- GeneLab
- ALSDA
- NIH GEO
- EBI PRIDE
- ANL MG-RAST

Data Type

- Study
- Experiment
- Subject
- Biospecimen
- Payload

Show more

Study Search Filters

Project Type

- Ground
- Spaceflight
- High Altitude


Assay Type

- Amplicon Sequencing Assay
- Bisulfite Sequencing
- ChIP-Seq
- Behavior (Gait)
- Gel Electrophoresis

Show more

Organism

Open Science Data Repository (https://osdr.nasa.gov/bio/repo/)



Version 2

OSD-596 Version 2

Artificial gravity partially protects space-induced neurological deficits in *Drosophila melanogaster* (Behavior, Climbing)

123.76 KB Study

Submitted Date: 08-Mar-2023
Initial Release Date: 25-May-2023

ALSDA ID: LSDS-43
DOI: 10.26030/evbf-lv71
Related Studies: OSD-514, OSD-595

Cite this Study

Description

Description

Spaceflight poses risks to the central nervous system (CNS), and understanding neurological responses is important for future missions. We report CNS changes in *Drosophila* aboard the International Space Station in response to spaceflight microgravity (SF μ g) and artificially simulated Earth gravity (SF1g) via inflight centrifugation as a countermeasure. While inflight behavioral analyses of SF μ g exhibit increased activity, postflight analysis displays significant climbing defects, highlighting the sensitivity of behavior to altered gravity. Multi-omics analysis shows alterations in metabolic, oxidative stress and synaptic transmission pathways in both SF μ g and SF1g; however, neurological changes immediately postflight, including neuronal loss, glial cell count alterations, oxidative damage, and apoptosis, are seen only in SF μ g. Additionally, progressive neuronal loss and a glial phenotype in SF1g and SF μ g brains, with pronounced phenotypes in SF μ g, are seen upon acclimation to Earth conditions. Overall, our results indicate that artificial gravity partially protects the CNS from the adverse effects of spaceflight. This study derives results from the climbing assay (locomotor assay).

Factor(s)

Factor	Ontology: Concept
Altered Gravity	Gravity, Altered
Spaceflight	Space Flight
Sex	Sex

Organism(s)


[Drosophila melanogaster](#)

Assay(s)

Measurement	Technology	Device Platform
Behavior	Locomotory Behavior	Manual Counting

Project

Payload Identifier MVP-Fly-01



Version 3

OSD-514 Version 3

Artificial gravity partially protects space-induced neurological deficits in *Drosophila melanogaster*

455.12 GB Study

Submitted Date: 21-Jun-2022
Initial Release Date: 08-Sep-2022

GeneLab ID: GLDS-514
DOI: 10.26030/vsge-cp98

Cite this Study

Description

Description

Spaceflight poses risks to the central nervous system (CNS), and understanding neurological responses is important for future missions. We report CNS changes in *Drosophila* aboard the International Space Station in response to microgravity (SF μ g) and artificially simulated Earth-gravity (SF1g) via inflight centrifugation as a countermeasure. While inflight behavioral analyses of SF μ g exhibit increased activity, postflight analysis displays significant climbing defects, highlighting the sensitivity of behavior to altered gravity. Multi-omics analysis shows alterations in metabolic, oxidative stress, and synaptic transmission pathways in both SF μ g and SF1g; however, neurological changes immediately postflight, including neuronal loss, glial cell count alterations, oxidative damage, and apoptosis, are seen only in SF μ g. Additionally, progressive neuronal loss and a glial phenotype in SF1g and SF μ g brains, with pronounced phenotypes in SF μ g, are seen upon acclimation to Earth conditions. Overall, our results indicate that artificial gravity partially protects the CNS from the adverse effects of spaceflight.

Factor(s)

Factor	Ontology: Concept
Sex	Sex
Spaceflight	Space Flight
Altered Gravity	Gravity, Altered

Organism(s)

[Drosophila melanogaster](#)

Assay(s)

Measurement	Technology	Device Platform
transcription profiling	RNA Sequencing (RNA-Seq)	Illumina
protein expression profiling	mass spectrometry	Orbitrap Fusion

Project

Payload Identifier MVP-FLY-01

GeneLab Data Visualization Portal (<https://visualization.genelab.nasa.gov/data/>)

GeneLab Open Science for Life in Space | Home About Data & Tools Research & Resources Working Groups Help BETA

Number of studies: 294

Factor

Assay technology type

Organism

Tissue

Filter

Assay technology type

- DNA microarray 149
- RNA Sequencing (RNA-Seq) 134
- nucleotide sequencing 75
- mass spectrometry 35
- microarray 9

Organism

- rodent 143
- bacteria 97
- fungus 79
- human 63
- plant 60

Tissue

- root 22
- liver 9
- leaf 8
- whole organism 6
- spleen 5

Factor

- spaceflight 191
- ionizing radiation 107
- time 69

Show 10 entries

GLDS	Title	Assay	Organism	Tissue	Factor
GLDS-1	Expression data from drosophila melanogaster	DNA microarray	Drosophila melanogaster	whole organism	infection, ionizing radiation, spaceflight,
GLDS-4	Microarray Analysis of Space-flown Murine Thymus Tissue	DNA microarray	Mus musculus	thymus	spaceflight,
GLDS-19	Transcription profiling of rat to study the effect of hindlimb unloading on healing of medial collateral ligaments 3 weeks after injury	DNA microarray	Rattus norvegicus	Medial collateral ligament	hindlimb unloading, treatment,
GLDS-21	Effects of spaceflight on murine skeletal muscle gene expression	DNA microarray	Mus musculus	calf muscle, gastrocnemius	spaceflight,
GLDS-25	STS-135 Liver Transcriptomics	DNA microarray	Mus musculus	liver	spaceflight,
GLDS-26	Microbiomes of the Dust Particles Collected from the International Space Station and Spacecraft Assembly Facilities	amplicon sequencing assay	cellular organisms	Cells	sample location,
GLDS-37	Comparison of the spaceflight transcriptome of four commonly	RNA Sequencing (RNA-Seq)	Arabidopsis thaliana	Seedlings	ecotvce, spaceflight,

Visualize Study

2D 3D X: PC1 Y: PC2 Color: mission Shape: Size: Update

Parameters: duration Add

GLDS	mission	library selection
GLDS-47	RR-1	mRNA enrichment
GLDS-48	RR-1	mRNA enrichment
GLDS-137	RR-1	Ribo-depletion
GLDS-168	RR-1, RR3	Ribo-depletion
GLDS-173	STS-135	Ribo-depletion
GLDS-242	RR-9	Ribo-depletion
GLDS-245	RR-6	Ribo-depletion

Hide expanded table

The expanded table will allow you to select the samples you want to use for the Differential Gene Expression Analysis. All samples are selected by default.

Re-normalize using only selected samples

GLDS	Sample	mission	library selection
GLDS-47	<input type="checkbox"/> Mmus_C57-6T_LVR_BSL_Rep1_B1	RR-1	mRNA enrichment
GLDS-47	<input type="checkbox"/> Mmus_C57-6T_LVR_BSL_Rep2_B2	RR-1	mRNA enrichment
GLDS-47	<input type="checkbox"/> Mmus_C57-6T_LVR_BSL_Rep3_B3	RR-1	mRNA enrichment
GLDS-47	<input checked="" type="checkbox"/> Mmus_C57-6T_LVR_FLT_Rep1_F1	RR-1	mRNA enrichment
GLDS-47	<input checked="" type="checkbox"/> Mmus_C57-6T_LVR_FLT_Rep2_F2	RR-1	mRNA enrichment

Visualize Studies

GeneLab Visualization

Home PCA Volcano Pair plot Group 1 Pair plot Group 2 Heatmap DGE GSEA

GLDS-47&GLDS-48&GLDS-137&GLDS-168&GLDS-173&GLDS-242&GLDS-245 DGE

Maximum p-value: Maximum adjusted p-value: Search:

Copy CSV Excel PDF Print

ENSEMBL	Symbol	LOG2FC	PVAL	ADJP
ENSMUSG00000031543	Ank1	-1.8700695226	5e-10	0.0000059515
ENSMUSG00000044309	Apo17c	0.956731259	1.42e-8	0.0000561777
ENSMUSG00000020108	Ddit4	-1.2769254691	3.85e-8	0.0001143909
ENSMUSG000000032080	Apoa4	-1.2372393607	1.34e-7	0.0003180109
ENSMUSG00000030302	Atp2b2	-1.1445791174	1.526e-7	0.0003293267
ENSMUSG00000029188	Slc34a2	-1.5189323708	2.281e-7	0.0004512668
ENSMUSG00000030244	Gys2	0.5579928642	0.000010983	0.0017383056
ENSMUSG00000021579	Lrrc14b	-0.9323304158	0.000012277	0.0017388094
ENSMUSG00000048521	Cxcr6	0.6073468945	0.000012451	0.0017388094
ENSMUSG00000006574	Slc4a1	-1.4885038062	0.000015841	0.0020765442

Showing 1 to 10 of 10,514 entries

Previous 1 2 3 4 5 ... 1,052 Next

ORGANISM(S): FACTORS: Spaceflight No OF SAMPLES: 112 No OF GENES: 10517 Group Selection Group 1: Ground Control Group 2: Space Flight Modify groups

GeneLab Visualization

Home PCA Volcano Pair plot Group 1 Pair plot Group 2 Heatmap DGE GSEA

GLDS-47&GLDS-48&GLDS-137&GLDS-168&GLDS-173&GLDS-242&GLDS-245 GSEA

Gene sets: MSigDB_Hallmark_2020 Permutations: 1000 Permutation type: Gene set Min size: 15

Max size: 500 Weighted score type: 1 Method: T test Update

Plot type: NES Table Enrichment Score Plot NES Plot Dot Plot Ridge Plot Network Plot GSEA Info

Number of genesets: 50

FDR color threshold: .25

Plot info: FDR ≥ .25 FDR < .25

Pathway NES from GSEA

Showing 50 gene sets sorted by ascending NES.

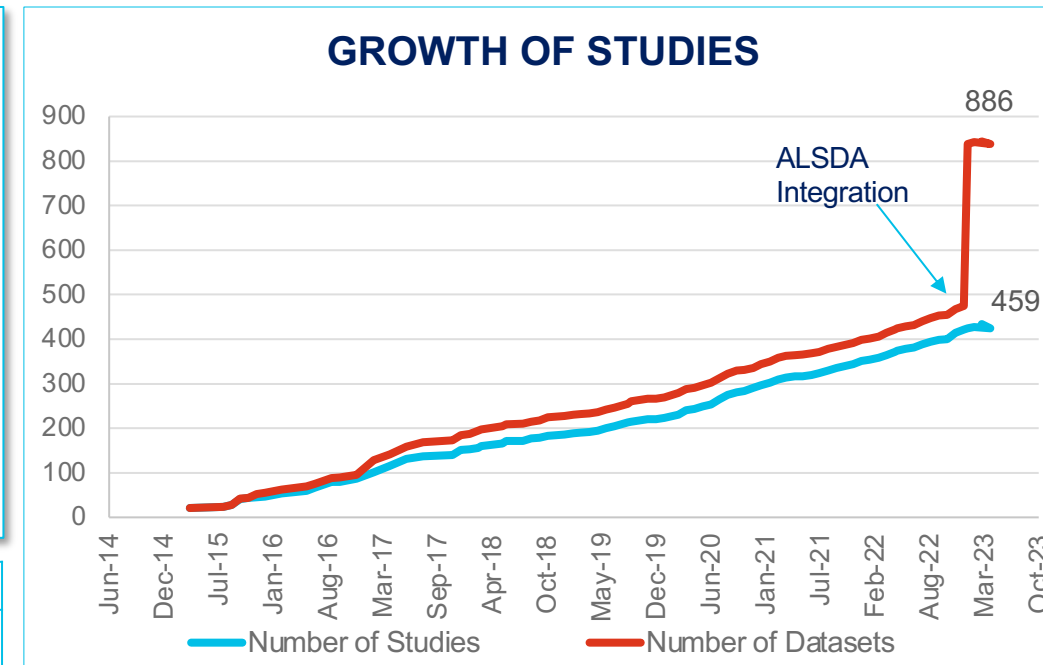
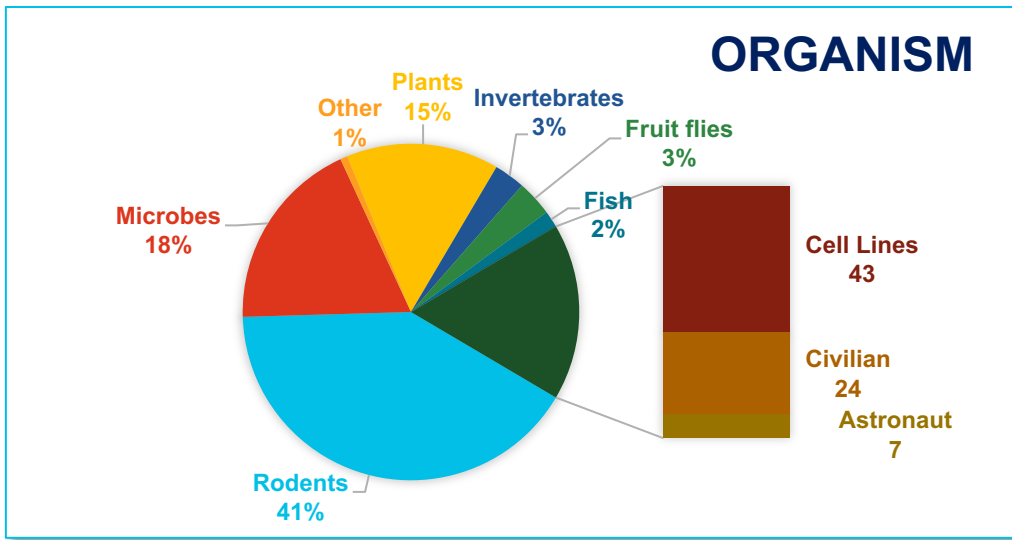
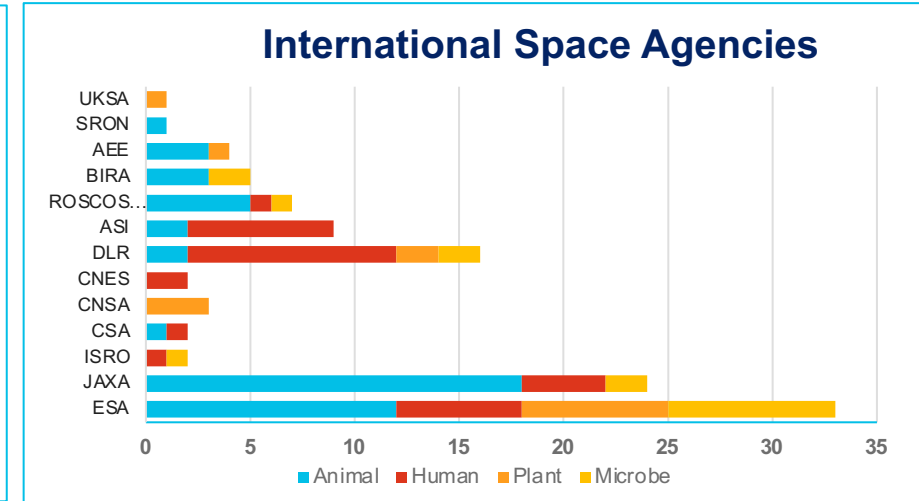
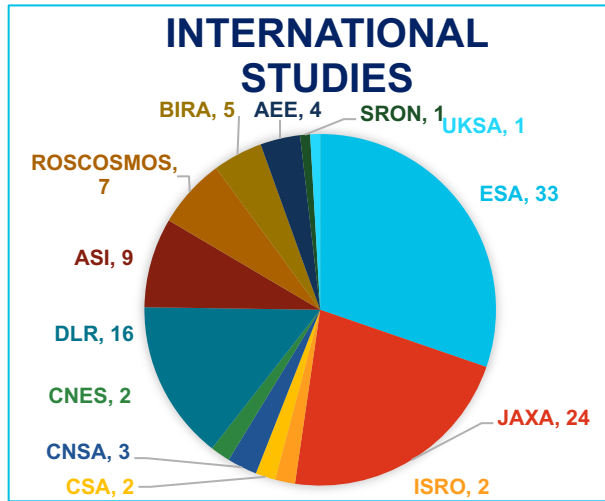
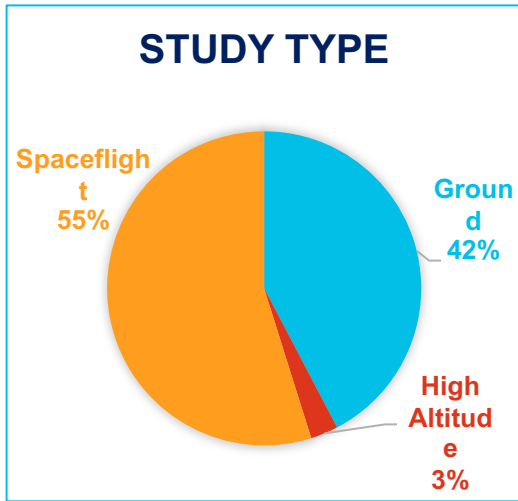
IL-6/JAK/STAT3 Signaling
Inflammatory Response
Coagulation
IL-2/STAT5 Signaling
Interferon Gamma Response
TGF-beta Signaling
p53 Pathway
Apical Surface
DNA Repair
Peroxisome
Notch Signaling
Protein Secretion
Fatty Acid Metabolism
Spermatogenesis
UV Response Dn
Glycolysis
Androgen Response
Pancreas Beta Cells
Hypoxia
Estrogen Response Late
G2-M Checkpoint
Estrogen Response Early
Unfolded Protein Response
KRAS Signaling Dn
Myogenesis

Normalized Enrichment Score

ORGANISM(S): FACTORS: Spaceflight No OF SAMPLES: 112 No OF GENES: 10517 Group Selection Group 1: Ground Control Group 2: Space Flight Modify groups

OSDR Database (GeneLab and ALSDA)

464 Studies
 893 Datasets
 45 Species
 >60 Assays
 >150TB Data

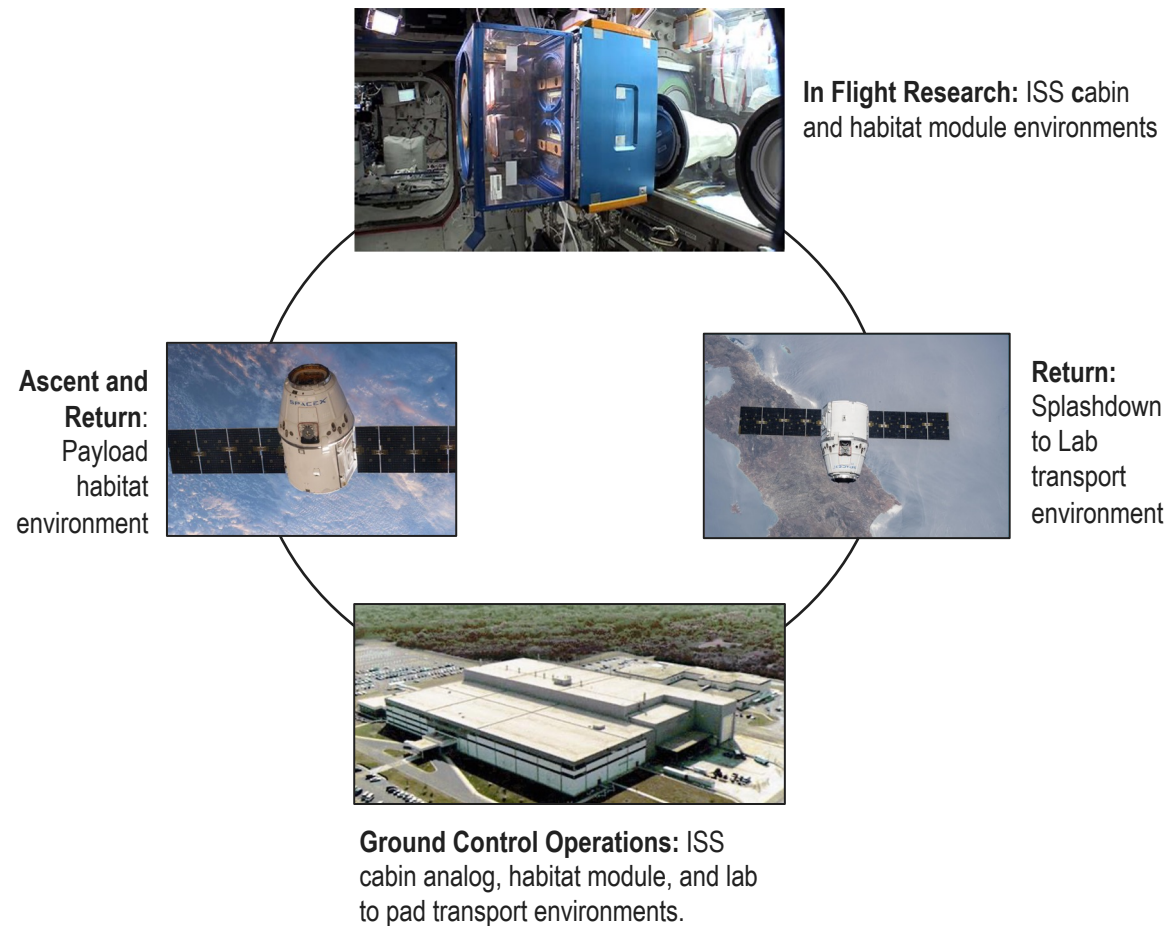


Civilian and Astronaut	Bed Rest, Spaceflight, Mars simulation
Cell Lines	Radiation (Ground), Simulated uG, Spaceflight, Parabolic Flight

Environmental Data Collected From Rodent Research Missions

Developing a data processing pipeline for Mission, Vehicle, and Hardware data. A pipeline that allows data repositories to easily collect, process, and store mission data in both CSV and JSON formats.

Data Collection Points



Relative Humidity
Rodent Transporter
Rodent Habitat
ISS Cabin
ISSES Cabin
LAR Sensor Package

Temperature
Rodent Transporter
Rodent Habitat
ISS Cabin
ISSES Cabin
LAR Sensor Package

Carbon dioxide
ISS Cabin
LAR Sensor Package

Acceleration
Rodent Transporter
LAR Sensor Package

Oxygen
LAR Sensor Package

Radiation
ISS Cabin

Environmental Data Visualization Portal (<https://visualization.osdr.nasa.gov/eda/>)

Environmental Data App

Mission Dashboard

RR-1 ▾

Mission Info
Telemetry Data
Radiation Data

Mission Comparison

Data Tables

Environmental Data App

The Environmental Data App (EDA) is a portal where users can visualize and compare ISS (International Space Station) cabin environmental telemetry data and radiation data gathered from spaceflight missions. On the Mission Dashboard the users can visualize Temperature, Carbon dioxide (CO2), and Relative Humidity measurements and Radiation doses recorded in the ISS cabin and ground control for the duration of the mission. The Mission Comparison feature can be used to visualize data across multiple missions. Raw and summary data is available for download on the Data Tables tab through the Data Tables section.

Rodent Research 1 (SpaceX-4)

NASA's Rodent Research Hardware System provides a research platform aboard the International Space Station (ISS) for long-duration rodent experiments in space. Scientists and engineers at Ames Research Center (ARC) developed the new system for the space station based on the Animal Enclosure Module (AEM) that flew aboard 27 space shuttle missions between 1983 and 2011 and supported studies ranging from four to 18 days. Rodent spaceflight experiments have contributed significantly to our understanding of the effects of microgravity on biological processes that are directly relevant to humans in space. The maiden voyage of the system, Rodent Research-1 (RR1), launched on SpaceX-4 on September 21, 2014 and returned October 25, 2014, during ISS expeditions 41/42. Lasting 37 days, RR1 was the longest duration spaceflight rodent study to date conducted in a NASA facility. RR1 was the first mission to transport rodents aboard an unmanned commercial vehicle. The primary goal is for NASA to validate hardware and demonstrate critical research operations, while supporting the Center for the Advancement of Science in Space (CASIS) in sponsoring the first commercial research study.

Mission Milestone Dates ⓘ

Display Mission milestones

Launch ⓘ: 2014-09-21 09:49

Reached orbit ⓘ: 2014-09-23 17:19

Animal transfer ⓘ: 2014-09-25 16:05

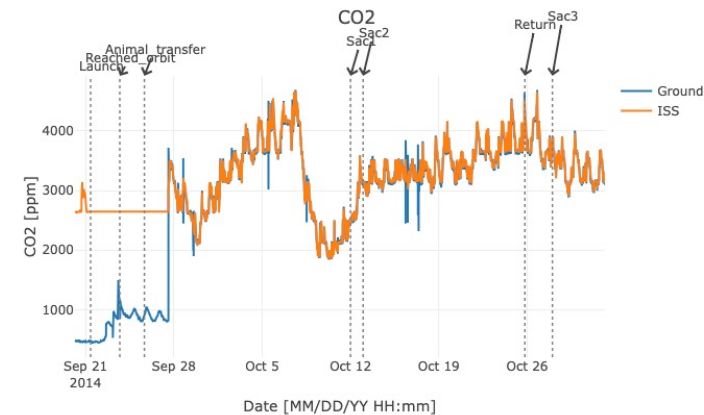
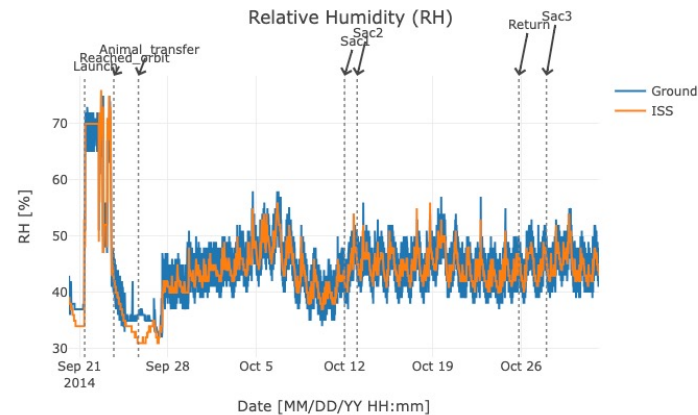
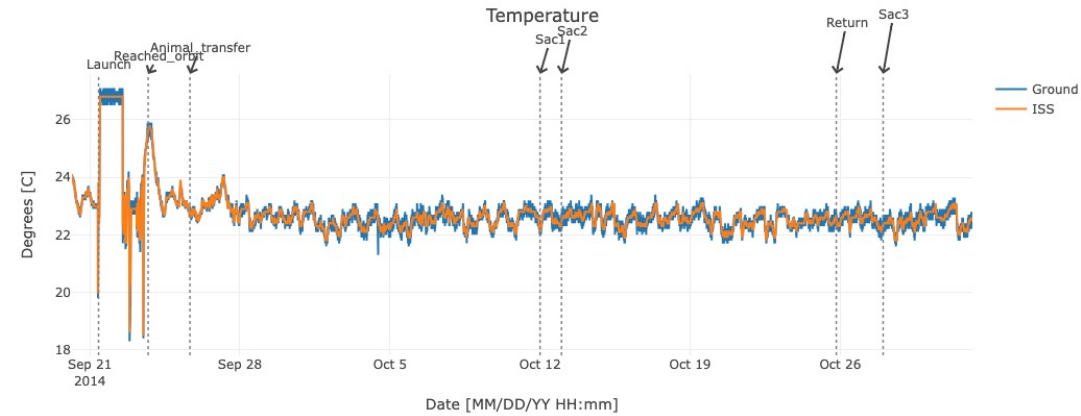
Sac1 ⓘ: 2014-10-12 00:00

Sac2 ⓘ: 2014-10-13 00:00

Return ⓘ: 2014-10-25 19:36

Sac3 ⓘ: 2014-10-28 00:01

Telemetry data ⓘ



Environmental Data Visualization Portal (<https://visualization.osdr.nasa.gov/eda/>)

Environmental Data App

Mission Dashboard

RR-1 ▾

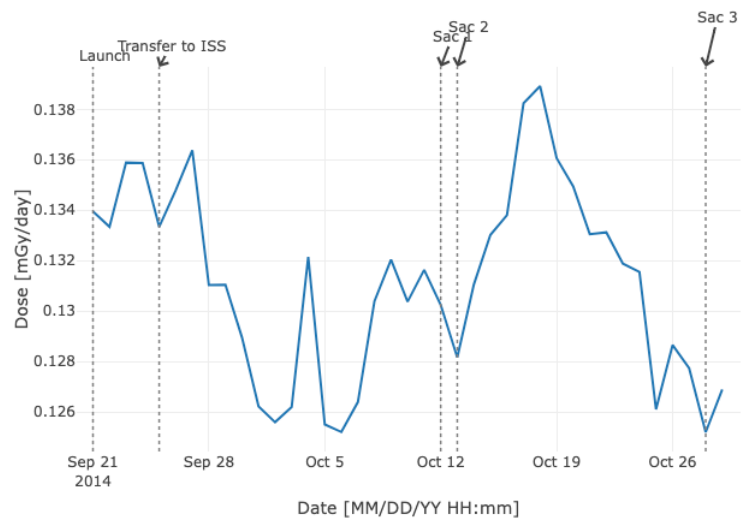
- Mission Info
- Telemetry Data
- Radiation Data

Mission Comparison

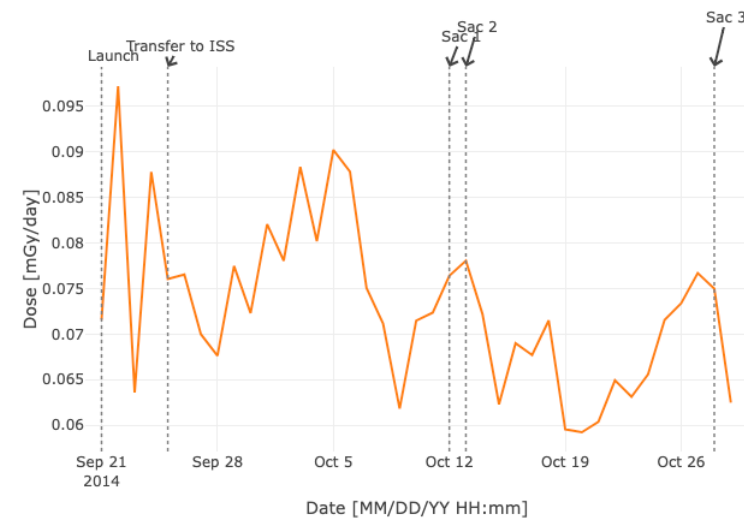
Data Tables

Radiation data ⓘ

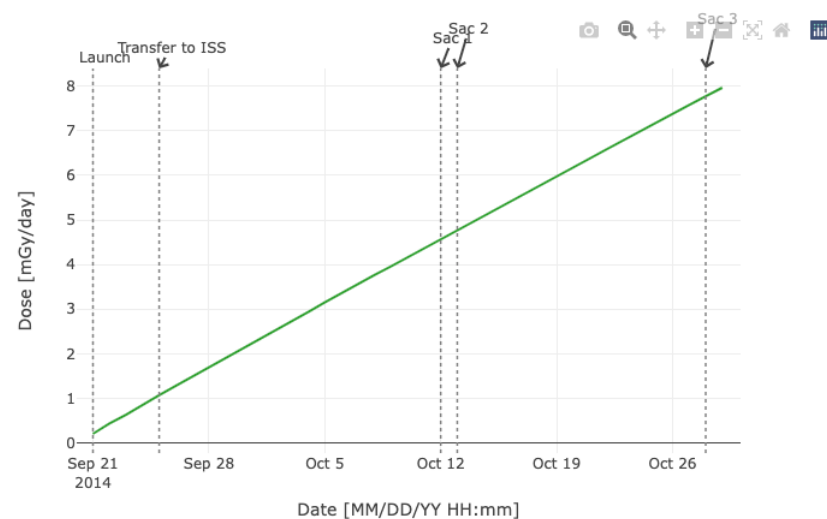
Galactic Cosmic Ray (GCR) ⓘ



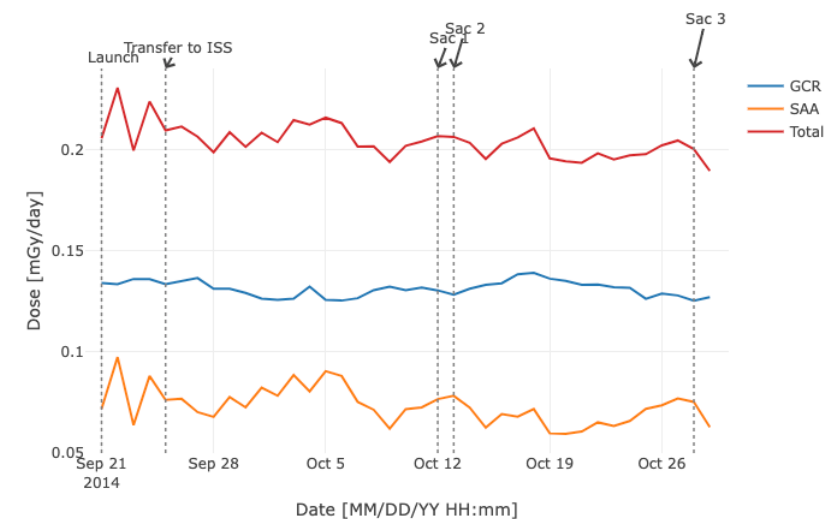
South Atlantic Anomaly (SAA) ⓘ



Accumulated Radiation Dose ⓘ



Total Radiation Dose ⓘ



Open Science Analysis Working Groups (AWGs)

We invite you to join - <https://osdr.nasa.gov/bio/awg/join.html>

ANIMAL

128 members



Facilitates the use of omics in understanding basic mechanisms by which animals and constituent tissues and cells adapt to the spaceflight environment.

MULTI-OMICS

342 members



Interactions between the different omics to provide complete understanding of the entire system being studied.

MICROBES

130 members



Focuses on analyzing microbial datasets within GeneLab that includes gene-expression, proteomic, metabolomic and environmental metagenomic datasets.

PLANTS

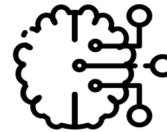
121 members



Share and discuss the latest developments in Astrobotany – the discipline of botany concerned with interactions between plant biology and space environment.

AI/ML

138 members



Focuses on developing data AI-readiness guidelines, algorithm and automation development, and developing ethical guidelines to increase trust and explainability surrounding AI in space biology.

ALSDA

249 members



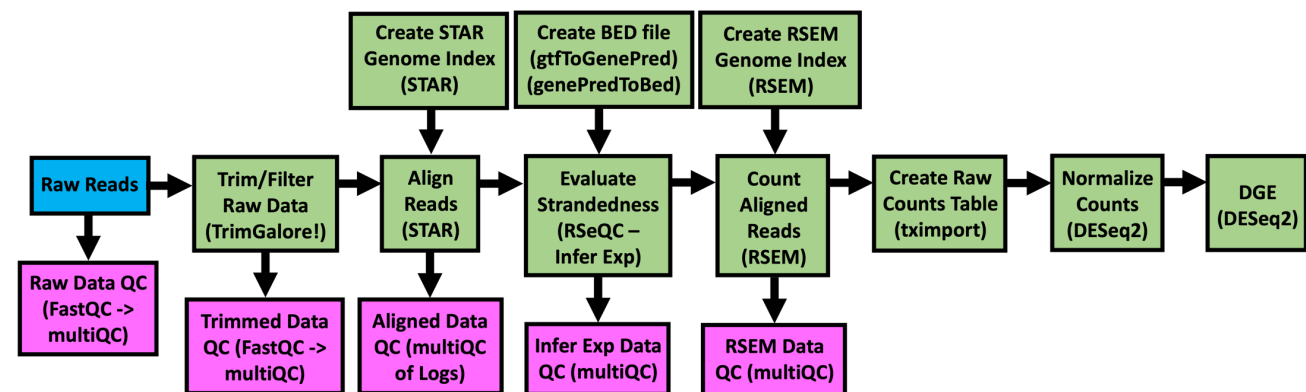
Feedback on science data and metadata standards for physiological, phenotypic, and behavioral datasets to be reusable. Datasets span from raw to processed-results data, and across tabular, bioimaging, and video formats.

Consist of **500+ scientists** from multiple space agencies, international institutions, and industry. Scientists meet monthly with each group to provide feedback, develop standards, and analyze data.

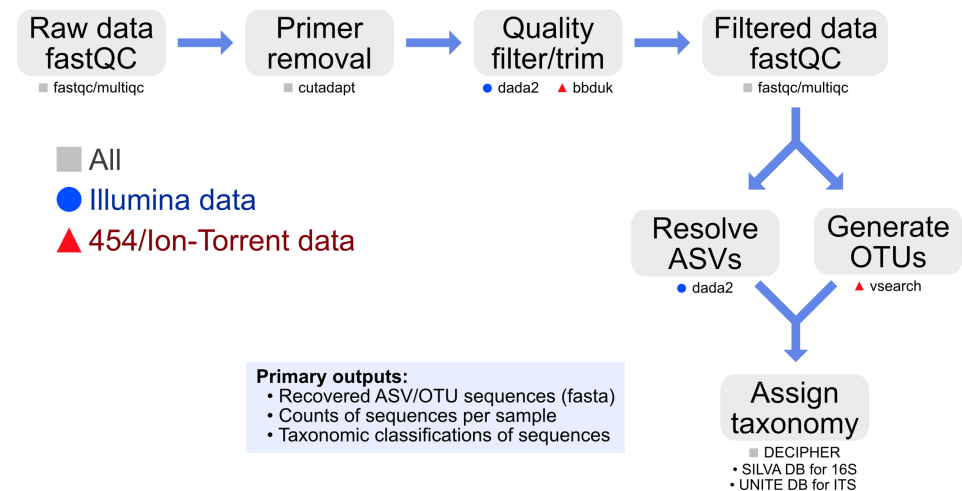
GeneLab Data Processing Pipelines

Build consensus data processing pipelines with the scientific community

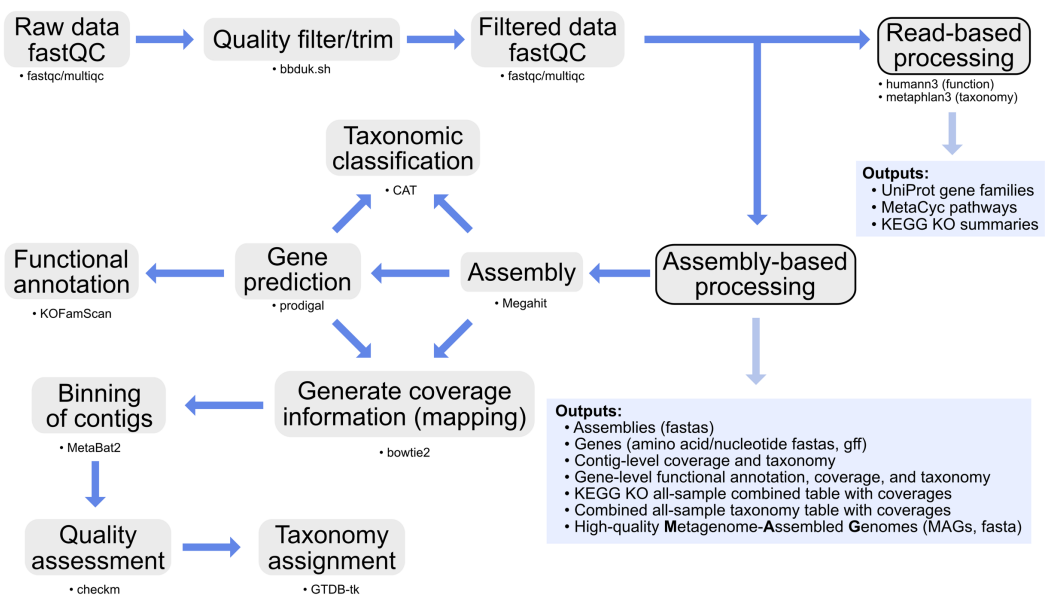
RNA Sequencing Data



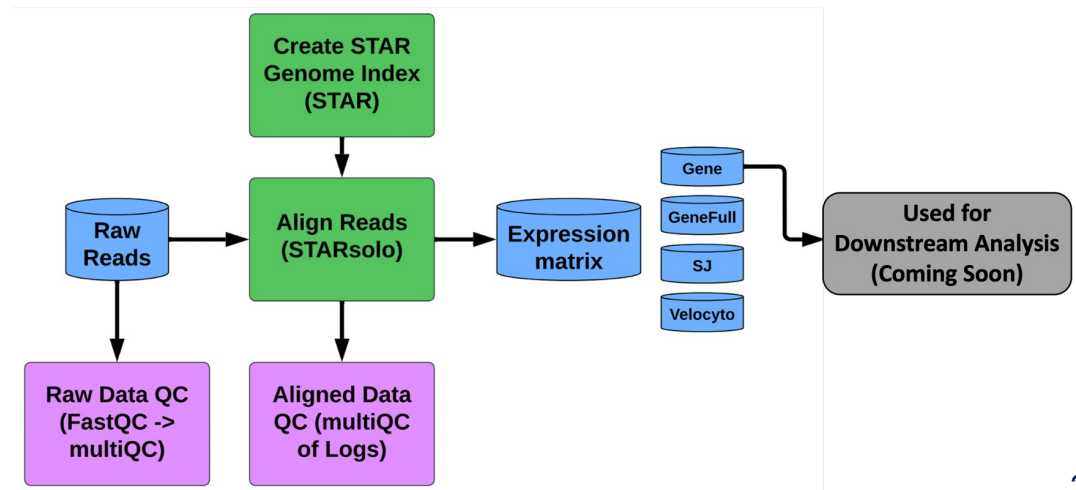
Amplicon Sequencing Data



Metagenomics Data



Single Cell RNA Sequencing Data






GeneLab Data Processing GitHub Repo

The screenshot shows the GitHub repository page for 'nasa / GeneLab_Data_Processing'. The repository is public and contains a README.md file. The README content includes the GeneLab logo, the tagline 'Open Science for Life in Space', the repository name 'GeneLab_Data_Processing', an 'About' section, and an 'Assay Types' section with a list of links to specific data processing pipelines.

Search or jump to... Pull requests Issues Marketplace Explore

nasa / GeneLab_Data_Processing Public Edit Pins

README.md

 Open Science for Life in Space

GeneLab_Data_Processing

About

The [NASA GeneLab](#) Data Processing team and [Analysis Working Group](#) members have created standard pipelines for processing omics data from spaceflight and space-relevant experiments. This repository contains the processing pipelines that have been standardized to date for the assay types indicated below. Each subdirectory in this repository holds current and previous pipeline versions for the respective assay type, including detailed descriptions and processing instructions as well as the exact processing commands used to generate processed data for datasets hosted in the [GeneLab Data Repository](#).

Assay Types

Click on an assay type below for data processing information.

- [Create GeneLab Reference Annotations](#)
- [Amplicon Sequencing](#)
 - [Illumina](#)
 - [454 and Ion-Torrent](#)
- [Metagenomics](#)
 - [Removing human reads](#)
 - [Illumina](#)
- [\(bulk\) RNAseq](#)
- [single cell RNAseq](#)

https://github.com/nasa/GeneLab_Data_Processing



Data Reuse Enables New Discoveries

Return on Investment

90

Original Publication linked to OSDR

60

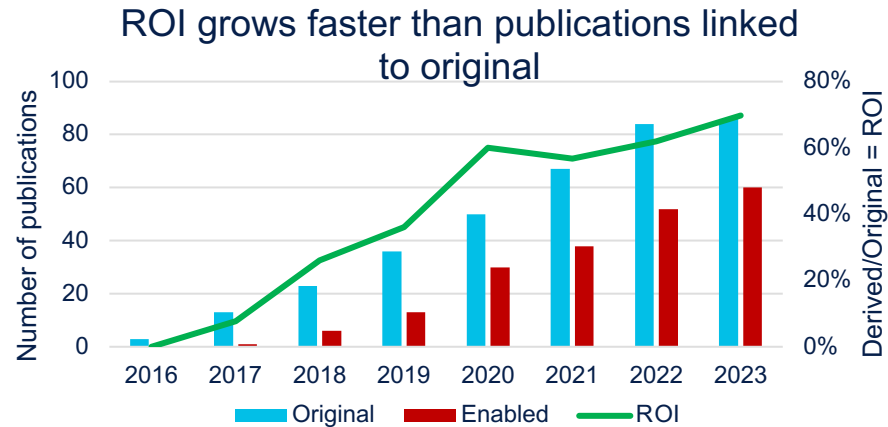
Enabled Publication linked to OSDR

80+

Presentations linked to OSDR

142+

Datasets used in enabled publications



Latest Publications

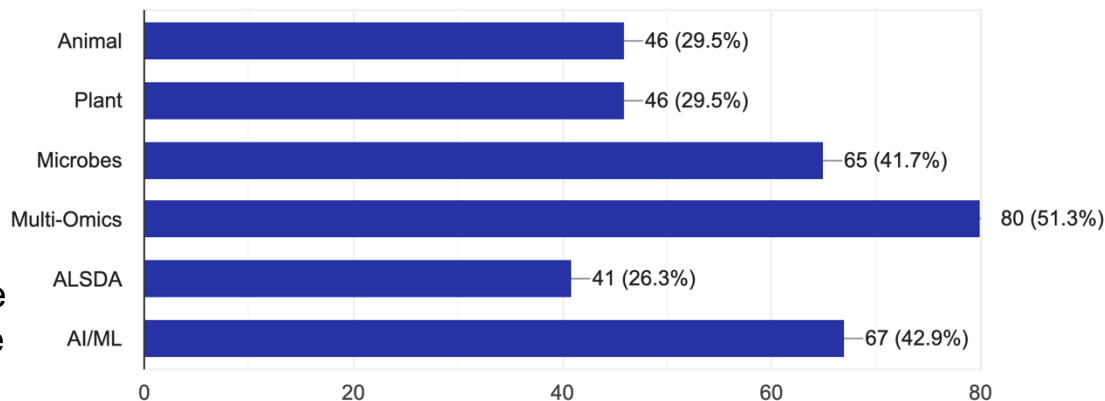
- Microbiomes AWG: [Multidrug-resistant Acinetobacter pittii is adapting to and exhibiting potential succession aboard the International Space Station](#)
- Microbiomes AWG: [Machine learning algorithm to characterize antimicrobial resistance associated with the International Space Station surface microbiome](#)
- Plant AWG: [Meta-analysis of the space flight and microgravity response of the Arabidopsis plant transcriptome](#)
- AI/ML AWG: [Research Prospects of Biomonitoring and Precision Health in Deep Space Supported by Artificial Intelligence](#)
- AI/ML AWG: [Recommendations on Biological Research and Self-driving Labs in Deep Space Supported by Artificial Intelligence](#)
- GeneLab: [Batch effect correction methods for NASA GeneLab transcriptomic datasets](#)

Analysis Working Groups

156

New AWG Members

Members have the option to join more than one AWG



INTERNATIONAL STATUS

ESA OMICS data arriving in OSDR: ~318 TB
Federation of all non-omics data with ESA

Training The Next Generation Of Scientists



GeneLab for High Schools (GL4HS):

A four-week intensive training summer program for rising high school juniors and seniors to learn bioinformatics and computational biology methods and techniques to analyze space omics data.

Learn more and apply at: <https://www.nasa.gov/ames/genelab-for-high-schools>



GeneLab for Colleges/Universities (GL4U):

For educators and students to learn how analyze omics data using GeneLab standard pipelines and space-relevant data

Access to course materials: <https://github.com/nasa/GeneLab-Training/tree/main/GL4U>



Space Life Sciences Training Program (SLSTP):

Provides undergraduate students entering their junior or senior years, and entering graduate students, with professional experience in space life science disciplines.

Learn more and apply at: <https://www.nasa.gov/ames/research/space-life-sciences-training-program>



Spaceflight Technology, Applications, and Research (STAR):

Annual course for PIs, senior research scientists, and postdoctoral scholars and aims to facilitate their entry to space biology and preparation for conducting spaceflight experiments using NASA and commercial platforms.



NASA Postdoctoral Program:

If you're an early-career or senior scientist, apply to the NASA Postdoctoral Program to help pursue NASA's mission and experience the world's most diverse technology and expertise.



Complete NASA's open science curriculum!

Open Science 101:

A community-developed introduction to **core open science skills**

- Know how to write a NASA open science and data management plan
- Learn about tools and best practices
- Increase the impact & visibility of your science
- Earn your digital NASA open science badge

All 5 modules now available through a self-paced online course and through in-person and virtual instructor-led workshops.

Take OS101!



<https://go.nasa.gov/40pPQMx>



TOPS @ AGU 2023



Sessions:

NASA's Transform to Open Science (TOPS) initiative

- Thursday, Dec 14, 2:12 - 2:22 p.m.
- 2014 - West (Level 2, West)

U51A - Science for All - Organizing Agile Openness

- Friday, Dec. 15, 8:30 - 10 a.m.
- 303-304 - South (Level 3, South)

Voices from the Global Open Science Community: Showcasing Progress, Challenges, and Lessons Learned

- Friday, Dec. 15, 8:30 A.M. PST
- 301-302 - South (Level 3, South)

Exhibit Hall:

NASA Exhibit Booth #531

Hyperwall Talks - Tuesday, Dec. 12

2:45 - 3:00 PM - Open Source Science Initiative

3:00 - 3:15 PM - Transforming to Open Science

Town Hall:

NASA Earth Science Division

- Tues, Dec. 12, 6:30 - 8:30 p.m.
- 2020 - West (Level 2, West)



Dec. 11-15, 2023



Moscone Center
San Francisco, CA



Accessing OSDR Data

Open Science Data Repository (https://osdr.nasa.gov/bio/repo/)



Open Science Data Repository Search

Search Datasets

Sort By: Release Date

Items per page: 25 1 - 25 of 435



Persistence of Escherichia coli in the microbiomes of red Romaine lettuce (Lactuca sativa cv. 'Outregeous')- does seed sanitization matter?

Study
OSD-385

Organisms	Factors	Assay Types	Release Date	Description
Microbiota	Treatment Seed Sanitization Tissue	Amplicon Sequencing	19-Apr-2024	Seed sanitization via chemical processes removes/reduces microbes from the external surfaces of the seed and thereby could have an impact on the plants,health or productivity. To determine the impact ...

Highlights: *cgene*



Transcriptional profiling of heart tissue from mice flown on the RRRM-2 mission

Study
OSD-580

Organisms	Factors	Assay Types	Release Date	Description
Mus musculus	Spaceflight Age Euthanasia Location	transcription profiling	03-Jan-2024	In the Rodent Research Reference Mission (RRRM-2), forty female C57BL/6NTac mice were flown on the International Space Station. To assess differences in outcomes due to age, twenty 12 week-old and twe...

Highlights: *cgene*



Transcriptional profiling of tibialis anterior muscle from mice flown on the RR-23 mission

Study
OSD-576

Organisms	Factors	Assay Types	Release Date	Description
Mus musculus	Spaceflight	transcription profiling	12-Dec-2023	The objective of the Rodent Research-23 mission (RR-23) was to better understand the effects of spaceflight on the eyes, specifically on the structure and function of the arteries, veins, and lymphati...

Highlights: *cgene*



Ionizing radiation induces transgenerational effects of DNA methylation in zebrafish

Study
OSD-524

Organisms	Factors	Assay Types	Release Date	Description
Danio rerio	Ionizing Radiation Generation	DNA methylation profiling	31-Aug-2023	Ionizing radiation is known to cause DNA damage, yet the mechanisms underlying potential transgenerational effects of exposure have been scarcely studied. Previously, we observed effects in offspring ...

General Search Filters

Data Source

- GeneLab
- ALSDA
- NIH GEO
- EBI PRIDE
- ANL MG-RAST

Data Type

- Study
- Experiment
- Subject
- Biospecimen
- Payload

Show more

Study Search Filters

Project Type

- Ground
- Spaceflight
- High Altitude


Assay Type

- Amplicon Sequencing Assay
- Bisulfite Sequencing
- ChIP-Seq
- Behavior (Gait)
- Gel Electrophoresis

Show more


Organism

OSD-249 (<https://osdr.nasa.gov/bio/repo/data/studies/OSD-249>)



Version 14

- Description
- Experiments
- Payloads
- Missions
- Protocols
- Samples
- Assays
- Publications
- Files
- Version History
- Visualization





OSD-249 Version 14
Metagenomic analysis of feces from mice flown on the RR-6 mission
Study
1013.43 GB

GeneLab ID: GLDS-249
DOI: [10.26030/h713-bd02](https://doi.org/10.26030/h713-bd02)

[Cite this Study](#)

Submitted Date: 13-Sep-2019
Initial Release Date: 27-Sep-2019

Description

Description

The objective of the Rodent Research-6 (RR-6) study was to evaluate muscle atrophy in mice during spaceflight and to test the efficacy of a novel therapeutic to mitigate muscle wasting. The experiment involved an implantable subcutaneous nanochannel delivery system (nDS; between scapula), which delivered the drug formoterol (FMT; a selective β_2 adrenoceptor agonist), over the course of time. To this end, a cohort of forty 32-weeks-old female C57BL/6NTac mice were either sham operated or implanted with vehicle or treatment-filled nDS, launched in two Transporters (20 mice per Transporter) on SpaceX-13 on December 15, 2017. They were transferred to Rodent Habitats onboard the International Space Station (ISS), and maintained in microgravity for 29 days (N=20, Live Animal Return Spaceflight [LAR FLT]), or >50 days (N=20, ISS Terminal Spaceflight [ISS-T FLT]). After 29 days, the 20 LAR FLT animals were returned live to back to Earth on January 13, 2018. After splashdown, the animals were ambulatory on-ground for ~4 days, until all subjects were processed during one day of dissections. There were two Basal (BSL) groups of animals sacrificed (LAR BSL & ISS-T BSL; N=20; 40 animals; ~36 weeks old) at Kennedy Space Center (KSC; 12/9/17). LAR BSL animals were dissected, and samples were collected upon euthanasia. A Ground Control (GC) group, LAR GC, mimicked the LAR FLT group, which was housed at KSC, then shipped alive, to Novartis's Facilities, where both the LAR FLT and LAR GC groups were processed (~41 weeks old; 1/16/18). All were anesthetized with isoflurane, blood samples were obtained by closed-chest cardiac puncture, and the animals were euthanized by exsanguination and thoracotomy. The 20 ISS-T FLT mice were anesthetized via intraperitoneal injection of ketamine/xylazine/acepromazine over the course of a four days of dissections (2/6/18 until 2/9/18; 53-56 days after launch; 44 weeks old at time of on-orbit dissections). Blood samples and euthanasia were conducted the same as LAR groups. Following blood draw and hind limb dissection, the ISS-T FLT animal carcasses were wrapped in aluminum foil, placed in a ziploc bag and placed in storage at -80°C or colder until return. The ISS-T Ground Control (ISS-T GC) (at KSC) followed the same euthanasia timeline, methods, and preservation. The final processing of frozen ISS-T FLT, frozen ISS-T GC and frozen 0-day ISS-T BSL animals were completed at Houston Methodist Research Institute, in Houston, TX (5/21/18 until 5/24/18). GeneLab received feces from only sham treated animals (no drug treated animals) from the following groups. FLT: LAR (n=9), ISS-T (n=7); GC: LAR (N=7), ISS-T (N=9); BSL: LAR (n=7), ISS-T (n=9). DNA was extracted and analyzed by sequencing using a variety of different targeted and untargeted metagenome profiling assays.

Factor(s)

Factor	Ontology: Concept
Spaceflight	Space Flight
Duration	duration
Dissection Condition	dissection
Euthanasia Location	Euthanasia Location
Ionizing Radiation	ionizing radiation

Organism(s)

[Microbiota](#)

OSD-249: Export Specific Sample Metadata Table Columns

The screenshot displays the OSD-249 interface. On the left is a sidebar with a navigation menu. The 'Samples' menu item is highlighted with a red rectangle. The main content area shows the 'Samples' panel with a 'Select Export Columns' dialog box open. A mouse cursor is pointing at the 'Export CSV' button. The dialog box contains a grid of 32 checkboxes, each corresponding to a specific metadata column. Below the grid are buttons for 'close', 'select all', and 'unselect all'. At the bottom of the dialog, a preview table shows the first two rows of the sample metadata table.

Sidebar Navigation:

- Description
- Experiments
- Payloads
- Missions
- Protocols
- Samples**
- Assays
- Publications
- Files
- Version History
- Visualization

Export CSV Dialog:


Select columns of the samples panel below to export **Export CSV**

- Source Name
- Sample Name
- Characteristics: Organism
- Characteristics: Host Organism
- Characteristics: Host Strain
- Characteristics: Animal Source
- Characteristics: Sex
- Characteristics: Age at Launch
- Characteristics: Diet
- Characteristics: Feeding Schedule
- Characteristics: Material Type
- Factor Value: Spaceflight
- Factor Value: Ionizing Radiation
- Factor Value: Duration
- Factor Value: Euthanasia Location
- Factor Value: Dissection Condition
- Protocol REF
- Parameter Value: NuRFB Lot Number used during Acclimation
- Parameter Value: NuRFB Lot Number used in the Habitat
- Parameter Value: NuRFB Nutrition Values
- Parameter Value: Euthanasia Method
- Parameter Value: Carcass Preservation Method
- Protocol REF
- Parameter Value: Carcass Weight
- Parameter Value: Sample Preservation Method
- Parameter Value: Sample Storage Temperature
- Comment: ALSDA Source Name
- Comment: ALSDA Biospecimen Subject ID
- Comment: ALSDA Biospecimen ID
- Comment: Launch Date
- Comment: Euthanasia Date
- Comment: Dissection Date
- Comment: LSDA Sample Name
- Protocol REF
- Parameter Value: absorbed radiation dose
- Parameter Value: exposure duration
- Parameter Value: absorbed radiation dose rate
- Parameter Value: ionizing radiation source

Preview Table:

Source Name	Sample Name	Characteristics: Organism	Characteristics: Host Organism	Characteristics: Host Strain	Characteristics: Animal Source	Characteristics: Sex	Characteristics: Age at Launch	Characteristics: Diet	Characteristics: Feeding Schedule
LAR BSL 1	Mmus_C57-6T_FCS_BSL_LAR_Rep1_B1	Microbiota	Mus musculus	C57BL/6NTac	Taconic Biosciences	female	36 week	Nutrient Upgraded Rodent Food Bar (NuRFB)	ad libitum
LAR BSL 2	Mmus_C57-	Microbiota	Mus musculus	C57BL/6NTac	Taconic Biosciences	female	36 week	Nutrient Upgraded Rodent Food Bar	ad libitum

OSD-249: Export Specific Assay Metadata Table Columns



Version 14

- [Description](#)
- [Experiments](#)
- [Payloads](#)
- [Missions](#)
- [Protocols](#)
- [Samples](#)
- [Assays](#)
- [Publications](#)
- [Files](#)
- [Version History](#)
- [Visualization](#)

Assays
^

Assay Name: Metagenomic sequencing - Whole-Genome Shotgun Sequencing - Nextera DNA Flex

Technology Ty: Metagenomic sequencing - Whole-Genome Shotgun Sequencing - Swift 1S

Technology Pl: Metagenomic sequencing - Whole-Genome Shotgun Sequencing - LoopSeq


Select Export: Amplicon Sequencing - 16S and ITS - Swift 16S + ITS

Select columns: Amplicon Sequencing - 16S - Fluidigm Access Array

<input checked="" type="checkbox"/> Sample Name	<input checked="" type="checkbox"/> Protocol REF	<input checked="" type="checkbox"/> Parameter Value: QA Instrument	<input checked="" type="checkbox"/> Parameter Value: QA Assay
<input checked="" type="checkbox"/> Parameter Value: QA Score	<input checked="" type="checkbox"/> Parameter Value: DNA size	<input checked="" type="checkbox"/> Extract Name	<input checked="" type="checkbox"/> Protocol REF
<input checked="" type="checkbox"/> Parameter Value: Primer Info	<input checked="" type="checkbox"/> Parameter Value: Library Kit	<input checked="" type="checkbox"/> Parameter Value: library selection	<input checked="" type="checkbox"/> Parameter Value: library layout
<input checked="" type="checkbox"/> Protocol REF	<input checked="" type="checkbox"/> Parameter Value: sequencing instrument	<input checked="" type="checkbox"/> Parameter Value: Read Length	<input checked="" type="checkbox"/> Parameter Value: Read Depth
<input checked="" type="checkbox"/> Assay Name	<input checked="" type="checkbox"/> Raw Data File	<input checked="" type="checkbox"/> Parameter Value: Fastqc File Names	<input checked="" type="checkbox"/> Parameter Value: Multiqc File Names
<input checked="" type="checkbox"/> Protocol REF	<input checked="" type="checkbox"/> Parameter Value: README	<input checked="" type="checkbox"/> Parameter Value: Raw Data	<input checked="" type="checkbox"/> Parameter Value: Read Depth
<input checked="" type="checkbox"/> Parameter Value: Trimmed Sequence Data	<input checked="" type="checkbox"/> Parameter Value: Filtered Sequence Data	<input checked="" type="checkbox"/> Parameter Value: FastQC Outputs	<input checked="" type="checkbox"/> Parameter Value: Final Outputs
<input checked="" type="checkbox"/> Parameter Value: Processing Info			

Sample Name	Protocol REF	Parameter Value: QA Instrument	Parameter Value: QA Assay	Parameter Value: QA Score	Parameter Value: DNA size	Extract Name	Protocol REF	Parameter Value: Primer Info	Pararr Kit
Mmus_C57-6T_FCS_BSL_LAR_Rep1_B1	nucleic acid extraction	Agilent 4200 TapeStation	gDNA ScreenTape Assay	6.6 DNA Integrity Number	10118 base pair	Mmus_C57-6T_FCS_BSL_LAR_Rep1_B1	library construction-Fluidigm Access	515F, 5'-GTGYCAGCMGCCGCGGTAA-3', 806R, 5'-GGACTACNVGGGTWTCTAAT-3'	Fluid

OSD-249: Download Files



Version 14

- Description
- Experiments
- Payloads
- Missions
- Protocols
- Samples
- Assays
- Publications
- Files**
- Version History
- Visualization

Files

Study Files Selected: 96

[Download](#)

Search Files

- OSD-249
 - Amplicon Sequencing
 - Supplemental Materials
 - Study Metadata Files
 - Metagenomics Data Files
 - GeneLab Processed Metagenomics Files
 - GeneLab Processed Diversity Amplicon Files
 - README
 - Raw Sequence Data
 - GLDS-249_GAmplicon_Mmus_C57-6T_FCS_BSL_ISS-T_FluidAA_Rep1_B1_R1_raw.fastq.gz 3.26 MB Tue Apr 26 2022
 - GLDS-249_GAmplicon_Mmus_C57-6T_FCS_BSL_ISS-T_FluidAA_Rep1_B1_R2_raw.fastq.gz 3.55 MB Tue Apr 26 2022
 - GLDS-249_GAmplicon_Mmus_C57-6T_FCS_BSL_ISS-T_FluidAA_Rep2_B2_R1_raw.fastq.gz 3.54 MB Tue Apr 26 2022
 - GLDS-249_GAmplicon_Mmus_C57-6T_FCS_BSL_ISS-T_FluidAA_Rep2_B2_R2_raw.fastq.gz 3.74 MB Tue Apr 26 2022
 - GLDS-249_GAmplicon_Mmus_C57-6T_FCS_BSL_ISS-T_FluidAA_Rep3_B3_R1_raw.fastq.gz 3.34 MB Tue Apr 26 2022
 - GLDS-249_GAmplicon_Mmus_C57-6T_FCS_BSL_ISS-T_FluidAA_Rep3_B3_R2_raw.fastq.gz 3.63 MB Tue Apr 26 2022
 - GLDS-249_GAmplicon_Mmus_C57-6T_FCS_BSL_ISS-T_FluidAA_Rep4_B4_R1_raw.fastq.gz 2.73 MB Tue Apr 26 2022
 - Trimmed Sequence Data
 - Filtered Sequence Data
 - FastQC Outputs
 - Final Outputs
 - Processing Info

Accessing OSDR Data – Public API (<https://genelab.nasa.gov/genelabAPIs>)



Open Science for Life in Space

[Home](#)

[About](#)

[Data & Tools](#)

[Research & Resources](#)

[Working Groups](#)

[Help](#)

Keywords



GeneLab Public API

NASA GeneLab provides a RESTful Application Programming Interfaces (API) to its full-text search, and data and metadata retrieval capabilities. The API provides a choice of standard web output formats, either JavaScript Object Notation (JSON) or Hyper Text Markup Language (HTML), of query results. The GeneLab Data Query API returns metadata on data files associated with dataset(s), including the location of these files for download via https. The GeneLab Metadata Query API returns entire sets of metadata for input dataset accession numbers. The GeneLab Dataset Search API can be used to search dataset metadata by keywords and/or metadata. It can also be used to provide search of three other omics databases: the National Institutes of Health (NIH) / National Center for Biotechnology Information's (NCBI) Gene Expression Omnibus (GEO); the European Bioinformatics Institute's (EBI) Proteomics Identification (PRIDE); the Argonne National Laboratory's (ANL) Metagenomics Rapid Annotations using Subsystems Technology (MG-RAST).

Contents:

- [Data Query API](#)
- [Metadata API](#)
- [Search API](#)

Data Query API

Syntax

https://genelab-data.ndc.nasa.gov/genelab/data/glds/files/{GLDS_STUDY_IDS}?page={CURRENT_PAGE_NUMBER}&size={RESULTS_PER_PAGE}

Returns: JSON-formatted response

Parameters	Data Type	Notes	Values	Required
{GLDS_STUDY_IDS}	Integers	Comma separated list with mixture of single GLDS accession numbers and ranges	ex. 87-95,137	Yes
{CURRENT_PAGE_NUMBER}	Integer	Current page number in pagination	Starts from 0	No
{RESULTS_PER_PAGE}	Integer	Number of results returned per page in pagination	Max 25 results per page	No

Example requests:

- Single study request using GLDS accession number
 - <https://genelab-data.ndc.nasa.gov/genelab/data/glds/files/87>
- Multiple studies request using combination of range and comma separated list
 - <https://genelab-data.ndc.nasa.gov/genelab/data/glds/files/137,87-95,13,20-50>

NOTE: `study_files` element in the JSON response has the `remote_url` attribute, which can be used to obtain the specific download URL for the file by prefacing with the GeneLab data server address, <https://genelab-data.ndc.nasa.gov>. In the example query/response below, the first study file for GLDS-87 study in the response below can be downloaded from https://genelab-data.ndc.nasa.gov/datamanager/file/Home/genelab/genelab-data/GLDS-87/metadata/GLDS-87_metadata.Zanello_STS135-ISA.zip

Example Requests:


Single Study Request:

<https://genelab-data.ndc.nasa.gov/genelab/data/glds/files/87>

Response:

```
{
  "hits": 1,
  "input": "87",
  "studies": {
```

Accessing OSDR Data Files – Public s3 Bucket (<https://registry.opendata.aws/nasa-osdr/>)

Registry of Open Data on AWS 

NASA Space Biology Open Science Data Repository (OSDR)

[biinformatics](#) [biology](#) [GeneLab](#) [genomic](#) [imaging](#) [life sciences](#) [space biology](#)

Description

NASA's Space Biology Open Science Data Repository (OSDR) introduces a one-stop site where users can explore and contribute a variety of NASA open science biological data. This site consolidates data from the Ames Life Sciences Data Archive (ALSDA) and GeneLab and includes information about the broader NASA Open Science and Open Data initiatives, all at one centralized location. Our mission is to maximize the utilization of the valuable biological research resources and enable new discoveries.

OSDR introduces access to data generated from spaceflight and space relevant experiments that explore the biological response of terrestrial biology through the AWS Open Data Registry page. The ALSDA is the official repository of non-human science data spanning a broad range of biological levels involving data from tissues, organs, whole organisms, physiology, and behavior. GeneLab is an open science repository hosting multiple types of 'omics including transcriptomics, metagenomics, epigenomics, proteomics, and metabolomics data. Studies comprise of data from model organisms including microbes, plants, fruit flies, rodents, as well as human cell culture, ground study, and commercial astronaut data. In addition, the data repository includes metadata searches across several external omics database.

Resources on AWS

Description
Biological research data from spaceflight and space relevant experiments

Resource type
S3 Bucket

Amazon Resource Name (ARN)
`arn:aws:s3:::~nasa-osdr`

AWS Region
`us-west-2`

AWS CLI Access (No AWS account required)
`aws s3 ls --no-sign-request s3://nasa-osdr/`

[Explore](#)
[Browse Bucket](#)

Update Frequency

New research data is added as soon as it is available.


License

There are no restrictions on the use of this data.

Documentation

<https://osdr.nasa.gov/bio/repo>

Managed By



See all datasets managed by [NASA](#).

Contact

<https://osdr.nasa.gov/bio/help/contact.html>

How to Cite

NASA Space Biology Open Science Data Repository (OSDR) was accessed on `DATE` from `https://registry.opendata.aws/nasa-osdr`.

Usage Examples

Publications

- [Advancing the Integration of Biosciences Data Sharing to Further Enable Space Exploration](#) by Ryan T. Scott, Kirill Grigorev, Graham Mackintosh, Samrawit G. Gebre, Christopher E. Mason, Martha E. Del Alto, Sylvain V. Costes
- [GeneLab: Omics database for spaceflight experiments](#) by Shayoni Ray, Samrawit Gebre, Homer Fogle, Daniel C Berrios, Peter B Tran, Jonathan M Galazka, Sylvain V Costes
- [NASA GeneLab: interfaces for the exploration of space omics data](#) by Daniel C Berrios, Jonathan Galazka, Kirill Grigorev, Samrawit Gebre, Sylvain V Costes

Accessing OSDR Data Files – Public s3 Bucket (https://registry.opendata.aws/nasa-osdr/)

AWS S3 Explorer interface showing the root of the 'nasa-osdr' bucket. The breadcrumb path is 'nasa-osdr'. The search bar is empty. The table lists 452 objects, all of which are folders (indicated by a trailing slash).

Object	Last Modified	Size
OSD-1/		
OSD-100/		
OSD-101/		
OSD-102/		
OSD-103/		
OSD-104/		
OSD-105/		
OSD-106/		
OSD-107/		
OSD-108/		
OSD-109/		
OSD-11/		
OSD-110/		
OSD-111/		
OSD-112/		
OSD-113/		
OSD-114/		
OSD-115/		
OSD-116/		
OSD-117/		
OSD-118/		
OSD-119/		
OSD-12/		
OSD-120/		
OSD-121/		
OSD-122/		
OSD-123/		

AWS S3 Explorer interface showing a filtered view of the 'nasa-osdr' bucket. The breadcrumb path is 'nasa-osdr / OSD-249 / version-14 / GAmplicon'. The search bar contains the filter '_raw.fastq.gz' and is highlighted with a red box. The table lists 303 objects, all of which are files (indicated by a file icon).

Object	Last Modified	Size
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_LAR_FluidAA_Rep7_G10_R2_raw.fastq.gz	5 months ago	3 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_LAR_FluidAA_Rep7_G10_R1_raw.fastq.gz	5 months ago	3 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_LAR_FluidAA_Rep6_G9_R2_raw.fastq.gz	5 months ago	4 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_LAR_FluidAA_Rep6_G9_R1_raw.fastq.gz	5 months ago	3 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_LAR_FluidAA_Rep5_G8_R2_raw.fastq.gz	5 months ago	4 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_LAR_FluidAA_Rep5_G8_R1_raw.fastq.gz	5 months ago	3 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_LAR_FluidAA_Rep4_G5_R2_raw.fastq.gz	5 months ago	3 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_LAR_FluidAA_Rep4_G5_R1_raw.fastq.gz	5 months ago	3 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_LAR_FluidAA_Rep3_G4_R2_raw.fastq.gz	5 months ago	3 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_LAR_FluidAA_Rep3_G4_R1_raw.fastq.gz	5 months ago	3 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_LAR_FluidAA_Rep2_G3_R2_raw.fastq.gz	5 months ago	3 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_LAR_FluidAA_Rep2_G3_R1_raw.fastq.gz	5 months ago	3 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_LAR_FluidAA_Rep1_G2_R2_raw.fastq.gz	5 months ago	3 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_LAR_FluidAA_Rep1_G2_R1_raw.fastq.gz	5 months ago	3 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_ISS-T_FluidAA_Rep9_G10_R2_raw.fastq.gz	5 months ago	3 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_ISS-T_FluidAA_Rep9_G10_R1_raw.fastq.gz	5 months ago	2 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_ISS-T_FluidAA_Rep8_G9_R2_raw.fastq.gz	5 months ago	3 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_ISS-T_FluidAA_Rep8_G9_R1_raw.fastq.gz	5 months ago	2 MB
GLDS-249_GAmplicon_Mmus_CS7-6T_FCS_GC_ISS-T_FluidAA_Rep7_G8_R2_raw.fastq.gz	5 months ago	3 MB

Amplicon (Illumina) Sequencing Overview

Common Sequencing Approaches in Microbial Ecology

Amplicon sequencing involves the use of “primers” – short sequences that match highly conserved stretches of DNA

These primers bind to and help amplify DNA from microbes in our sample, usually targeting a portion of a single gene

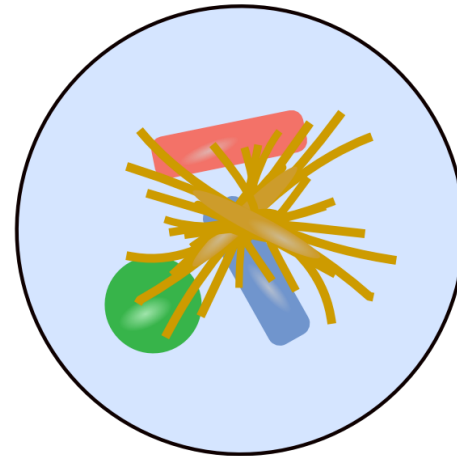
We then use these recovered sequences from our sample to try to get information about the microbial community that is present (often trying to figure out “who” is there)

Metagenomics sequencing involves attempting to sequence “all” the DNA present in a system

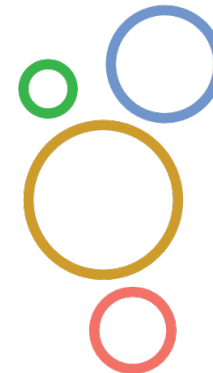
We also then try to use these recovered DNA sequences to try to get information about the microbial community present in our sample

With metagenomics sequencing, in addition to getting insight into which microbes are present, we will often also try to understand “what” they might be doing functionally

Mixed microbial community



DNA
Extraction



Amplicon sequencing



Multiple copies of a portion
of 1 target gene

Metagenomics sequencing



Short sequences from “all”
DNA

Amplicon Sequencing Overview: Select the Target Region

Amplicon sequencing involves the use of “primers” – short sequences that match highly conserved stretches of DNA

These primers most often target a gene that is universally present across all known life, while also being useful for delineating between different organisms. These target genes are often referred to as “marker-genes”.

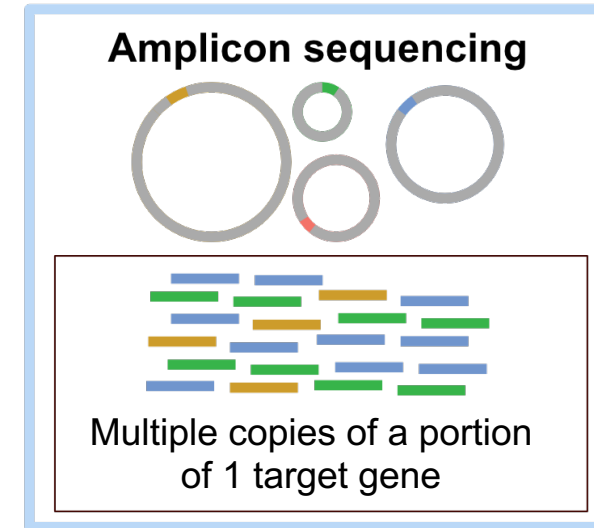
The most common marker-gene target is part of the “ribosome”, specifically the small subunit ribosomal RNA gene (SSU rRNA), known as the “16S” in bacteria and archaea and the “18S” in eukarya

The “ribosome” is an assemblage of multiple proteins and RNA molecules, and it is essential for protein synthesis in all known bacteria, archaea, and eukarya

Because many molecules of proteins and RNA (including the SSU rRNA) need to interact together to function properly as a full ribosome, there is a lot of evolutionary pressure functionally constraining them, which results in the DNA that encodes for these molecules having relatively conserved sequences

The SSU rRNA itself is particularly suited for amplicon analysis because it has both:

1. regions that are highly conserved sequence-wise (making them good targets for the primers to find them and amplify them)
2. regions that are highly variable sequence-wise (making them good for delineating between different organisms)



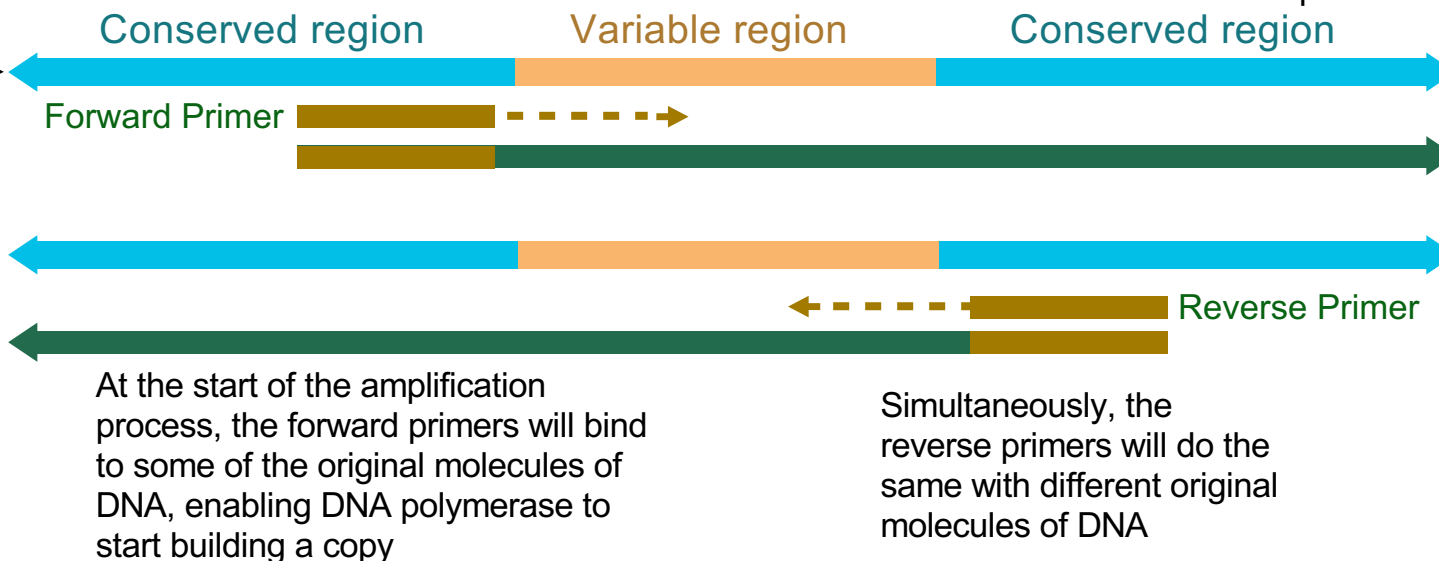
Amplicon Sequencing Overview: Amplifying the Target Region

A typical example of amplifying an SSU rRNA gene via PCR:

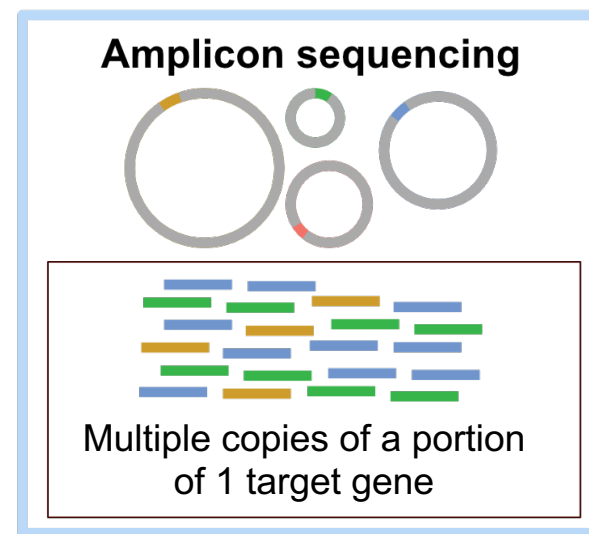
This represents a portion of the common "16S" marker-gene

It starts out as part of a longer portion of DNA that came from an organism's genome

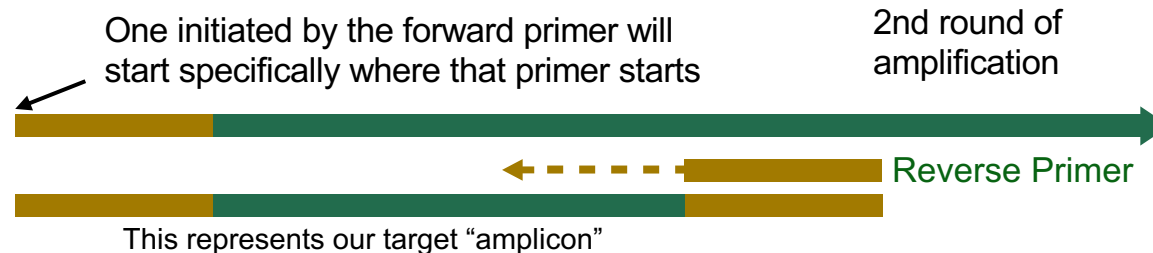
The primers we put in target the region we want to amplify



1st round of amplification



In the next round of amplification, there are going to be many copies that were initiated by either the forward or reverse primer



2nd round of amplification

When a copy like this is then amplified via the "other" primer, the reverse primer in this case, we will start getting copies of just our target "amplicon"

After multiple rounds of amplification, there will be many copies of just our target amplicon, which is what we then actually sequence

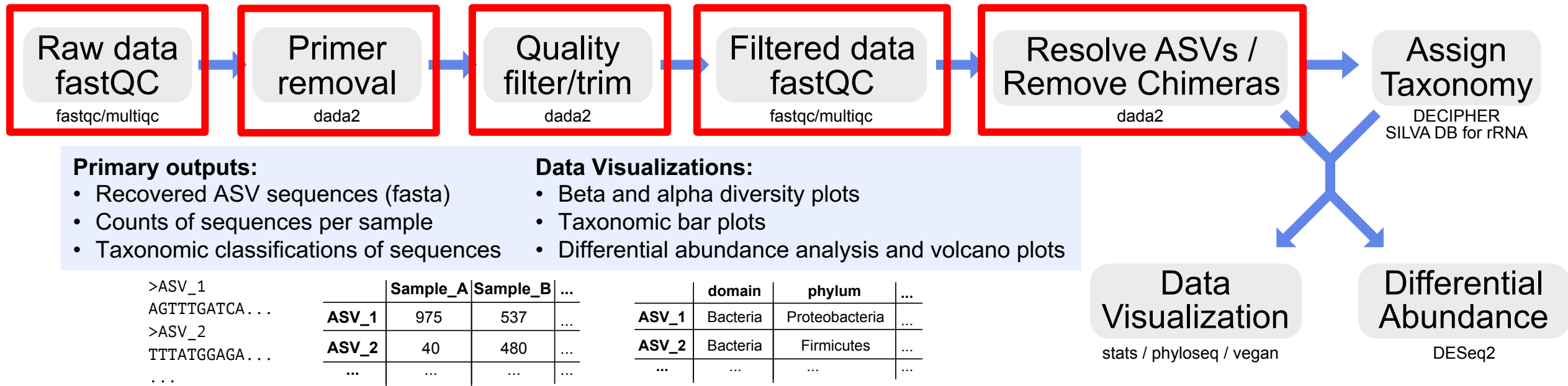
They are all colored the same here, but the initial DNA molecules that were first amplified came from different initial biological sequences, so there are actually many copies of different initial DNA sequences here



After multiple rounds of amplification

GeneLab Amplicon (Illumina) Sequencing Data Processing Pipeline

GeneLab Amplicon (Illumina) Sequencing Data Processing Pipeline

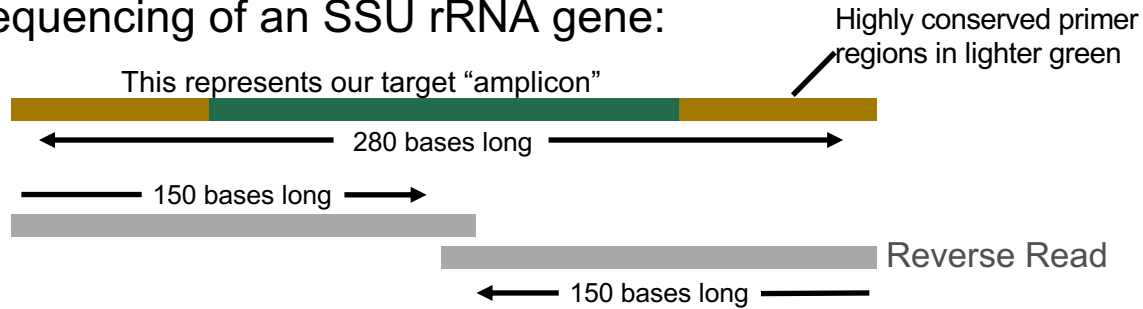


Processing AmpSeq Data: Resolve and Quantify ASVs

A typical example of “paired-end” sequencing of an SSU rRNA gene:

It’s all of these copies of our target amplicon, from organisms in our original sample, that we then sequence

Often, the target amplicon is longer than the “reads” the sequencer is generating



Part of processing amplicon data usually involves combining these forward and reverse reads in order to reconstruct the full initial target amplicon



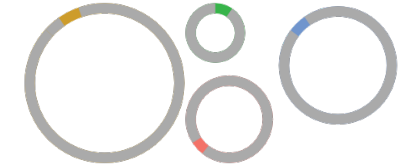
“Paired-end” sequencing is commonly employed to enable us to span the entire target amplicon with a forward and reverse read

In order to be able to combine our reads, we need to be sure that together they span the full length of the target amplicon

Following primer-region removal, these recovered unique amplicon sequences (known as **Amplicon Sequence Variants, ASVs**) would typically then be counted to see how many times each unique sequence appears in each sample, creating a table such as this:

	Sample_A	Sample_B	...
Seq_1	975	537	...
Seq_2	40	480	...
...

Amplicon sequencing



Multiple copies of a portion of 1 target gene

In this example, we have 2×150 reads = 300 bases
 And a nominal target amplicon length of 280
 $300 - 280 = 20$
 So we have ~20 bases of expected overlap
 (there is no “one” cutoff, but a minimum of 8 is reasonable)

These counts we get from processing amplicon data represent *counts of recovered gene-copies*.

They do not represent counts of organisms or counts of genomes.

What Are Chimeras?

“Chimeras” in the sequencing world are artificial sequences that are created during the amplification process when 2 or more biological sequences join together



As discussed earlier, a primer would be what enables these sequences to be amplified via PCR:



But sometime that amplification process can prematurely terminate:



That partial sequence may then act as a “primer” by attaching to a different sequence, even with some mismatches between them, and then being extended based on this different sequence:



Ultimately yielding a chimeric sequence, a mixture of both, that then continues to be amplified:



Amplicon sequencing

Multiple copies of a portion of 1 target gene

Processing AmpSeq Data: Remove Chimeras

“Chimeras” in the sequencing world are artificial sequences that are created during the amplification process when 2 or more biological sequences join together



Chimeras are often very common in amplicon datasets in terms of the *unique* sequences recovered, but not usually in terms of the *total* sequences recovered.

Remember our example count table from earlier:

	Sample_A	Sample_B	...
Seq_1	975	537	...
Seq_2	40	480	...
...

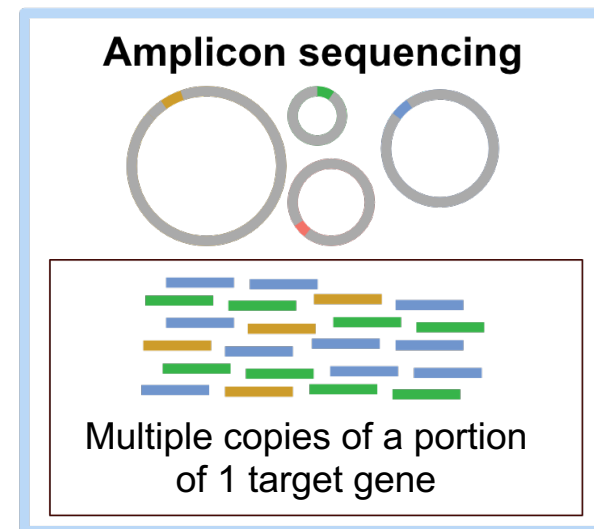
The total number of *unique* sequences would be however many rows are in this table (2 showing)

The total count of *all* sequences recovered would be the sum of all the values in the cells (2,032 showing)

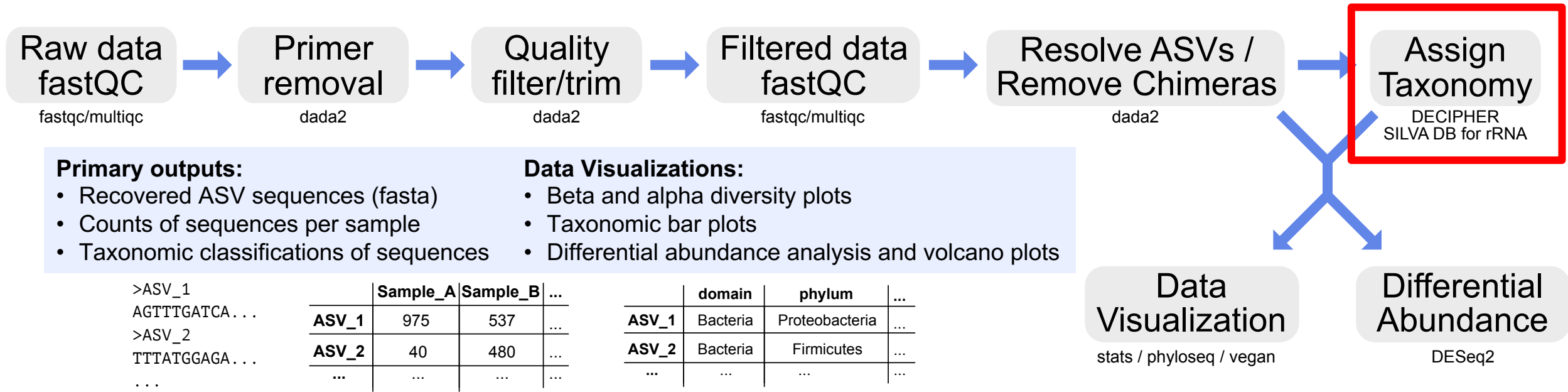
Chimeras may frequently make up a large proportion of *unique* sequences recovered, but those unique chimeric sequences are usually not seen that many times each, and therefore make up only a small fraction of the *total* sequences recovered.

It is important to perform a step that tries to identify and remove chimeras, as they otherwise can drastically, artificially increase the apparent diversity of sequences that appears to exist in a sample.

Automated methods of chimera detection typically involve checking if a sequence can be made exactly by mixing together other sequences.



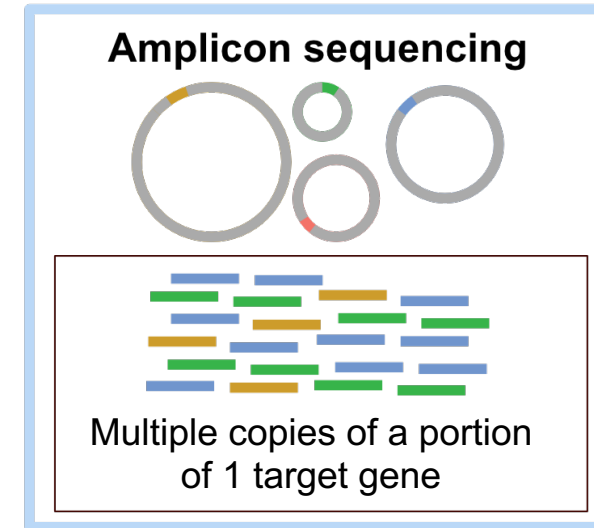
GeneLab Amplicon (Illumina) Sequencing Data Processing Pipeline



Processing AmpSeq Data: Assigning Taxonomy

Once we have some set of unique sequences recovered from our samples that we want to move forward with, it is common to try to assign a taxonomy to each of them, so we can get an idea of which microbes are present in our system

This is typically done by comparing our recovered sequences to a reference-sequence database in an effort to identify the taxonomy of the source organism each sequence likely came from



Like all things that use a reference database and any method, the results of a given taxonomic classification are entirely dependent on the reference database and method used

Moreover, taxonomic classification of a short portion of a single gene is difficult, and can often even be impossible assign a taxonomy at the species-level because there are cases where the portion we sequenced could be identical between different organisms

It's also worth remembering that taxonomic delineations are arbitrary; they are not fundamental units of biology that actually exist. Biological taxonomy is an ongoing process of trying to define things and attach agreed-upon labels to them.

That doesn't negate its utility, of course!

While there can be exceptions everywhere, and we don't really know the functions of some organism based on its taxonomy alone, taxonomic labels still provide a requisite means of communication and can help guide hypotheses.

GeneLab Amplicon (Illumina) Sequencing Data Processing Pipeline



Primary outputs:

- Recovered ASV sequences (fasta)
- Counts of sequences per sample
- Taxonomic classifications of sequences

Data Visualizations:

- Beta and alpha diversity plots
- Taxonomic bar plots
- Differential abundance analysis and volcano plots

```

>ASV_1
AGTTTGATCA...
>ASV_2
TTTATGGAGA...
...
  
```

	Sample_A	Sample_B	...
ASV_1	975	537	...
ASV_2	40	480	...
...

	domain	phylum	...
ASV_1	Bacteria	Proteobacteria	...
ASV_2	Bacteria	Firmicutes	...
...



Processing AmpSeq Data: Data Visualization – Beta Diversity

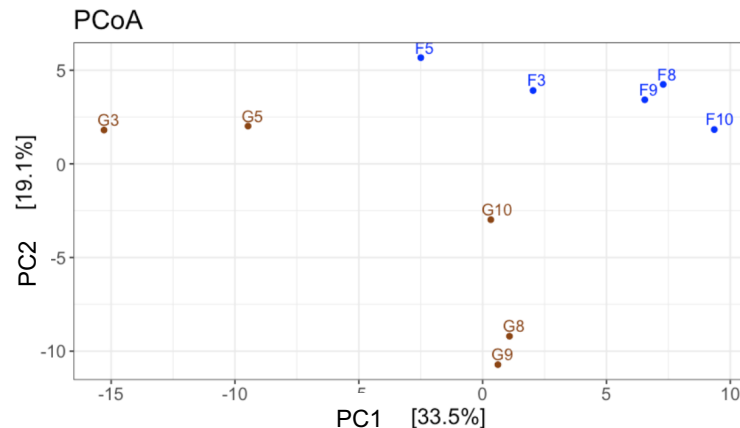
- Beta diversity involves calculating metrics such as distances or dissimilarities based on how samples relate to each other
- Important decisions to consider:
 - How the data are normalized (to account for differences in sample read-depth)
 - What is used for calculating distance/dissimilarity
 - How the data are ordinated and/or clustered

Normalization

- We use DESeq2's variance stabilizing transformation (VST) to normalize ASV counts

Ordination Method

- Ordinations provide visualizations of sample-relatedness based on dimension reduction. In our case, the dimensions are ASV counts.
- We perform **Principle Coordinates Analysis (PCoA)** on the normalized ASV counts using the *phyloseq* package

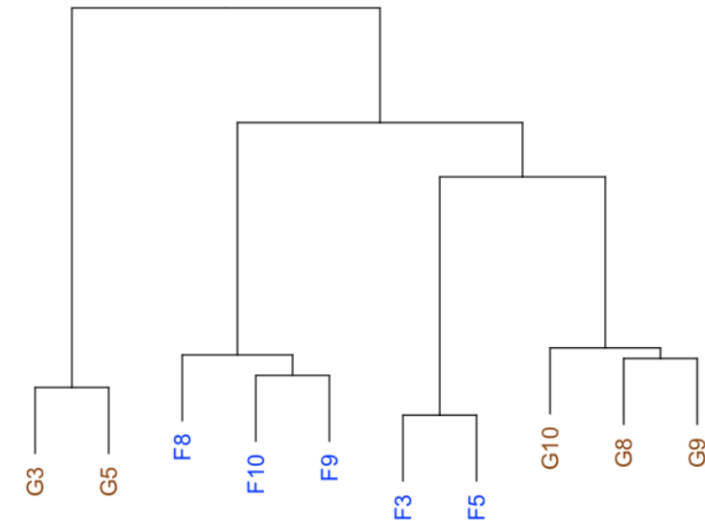


Distance/Dissimilarity Metric

- We calculate a Euclidean distance matrix of our samples using the normalized ASV counts

Clustering Method

- We perform hierarchical clustering, using Ward.D2, to group similar samples based on the Euclidean distances

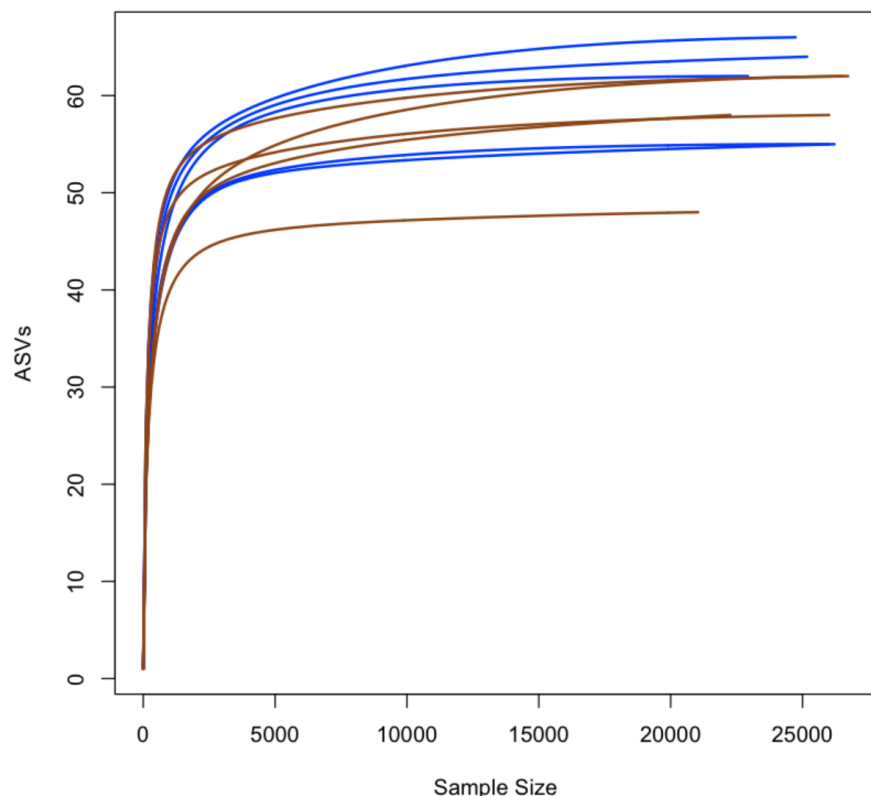


Processing AmpSeq Data: Data Visualization – Alpha Diversity

- Alpha diversity involves summary metrics that describe individual samples
 - Richness – the total number of distinct units in a sample (number of unique ASVs)
 - Evenness – how close in number each distinct unit is observed in a sample
 - Diversity – a combination of richness and evenness

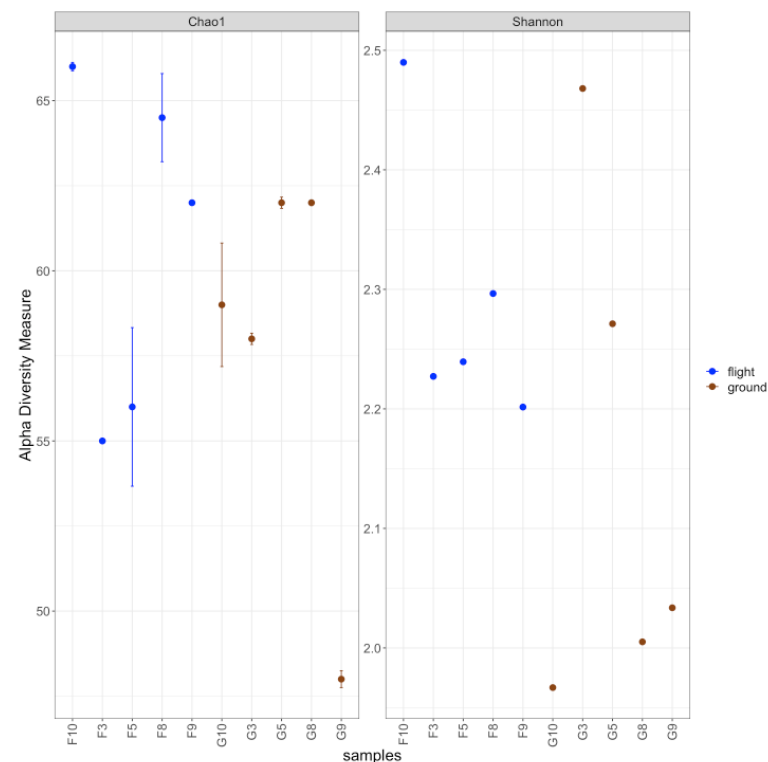
Rarefaction Curves

- A visual representation of the diversity that exists in samples, we use the `rarecurve()` function



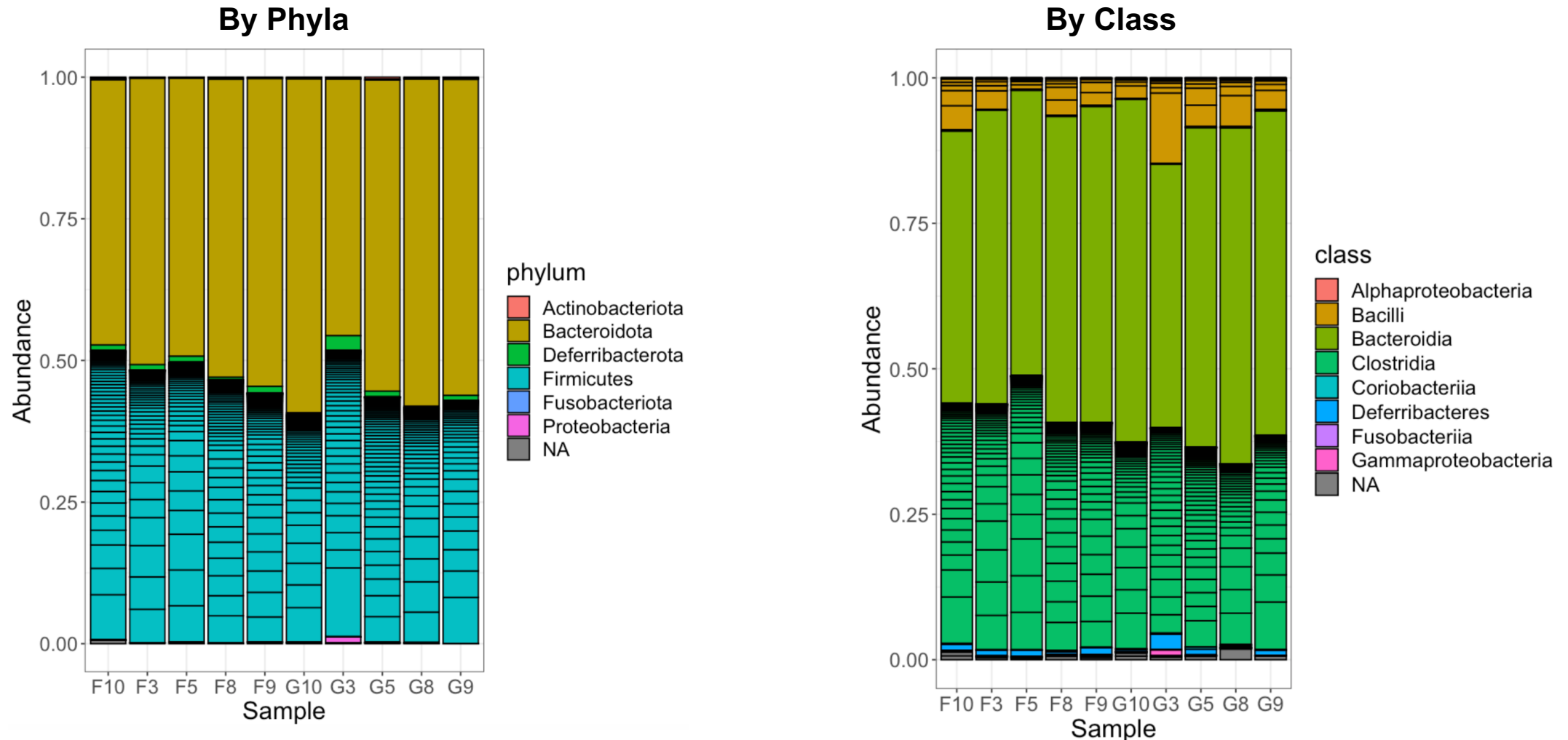
Richness and Diversity Estimates

- *Chao1* – A richness estimator
- *Shannon's diversity index* – A metric of diversity that includes richness and evenness
- We use the `plot_richness()` function of `phyloseq`



Processing AmpSeq Data: Data Visualization – Taxonomy

- Taxonomic figures can be used as a visual summary of the data
- These can be generated by normalizing all counts as a proportion of the total number of counts per sample, using the `transform_sample_counts()` function of the `phyloseq` package, then plotting those data with `plot_bar()`



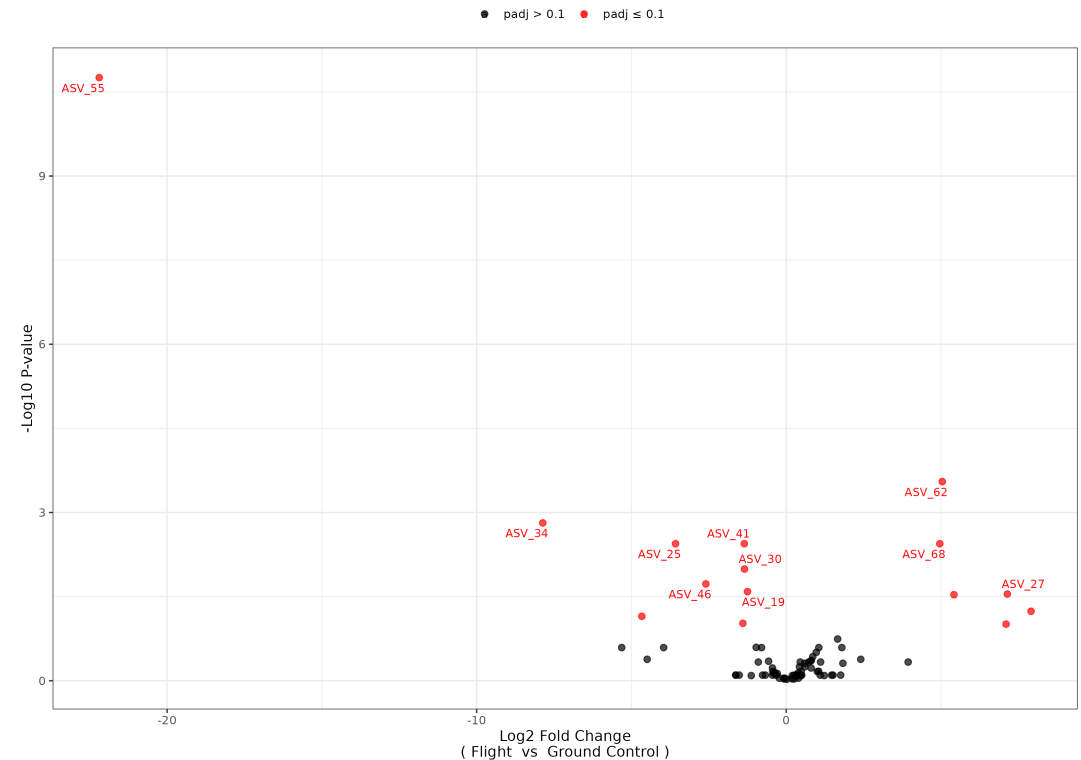
Processing AmpSeq Data: Data Visualization – Differential Abundance Analysis

- Differential abundance analysis (DAA) is used to test for differences in abundances of sequences (e.g. ASVs) or taxa recovered between groups
- We use DESeq2 to perform DAA, which runs 3 consecutive steps:
 - normalizing for sample read-depth and composition (using the median of ratios method to estimate size factors)
 - transforming the data (estimates dispersion and performs data normalization)
 - testing for differential abundance between groups (using the negative binomial GLM fitting and Wald statistics by default)

DAA Table

	baseMean	log2FoldChange	lfcSE	stat	pvalue	padj
	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
ASV_3	860.738994	-1.399241	0.5890130	-2.375569	1.752192e-02	9.932675e-02
ASV_19	150.237667	-1.249008	0.4272325	-2.923485	3.461365e-03	2.769092e-02
ASV_25	114.485574	-3.569227	0.9612061	-3.713280	2.045903e-04	4.091806e-03
ASV_27	97.950896	7.150070	2.4068052	2.970772	2.970522e-03	2.640464e-02
ASV_30	86.741651	-1.342332	0.4048856	-3.315337	9.153265e-04	1.046087e-02
ASV_33	71.259216	-7.870632	1.9236449	-4.091520	4.285545e-05	1.142812e-03
ASV_41	49.164451	-1.361315	0.3811661	-3.571450	3.550109e-04	4.733479e-03
ASV_46	43.039718	-2.612064	0.8335157	-3.133791	1.725639e-03	1.725639e-02
ASV_52	30.142809	-4.543276	1.8465438	-2.460421	1.387740e-02	8.539939e-02
ASV_54	26.325446	-22.201389	3.0274634	-7.333330	2.245029e-13	1.796023e-11

Volcano Plot



GeneLab Amplicon (Illumina) Sequencing Data Processing Pipeline

GeneLab's Amplicon Seq pipeline details are available on GitHub:

https://github.com/nasa/GeneLab_Data_Processing/blob/master/Amplicon/Illumina/Pipeline_GL-DPPD-7104_Versions/GL-DPPD-7104-A.md

Files

GeneLab_Data_Processing / Amplicon / Illumina / Pipeline_GL-DPPD-7104_Versions / GL-DPPD-7104-A.md

asaravia-butler Updating links from GLDS to OSDR and making files published on OSDR bold · 6637463 · 3 weeks ago · History

Preview Code Blame 560 Lines (321 loc) · 19.3 KB

Bioinformatics pipeline for amplicon Illumina sequencing data

This page holds an overview and instructions for how GeneLab processes Illumina amplicon datasets. Exact processing commands for specific datasets that have been released are available in the [GLDS_Processing_Scripts](#) sub-directory and/or are provided with their processed data in the [Open Science Data Repository \(OSDR\)](#).

Table of contents

- Software used
- Reference databases used
- General processing overview with example commands
 - 1. Raw Data QC
 - Compile Raw Data QC
 - 2. Trim Primers
 - 3. Quality filtering
 - 4. Filtered Data QC
 - Compile Filtered Data QC
 - 5. Calculate error model, apply DADA2 algorithm, assign taxonomy, and create output tables
 - Learning the error rates
 - Inferring sequences
 - Merging forward and reverse reads
 - Generating sequence table with counts per sample
 - Removing putative chimeras
 - Assigning taxonomy
 - Generating and writing standard outputs

Software used

Program	Version*	Relevant Links
FastQC	fastqc -v	https://www.bioinformatics.babraham.ac.uk/projects/fastqc/
MultiQC	multiqc -v	https://multiqc.info/
Cutadapt	cutadapt --version	https://cutadapt.readthedocs.io/en/stable/
DADA2	packageVersion("dada2")	https://www.bioconductor.org/packages/release/bioc/html/dada2.html
DECIPHER	packageVersion("DECIPHER")	https://bioconductor.org/packages/release/bioc/html/DECIPHER.html
biomformat	packageVersion("biomformat")	https://github.com/joey711/biomformat

* Exact versions are available along with the processing commands for each specific dataset.

Reference databases used

Program used	Database	Relevant Links
DECIPHER	SILVA SSU r138	http://www2.decipher.codes/Classification/TrainingSets/SILVA_SSU_r138_2019.RData
DECIPHER	UNITE v2020	http://www2.decipher.codes/Classification/TrainingSets/UNITE_v2020_February2020.RData

2. Trim Primers

The location and orientation of primers in the data is important to understand in deciding how to do this step. `cutadapt` has many options for primer identification and removal. They are described in detail on their documentation page here:

<https://cutadapt.readthedocs.io/en/stable/guide.html#adapter-types>

The following example commands show how it was done for some samples of [GLDS-200](#), which was 2x250 sequencing of the 16S gene using these primers:

- forward: 5'-GTGCCAGCMGCCGCGGTAA-3'
- reverse: 5'-GGACTACVSGGGTATCTAAT-3'

Due to the size of the target amplicon and the type of sequencing done here, both forward and reverse primers are expected to be on each of the forward and reverse reads. It therefore takes "linked" primers as input for forward and reverse reads, specified above by the `...` between them. It also expects that the primers start at the first position of the reads ("anchored"), specified with the leading `^` characters.

The following website is useful for reverse complementing primers and dealing with degenerate bases appropriately:

<http://arep.med.harvard.edu/labgc/adnan/projects/Utilities/revcomp.html>

```
cutadapt -a ^GTGCCAGCMGCCGCGGTAA...ATTAGATACCCSBGTAGTCC -A ^GGACTACVSGGGTATCTAAT...TTACCGGGCKGCTGGCAC \  
## Define what B represents; and define what K represents ##  
-o Primer-trimmed-R1.fq.gz -p Primer-trimmed-R2.fq.gz Input_R1_raw.fastq.gz Input_R2_raw.fastq.gz \  
--discard-untrimmed
```

Parameter Definitions:

- `-a` – specifies the primers and orientations expected on the forward reads (when primers are linked as noted above)
- `-A` – specifies the primers and orientations expected on the reverse reads (when primers are linked as noted above)
- `-o` – specifies output of forward, primer-trimmed reads
- `-p` – specifies output of reverse, primer-trimmed reads
- `Input_R1_raw.fastq.gz` – this and following "R2" version are positional arguments specifying the forward and reverse reads, respectively, for input
- `--discard-untrimmed` – this filters out those reads? where the primers were not found as expected

Input Data:

- `fastq`, compressed or uncompressed (original reads)

Output Data:

- `trimmed.fastq.gz`, compressed or uncompressed (trimmed reads)
- `trimmed-read-counts.tsv` (per sample read counts before and after trimming)
- `cutadapt.log` (log file of standard output and error from cutadapt)

GeneLab Amplicon (Illumina) Sequencing Data Processing Workflow

GeneLab's Amplicon Seq workflow details are available on GitHub:

https://github.com/nasa/GeneLab_Data_Processing/blob/amplicon-add-runsheet-visualizations/Amplicon/Illumina/Workflow_Documentation/SW_AmpIllumina-A

The screenshot shows a GitHub repository page for the file `SW_AmpIllumina-A/README.md`. The left sidebar displays the repository's file structure, with `SW_AmpIllumina-A/README.md` selected. The main content area shows the README text, which includes a title, general workflow information, and a list of steps for utilizing the workflow.

SW_AmpIllumina-A Workflow Information and Usage Instructions

General workflow info

The current GeneLab Illumina amplicon sequencing data processing pipeline (AmpIllumina), [GL-DPPD-7104-A.md](#), is implemented as a [Snakemake](#) workflow and utilizes [conda](#) environments to install/run all tools. This workflow (SW_AmpIllumina-A) is run using the command line interface (CLI) of any unix-based system. The workflow can be used even if you are unfamiliar with Snakemake and conda, but if you want to learn more about those, [this Snakemake tutorial](#) within [Snakemake's documentation](#) is a good place to start for that, and an introduction to conda with installation help and links to other resources can be found [here at Happy Belly Bioinformatics](#).

Utilizing the workflow

1. [Install conda, mamba, and `geneLab-utils` package](#)
2. [Download the workflow template files](#)
3. [Run the workflow using `run_workflow.py`](#)
 - 3a. [Approach 1: Run the workflow on a GeneLab Amplicon \(Illumina\) sequencing dataset with automatic retrieval of raw read files and metadata](#)
 - 3b. [Approach 2: Run the workflow on a non-OSD dataset using a user-created runsheet](#)
4. [Parameter Definitions](#)
5. [Additional output files](#)

1. Install conda, mamba, `geneLab-utils`, and `dp-tools` package

We recommend installing a Miniconda, Python3 version appropriate for your system, as exemplified in [the above link](#).

Once conda is installed on your system, we recommend installing [mamba](#), as it generally allows for much faster conda installations:

```
conda install -n base -c conda-forge mamba
```

You can read a quick intro to mamba [here](#) if wanted.

NASA EDGE GeneLab AmpSeq Workflow

<https://nasa.edgebioinformatics.org/home>

My Projects My uploads Job Queue AS

Home
Public Projects
Upload Files
Run Workflow

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 (Amplllumina) integrated.

Amplicon Processing Overview

Raw data fastQC (fastqc/multiqc) → Primer removal (cutadapt) → Quality filter/trim (dada2, bbdduk) → Filtered data fastQC (fastqc/multiqc)

Legend:
■ All
● Illumina data
▲ 454/Ion-Torrent data

Flowchart:
Filtered data fastQC → Resolve ASVs (dada2) → Differential abundance (DESeq2)
Filtered data fastQC → Generate OTUs (vsearch) → Assign taxonomy (DECIPHER, SILVA DB for 16S, UNITE DB for ITS, PR2 DB for 18S)

Primary outputs:

- Recovered ASV/OTU sequences (fasta)
- Counts of sequences per sample
- Taxonomic classifications of sequences
- Beta and alpha diversity plots
- Taxonomic bar plots
- Differential abundance analysis and visualizations

Los Alamos NATIONAL LABORATORY
Managed by Triad National Security, LLC for the U.S Dept. of Energy's NNSA
© Copyright Triad National Security, LLC. All Rights Reserved.

NASA ACCESS Advancing Innovation

NASA EDGE GeneLab AmpSeq Workflow

<https://nasa.edgebioinformatics.org/ampillumina>

GeneLab EDGE AS

My Projects My uploads Job Queue

Home Public Projects Upload Files Run Workflow

Ampillumina: GeneLab Illumina amplicon sequencing data processing pipeline

Run Workflow

Project/Run Name (required, at 3 but less than 30 characters)

Description (optional)

Input

- Run the workflow on a GeneLab Amplicon (Illumina) sequencing dataset with automatic retrieval of raw read files and metadata**
This approach processes data hosted on the [NASA Open Science Data Repository \(OSDR\)](#). Upon execution, the command downloads then parses the OSD ISA.zip file to create a runsheet containing link(s) to the raw reads and the metadata required for processing. The runsheet is then used to prepare the necessary configuration files before executing the workflow using the specified Snakemake run command.
- Run the workflow on a non-OSD dataset using a user-created runsheet**
If processing a non-OSD dataset, you must manually create the runsheet for your dataset to run the workflow. Specifications for creating a runsheet manually are described [here](#).

OSD Id OSD-
⚠ An OSD Id is required. Acceptable format: OSD-###

Target Region 16S

Trim Primers TRUE FALSE

Primers Linked TRUE FALSE

Anchor Primers TRUE FALSE

Discard Untrimmed TRUE FALSE

Left Trunc 0

Right Trunc 0

Left maxEE 1

Right maxEE 1

Minimum Cutadapt Length 130

Concatenate Reads Only TRUE FALSE

Output Prefix (optional)

Specify Runsheets (optinal, select a file...)

Submit

GeneLab EDGE AS

My Projects My uploads Job Queue

Home Public Projects Upload Files Run Workflow

Ampillumina: GeneLab Illumina amplicon sequencing data processing pipeline

Run Workflow

Project/Run Name (required, at 3 but less than 30 characters)

Description (optional)

Input

- Run the workflow on a GeneLab Amplicon (Illumina) sequencing dataset with automatic retrieval of raw read files and metadata**
This approach processes data hosted on the [NASA Open Science Data Repository \(OSDR\)](#). Upon execution, the command downloads then parses the OSD ISA.zip file to create a runsheet containing link(s) to the raw reads and the metadata required for processing. The runsheet is then used to prepare the necessary configuration files before executing the workflow using the specified Snakemake run command.
- Run the workflow on a non-OSD dataset using a user-created runsheet**
If processing a non-OSD dataset, you must manually create the runsheet for your dataset to run the workflow. Specifications for creating a runsheet manually are described [here](#).

Runsheet (required, select a file...)

Target Region 16S

Trim Primers TRUE FALSE

Primers Linked TRUE FALSE

Anchor Primers TRUE FALSE

Discard Untrimmed TRUE FALSE

Left Trunc 0

Right Trunc 0

Left maxEE 1

Right maxEE 1

Minimum Cutadapt Length 130

Concatenate Reads Only TRUE FALSE

Output Prefix (optional)

Submit

Acknowledgements

- **NASA OSDR**

- PM: Sylvain Costes, PhD
- PS: Lauren Sanders, PhD
- OSDR Team

- **NASA GeneLab**

- PM: Sylvain Costes, PhD
- DPM: Samrawit Gebre
- PS/DP: Amanda Saravia-Butler, PhD
- DP: Mike Lee, PhD
- GeneLab Team

- **NASA ALSDA**

- PM: Sylvain Costes, PhD
- DPM: Danielle Lopez
- PS: Ryan Scott
- ALSDA Team

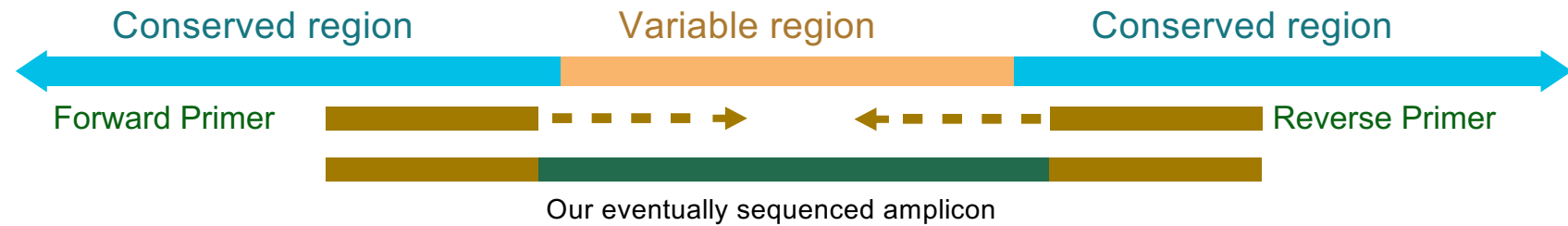
Funding

- GeneLab is funded by the NASA Space Biology program within the NASA Science Mission Directorate's (SMD) Biological and Physical Sciences (BPS) Division

**GeneLab Amplicon (Illumina)
Sequencing
Data Processing Pipeline
EXTRA SLIDES**

Why It's Important To Remove Primers

Remember from our sketch earlier that the part of the DNA our primers bind to is included in our final amplicon that we sequence



Why It's Important To Remove Primers

Primers often have “degenerate” bases, which are positions that can be multiple bases

Forward primer: GTG^YCAGC^MGCCGCGGTAA
 Reverse primer: GGACTAC^NVGGGT^WTCTAAT

Y = C or T

M = A or C

N = A, T, G, or C

V = A, C, or G

W = A or T

These varied primers will bind to a biological sequence with some frequency even if they are not identical

A primer written like this:

GTG^YCAGC^MGCCGCGGTAA

Actually means these are the sequences used during PCR amplification:

GTG^CCAGC^AGCCGCGGTAA

GTG^CCAGC^CGCCGCGGTAA

GTG^TCAGC^AGCCGCGGTAA

GTG^TCAGC^CGCCGCGGTAA

GTG^CCAGCAGCCGCGGTAA...

Biological sequence:

<the forward primer binds to the complement>

CAC^GGGTCGTCGGCGCCATT...

Primer:

GTG^TCAGCCGCCGCGGTAA

Now when this sequence is built, extending from the primer here, this molecule will have a 'T' in the 4th position, even though the biological source had a 'C'

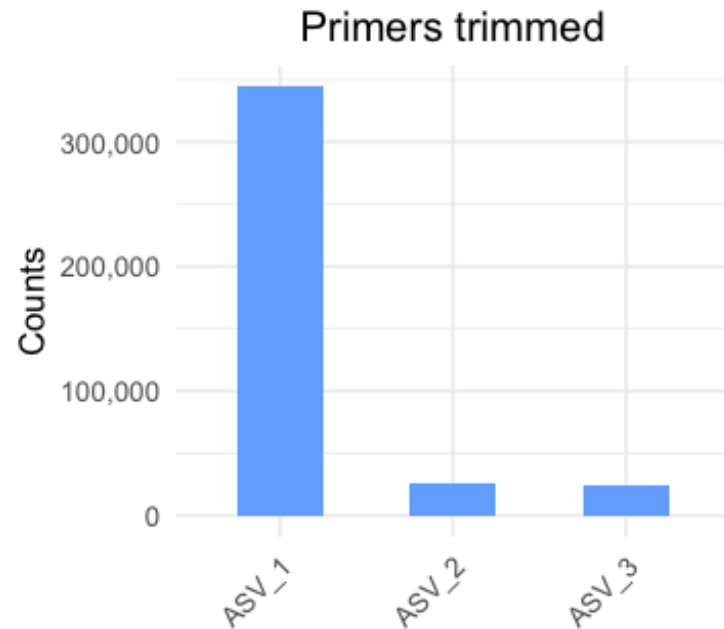
That sequence will get amplified further and show up in our data – presenting the true biological sequence between the primers, but with technically introduced variations in the primer regions.

This is why it is critical we remove the primers during processing, so that we are only working with true biological sequences.

Why It's Important To Remove Primers

Example with real data

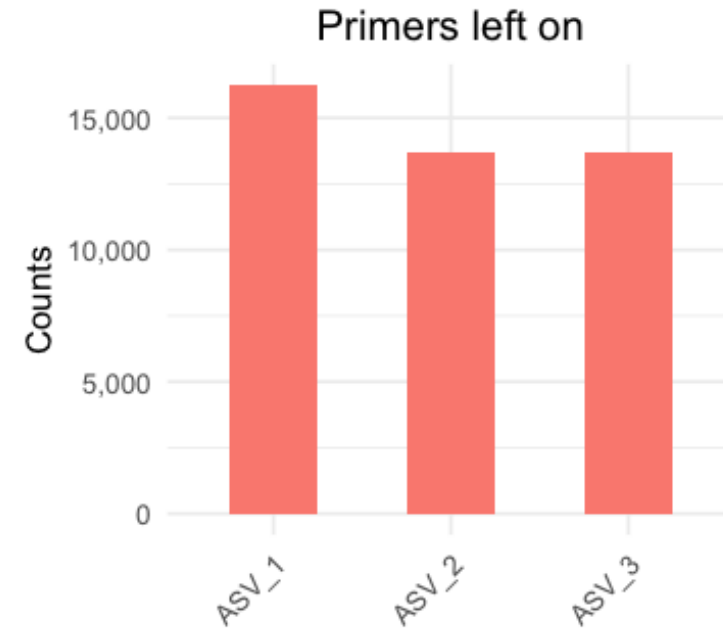
When trimming the primers off at the start of processing, the most abundant ASV (Amplicon Sequence Variant) recovered was sequenced over 300,000 times.



Assigned taxonomy

ASV_1: phylum Bacteroidota; genus *Parabacteroides*
ASV_2: phylum Firmicutes; genus *Erysipelatoclostridium*
ASV_3: phylum Firmicutes; genus *Ruminococcus*

When processed *without* trimming the primers, the most abundant ASV recovered was sequenced just over 15,000 times.

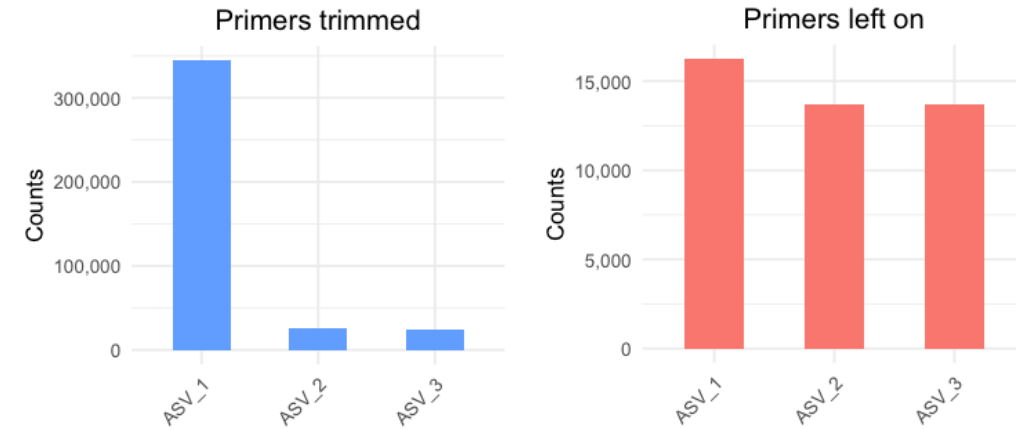


ASV_1: phylum Bacteroidota; genus *Parabacteroides*
ASV_2: phylum Bacteroidota; genus *Parabacteroides*
ASV_3: phylum Bacteroidota; genus *Parabacteroides*

***Parabacteroides* definitely dominates the sample in either case, but nothing else shows up in the top 3 when we left the primers on.**

Why It's Important To Remove Primers

We can look at these sequences and see exactly what is going on.



This is the super-abundant ASV_1 from when the primers were trimmed:

```
>ASV_1
TACGGAGGATGCGAGCGTTATCCGGATTTATTGGGTTTAAAGGGTGCGTAGGTGGTTAATTAAGTCAGCGGTGAAAGTTTGTGGCTCAACCATAAAAATTGCCGTTGAAACTGGTTGACTTGAGTATATTTGAGGTAGGCGGAATGCGTG
GTGTAGCGGTGAAATGCATAGATATCACGCAGAACTCCGATTGCGAAGGCAGCTTACTAACTATAACTGACACTGAAGCACGAAAGCGTGGGGATCAAACAGG
```

These are the sequences for the first 3 ASVs from when the primers were left on:

```
>ASV_1
GTGCCAGCAGCCGCGGTAATACGGAGGATGCGAGCGTTATCCGGATTTATTGGGTTTAAAGGGTGCGTAGGTGGTTAATTAAGTCAGCGGTGAAAGTTTGTGGCTCAACCATAAAAATTGCCGTTGAAACTGGTTGACTTGAGTATATTT
GAGGTAGGCGGAATGCGTGGTGTAGCGGTGAAATGCATAGATATCACGCAGAACTCCGATTGCGAAGGCAGCTTACTAACTATAACTGACACTGAAGCACGAAAGCGTGGGGATCAAACAGGATTAGAAACCCTAGTAGTCC
```

```
>ASV_2
GTGCCAGCAGCCGCGGTAATACGGAGGATGCGAGCGTTATCCGGATTTATTGGGTTTAAAGGGTGCGTAGGTGGTTAATTAAGTCAGCGGTGAAAGTTTGTGGCTCAACCATAAAAATTGCCGTTGAAACTGGTTGACTTGAGTATATTT
GAGGTAGGCGGAATGCGTGGTGTAGCGGTGAAATGCATAGATATCACGCAGAACTCCGATTGCGAAGGCAGCTTACTAACTATAACTGACACTGAAGCACGAAAGCGTGGGGATCAAACAGGATTAGAAACCCTGGTAGTCC
```

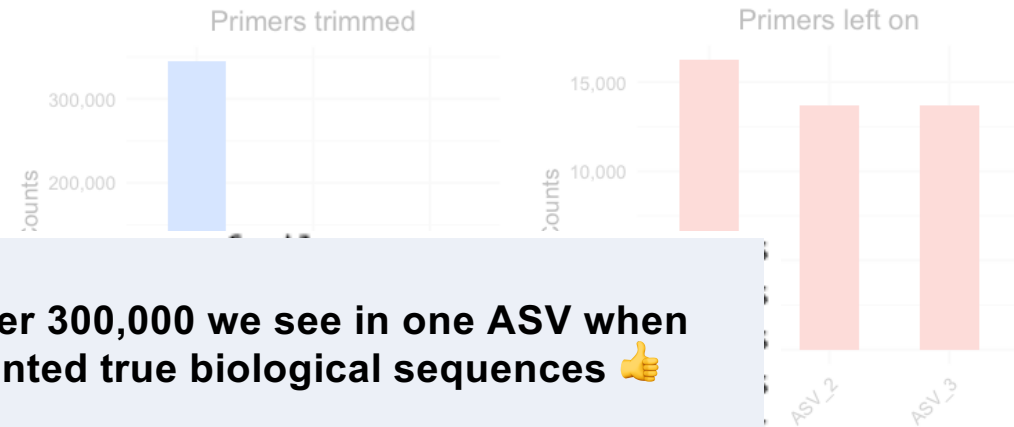
```
>ASV_3
GTGCCAGCCGCCGCGGTAATACGGAGGATGCGAGCGTTATCCGGATTTATTGGGTTTAAAGGGTGCGTAGGTGGTTAATTAAGTCAGCGGTGAAAGTTTGTGGCTCAACCATAAAAATTGCCGTTGAAACTGGTTGACTTGAGTATATTT
GAGGTAGGCGGAATGCGTGGTGTAGCGGTGAAATGCATAGATATCACGCAGAACTCCGATTGCGAAGGCAGCTTACTAACTATAACTGACACTGAAGCACGAAAGCGTGGGGATCAAACAGGATTAGAAACCCTAGTAGTCC
```

The highlighted portion for all 3 is identical to the above sequence recovered when the primers were trimmed.

The only difference in each is 1 base that was degenerate in the primers. All of these variants are technical, not biological.

Why It's Important To Remove Primers

In fact, when we failed to remove the primers, the top 36 most abundant ASVs were all the same biological sequence, just with single-nucleotide variations introduced by the degenerate primers.



This is the sequence:

```

>ASV_1 TACGGAGGATGCGGTAATACGGAGGATGCGAGCGTTATCCGGATTTATTGGGTTTAAAGGGTGCGTAGGTGGTTAATTAAGTCAGCGGTGAAAGTTTGTGGCTCAACCATAAAATTGCCGTTGAAACTGGTTGACTTGAGTATATTT
GAGGTAGGCGGAATGCGTGGTGTAGCGGTGAAATGCATAGATATCACGCAGAACTCCGATTGCGAAGGCAGCTTACTAACTATAACTGACACTGAAGCACGAAAGCGTGGGGATCAAACAGGATTAGAAACCTAGTAGTCC
ASV_2 GTGCCAGCAGCCGCGGTAATACGGAGGATGCGAGCGTTATCCGGATTTATTGGGTTTAAAGGGTGCGTAGGTGGTTAATTAAGTCAGCGGTGAAAGTTTGTGGCTCAACCATAAAATTGCCGTTGAAACTGGTTGACTTGAGTATATTT
GAGGTAGGCGGAATGCGTGGTGTAGCGGTGAAATGCATAGATATCACGCAGAACTCCGATTGCGAAGGCAGCTTACTAACTATAACTGACACTGAAGCACGAAAGCGTGGGGATCAAACAGGATTAGAAACCTAGTAGTCC
ASV_3 GTGCCAGCAGCCGCGGTAATACGGAGGATGCGAGCGTTATCCGGATTTATTGGGTTTAAAGGGTGCGTAGGTGGTTAATTAAGTCAGCGGTGAAAGTTTGTGGCTCAACCATAAAATTGCCGTTGAAACTGGTTGACTTGAGTATATTT
GAGGTAGGCGGAATGCGTGGTGTAGCGGTGAAATGCATAGATATCACGCAGAACTCCGATTGCGAAGGCAGCTTACTAACTATAACTGACACTGAAGCACGAAAGCGTGGGGATCAAACAGGATTAGAAACCTAGTAGTCC
  
```

If we sum these 36 variants, we get back to the over 300,000 we see in one ASV when we removed the primers, where our data represented true biological sequences 👍

If we didn't properly remove the primers, it would appear as though there were dozens more unique variants for *Parabacteroides* than there actually were.

This is why it is critical we remove the primers during processing, so that we are only working with true biological sequences.

The highlighted portion for all 3 is identical to the above sequence recovered when the primers were trimmed.

The only difference in each is 1 base that was degenerate in the primers. All of these variants are technical, not biological.