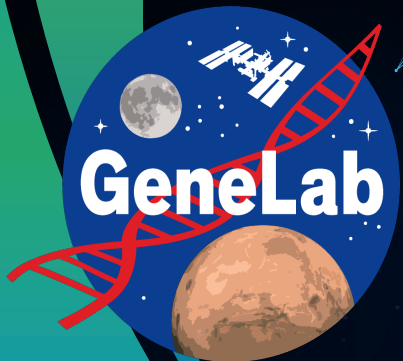
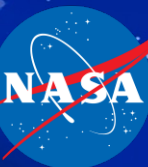


Developing Open-Source Training Materials for AI/ML and Space Biological Sciences using NASA Cloud-Based Data

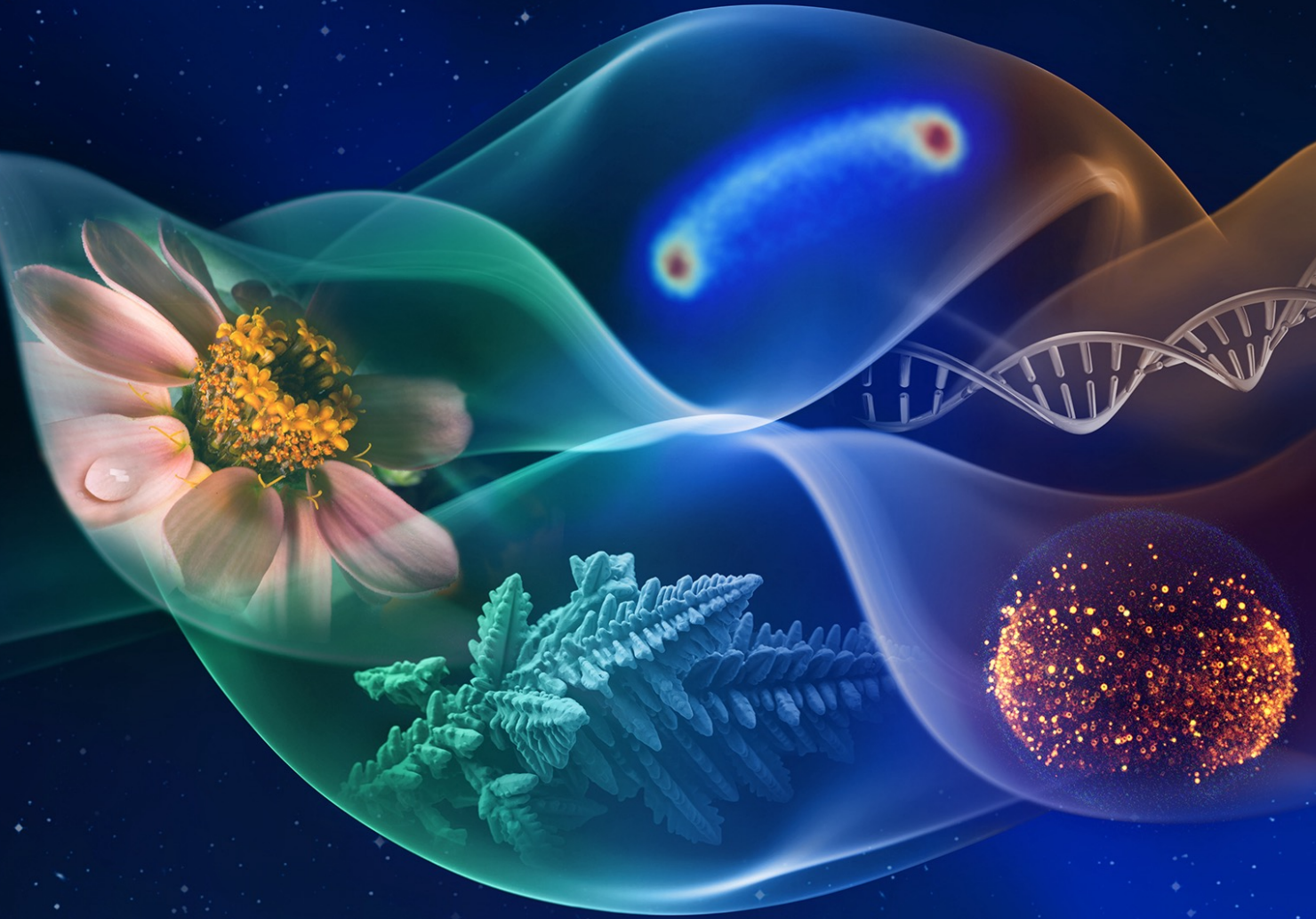


James Casaletto, Ph.D.
Staff Scientist
Blue Marble Space Institute of Science
NASA Ames Research Center

National Aeronautics and
Space Administration



What is the
motivation for
developing this
training?

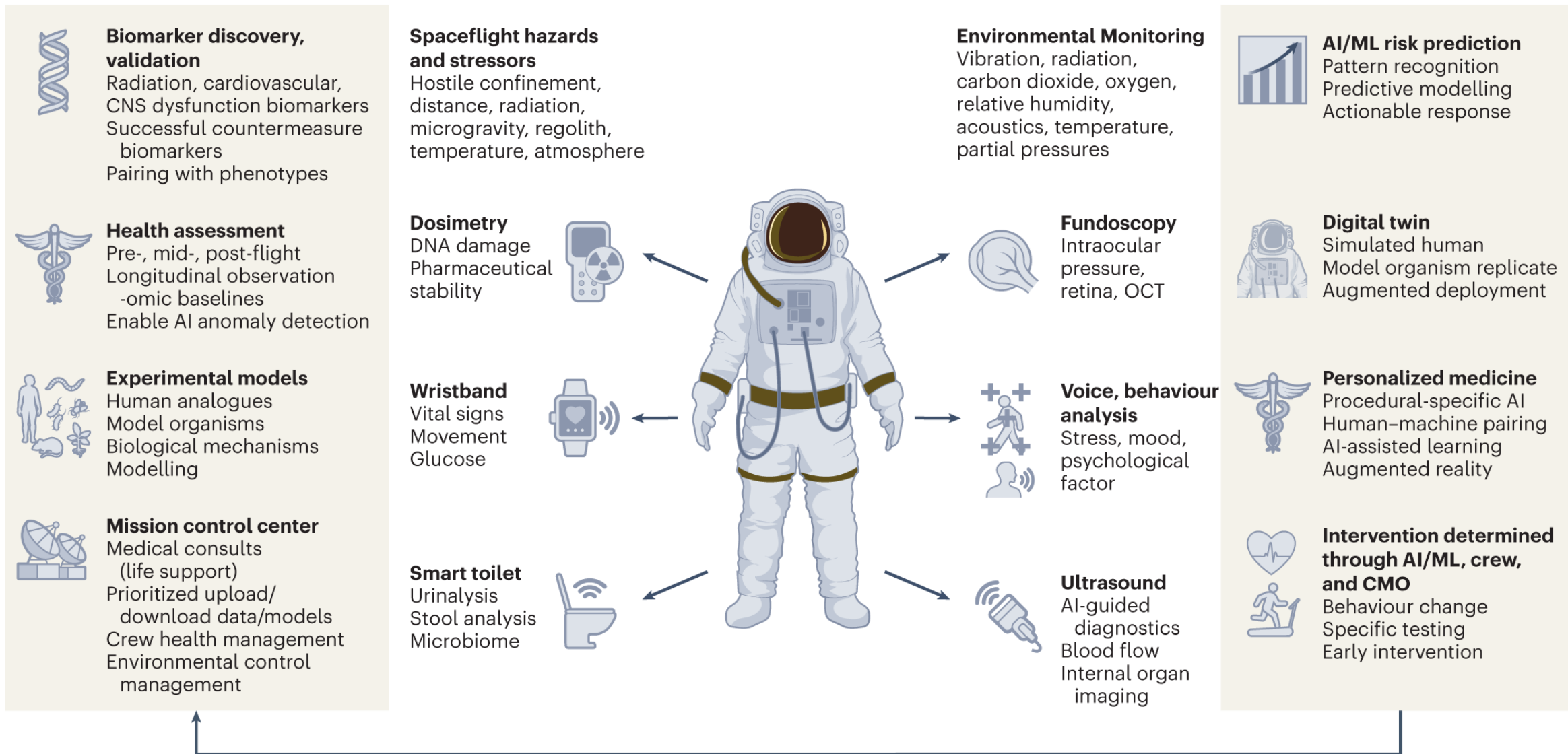


The moonshot: use AI/ML to help keep astronauts healthy

Research and terrestrial support

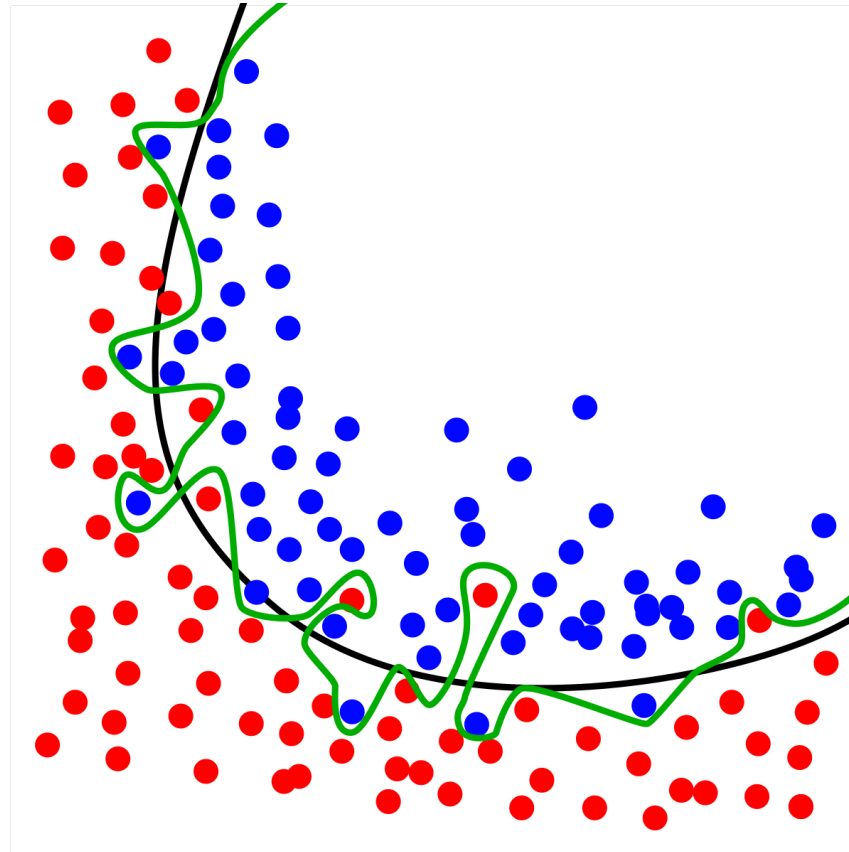
Real-time support and operations

Countermeasures



Space biology is under-understood in part for lack of data

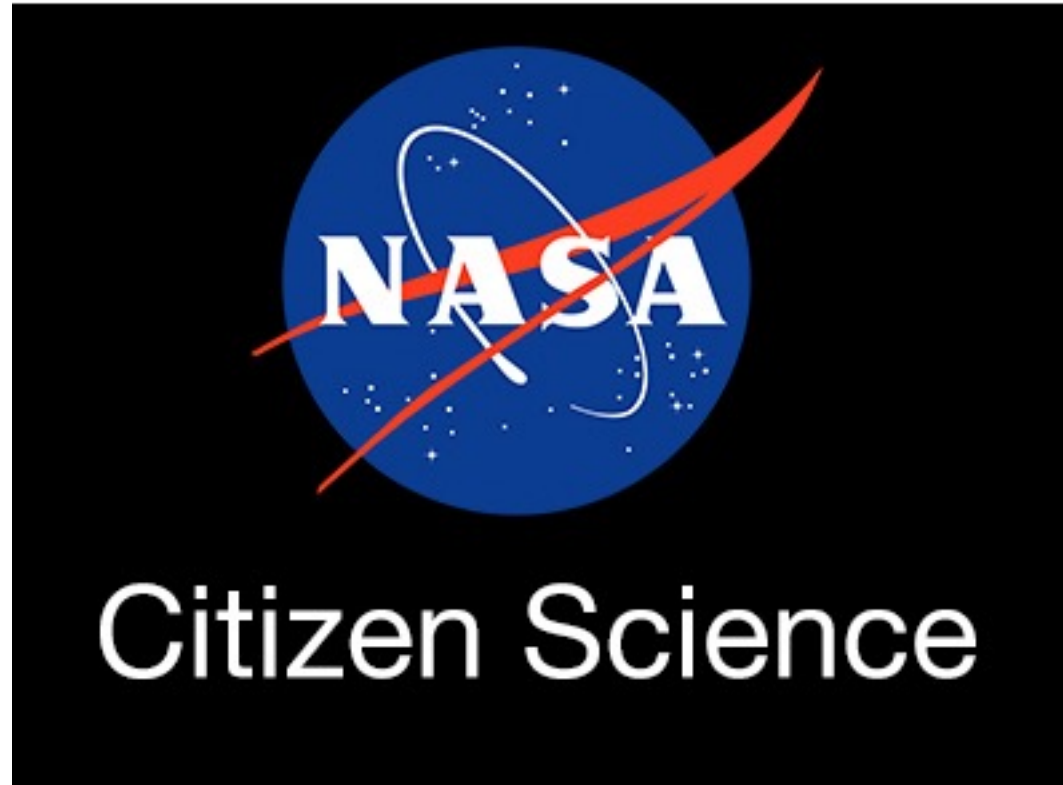
Molecular biology is difficult to study for sparse data



AI/ML is a promising approach to study sparse data

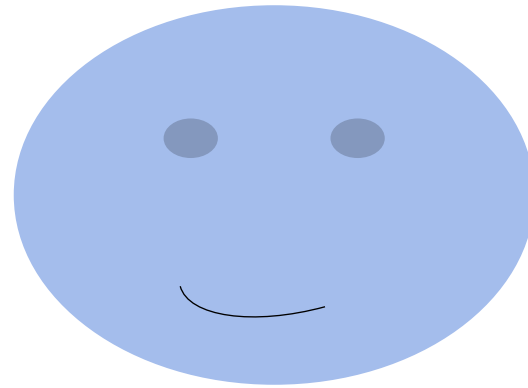
ML Model	Description
Support vector machine	SVM uses support vectors to arbitrarily partition the feature space.
Logistic regression	LR uses regularization to prevent overfitting.
K-nearest neighbors	KNN computes distances between data points.
Multi-layer perceptron	MLP uses gradient descent for optimization.
Decision tree & random forest	DT and RF can work well when coupled with gradient boosting.

Calling all citizen scientists!



AI/ML, while ubiquitous, is poorly understood by many

Dude –
how hard
can it be?



There are LOTS of options to learn machine learning

The collage features several educational options:

- Coursera:** "Applied Data Science Program: Leveraging AI for Effective Decision-Making" (12 weeks, MIT faculty) and "Machine Learning" (focus on creating systems for language processing and statistical pattern recognition).
- DataCamp:** "Understanding Machine Learning: An introduction to machine learning with no coding involved."
- Udacity:** "Intro to Machine Learning" (free course, part of the School of Artificial Intelligence).
- YouTube:** A list of 9 videos from StatQuest with Josh Starmer, including "A Gentle Introduction to Machine Learning", "Machine Learning Fundamentals: Cross Validation", "Machine Learning Fundamentals: The Confusion Matrix", "Machine Learning Fundamentals: Sensitivity and Specificity", "The Sensitivity, Specificity, Precision, Recall Sing-a-Long!!!", "Machine Learning Fundamentals: Bias and Variance", "ROC and AUC, Clearly Explained!", "ROC and AUC in R", and "Entropy (for data science) Clearly Explained!!!".
- Other:** "GPT-4 on Khan Academy" featuring a portrait of a man, and "Machine Learning Crash Course with TensorFlow APIs" from Google.

Machine Learning Crash Course with TensorFlow APIs

Google's fast-paced, practical introduction to machine learning, featuring a series of lessons lectures, real-world case studies, and hands-on practice exercises.

There are NO options to learn AI/ML for space biology

Until now

<https://github.com/nasa/Transform-to-Open-Science/>

A NASA OPEN-SOURCE SCIENCE INITIATIVE:
TOPS: TRANSFORM TO OPEN SCIENCE

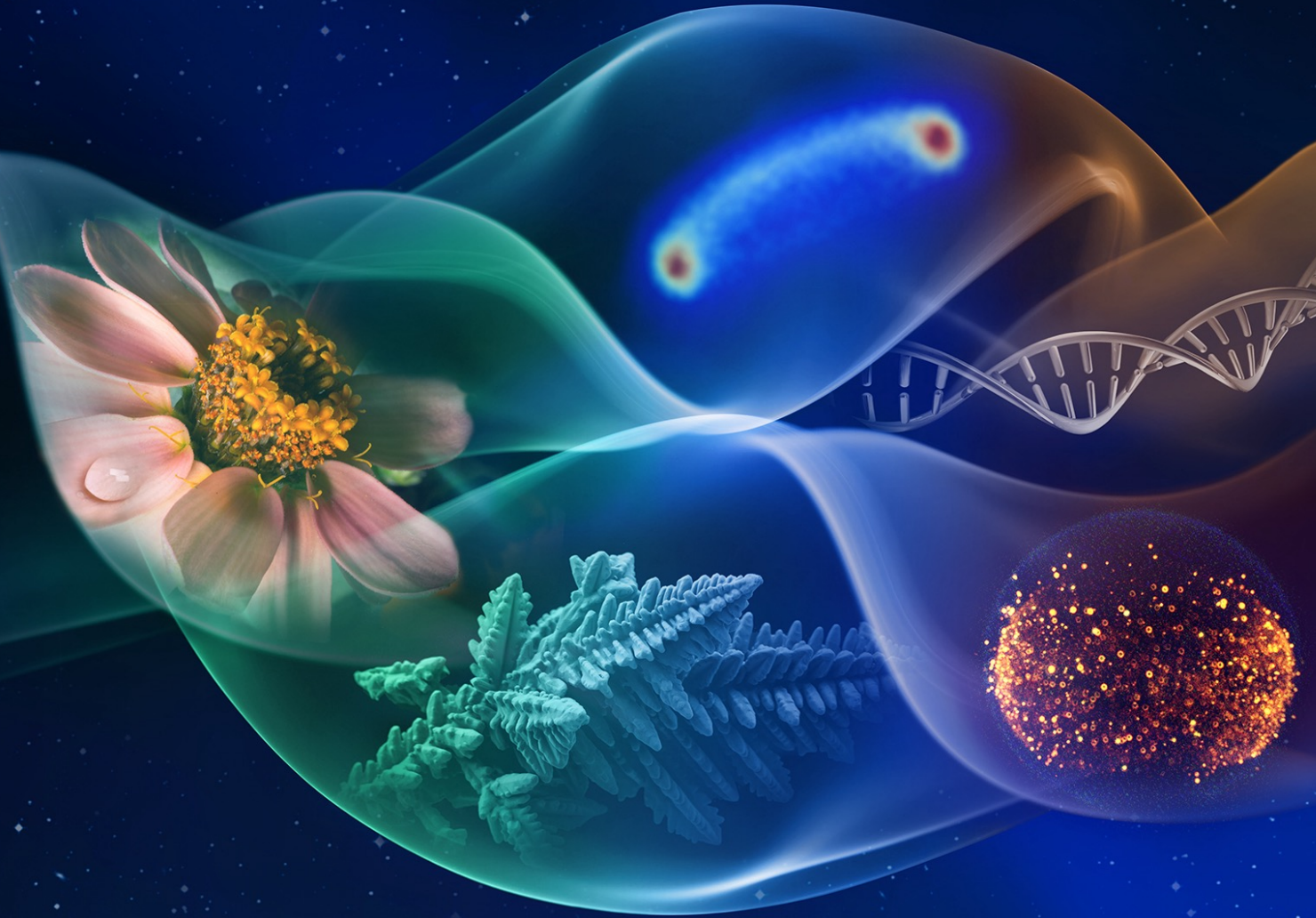
**F.14 Transform to Open Science Training:
A Science Mission Directorate
Cross Division Program**

Dr. Yaítza Luna-Cruz
Program Scientist | TOPS Program Officer

Dr. Chelle Gentemann
Program Scientist | TOPS Science Lead

Chief Science Data Office
NASA Headquarters 1

What are the choices for developing this training?



Choice 1: How to design the curriculum?

Options	Pros	cons
Design in a vacuum	<ul style="list-style-type: none">• simple• control deliverables	<ul style="list-style-type: none">• biased• lacks user perspective
Design by committee	<ul style="list-style-type: none">• fairly simple• less biased	<ul style="list-style-type: none">• lack of expertise• conflict resolution

Choice 2: How to develop the curriculum?

Options	Pros	cons
Build from scratch	<ul style="list-style-type: none">• control over product• consistent look/feel	<ul style="list-style-type: none">• time and effort• wheel re-invention
Curate collection	<ul style="list-style-type: none">• less time/effort• leverage what's good	<ul style="list-style-type: none">• inconsistent look/feel• no ML-for-space-bio

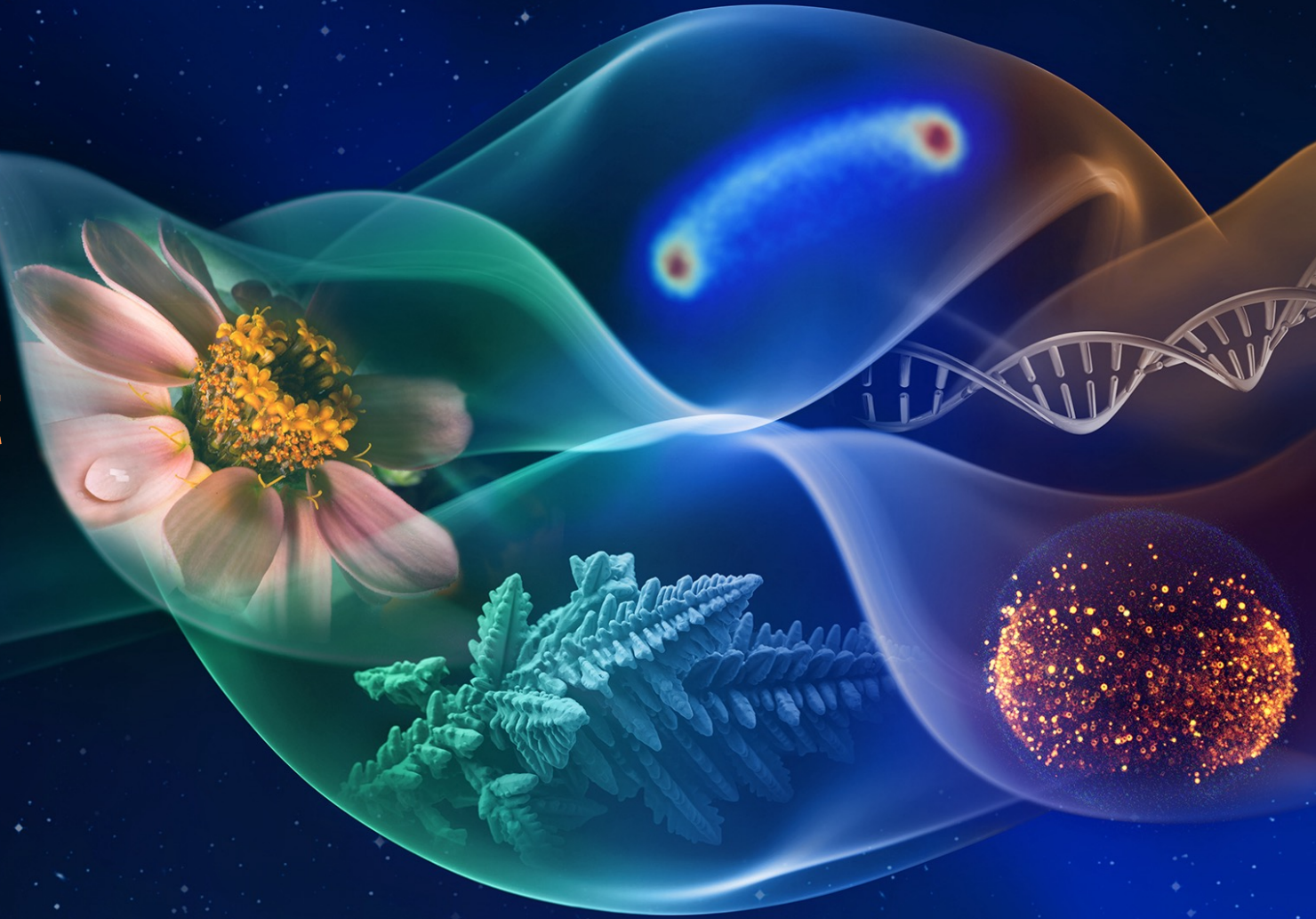
Choice 3: How to deploy the curriculum?

Options	Pros	cons
Instructor-led	<ul style="list-style-type: none">• high touch/interaction• dynamic	<ul style="list-style-type: none">• doesn't scale• small reach
Self-paced	<ul style="list-style-type: none">• scales• global reach	<ul style="list-style-type: none">• limited interaction• distractions/motivation

Choice 4: How to deploy the labs?

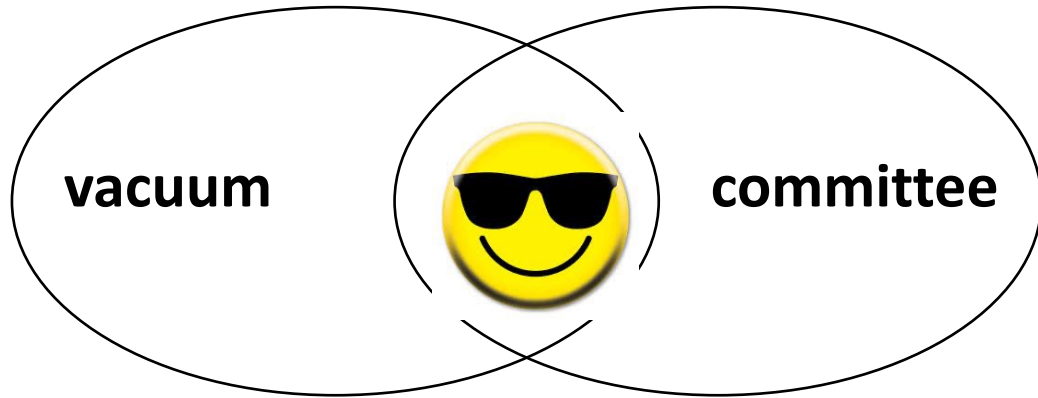
Options	Pros	cons
Student computer	<ul style="list-style-type: none">• zero cost• reproducible results	<ul style="list-style-type: none">• may exclude students• inconsistent experience
Cloud IaaS	<ul style="list-style-type: none">• scales well• consistent experience	<ul style="list-style-type: none">• cost• internet requirement

So what's the plan, at least to start?

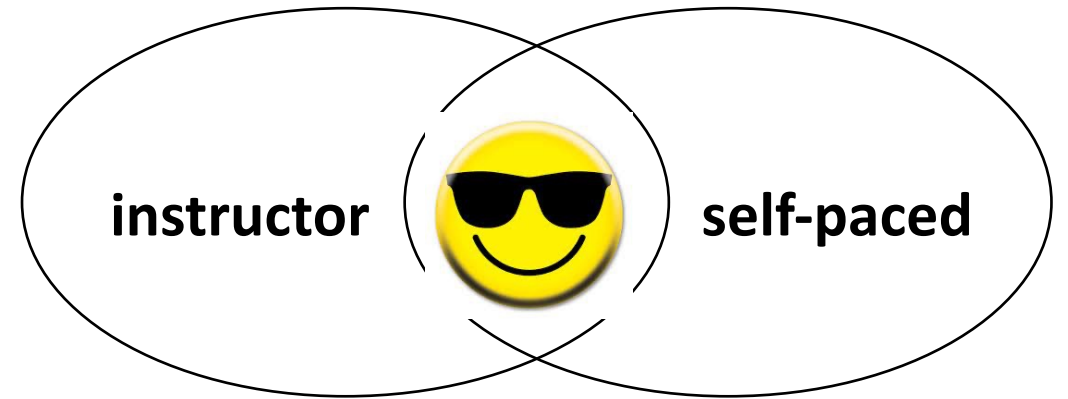


Our decision: meet in the middle

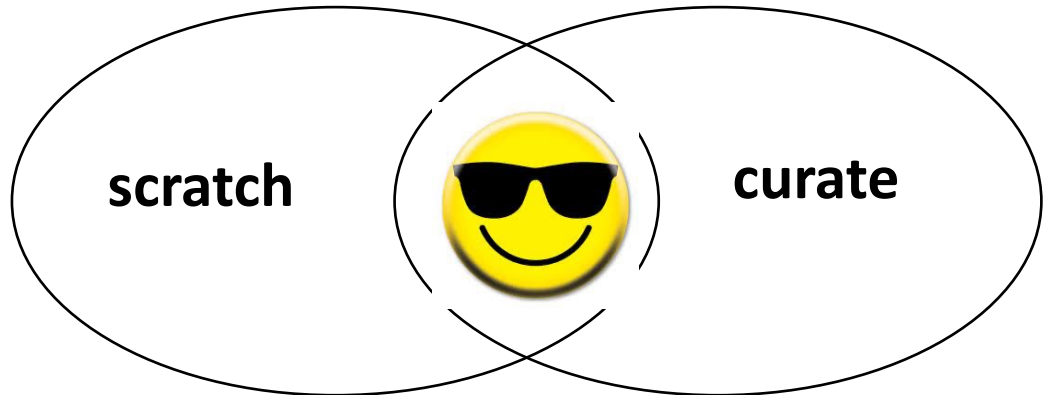
curriculum design



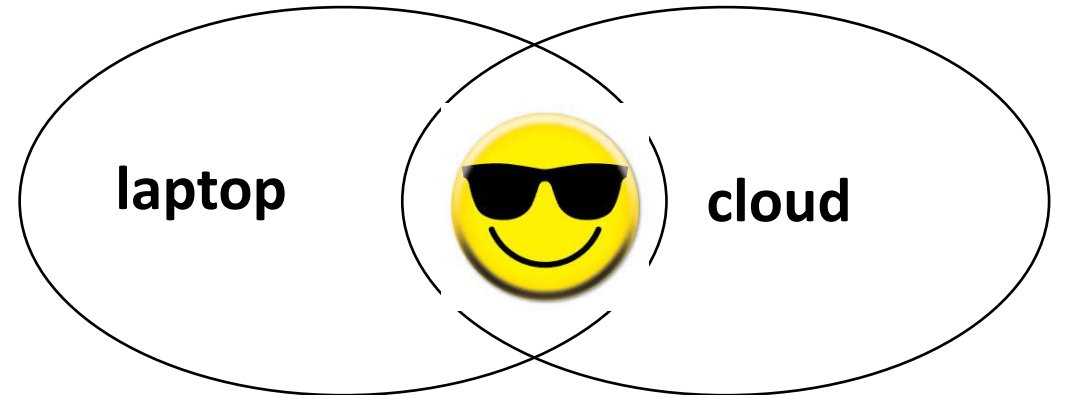
curriculum deployment



curriculum development



lab deployment



Design: start in vacuum, pivot with users

I. intro to course

- A. intro to space biology
- B. intro to data engineering
- C. intro to AI/ML
- D. intro to lab environment

II. building basic models

- A. using scikit-learn
- B. clustering
- C. classification
- D. regression

III. working with data

- A. OSDR & FAIR
- B. processing tabular data
- C. processing image data
- D. visualizing data

IV. interpreting results

- A. bioinformatic tools
- B. literature search
- C. ethical considerations
- D. reproducibility

V. advanced topics

- A. transfer learning
- B. causal inference
- C. explainable AI
- D. biological interpretations

VI. capstone projects

- A. use deep learning to find DNA damage in irradiated immune cells (microscopy)
- B. use causal inference machine learning ensemble to find genes correlated to lipid dysfunction (RNA-seq)

Development: build space bio, borrow AI/ML (CC)

coursera



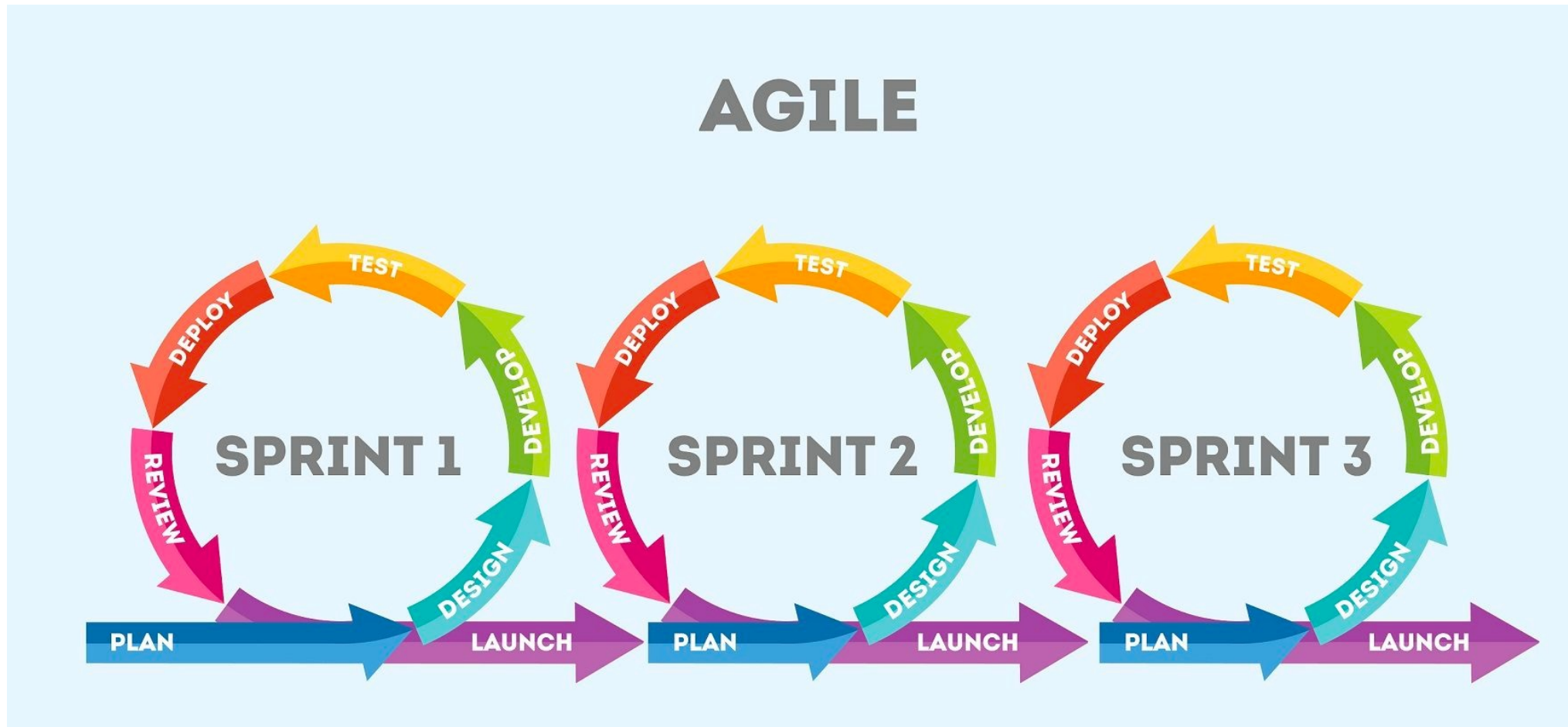
StatQuest!!!

towards
data science

Curriculum and lab deployment: Python notebooks



We are borrowing from Agile to manage the project



Which data will we use?

RNA-seq benchmarking data



https://registry.opendata.aws/bps_rnaseq

Microscopy benchmarking data



https://registry.opendata.aws/bps_microscopy



Extra credit: use AI/ML to design, develop, deliver

1. Syllabus design and course structure
2. Course overview, module introductions, and summaries
3. Discussion forums with chatbots
4. Code walkthroughs and commenting
5. Integrate context in case studies
6. Resource recommendations
7. Quiz/assessment generation
8. Design and generate feedback surveys

Join our user community!

<https://bit.ly/tops-aiml-users>



Acknowledgements

AI for Life in Space

- Lauren Sanders
- Sylvain Costes



Compute

NASA Center for Climate Simulation Science Managed Compute Environment

- Aaron Skolnik
- Andre Avelino Paniagua
- Ellen Salmon
- Daniel Duffy



Open Science for Life in Space Teams



Support

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