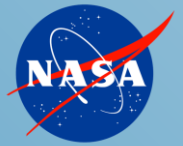# Architecting the Future: NASA's Use of Large Language Models to Enable Open Science

**Kaylin Bugbee**
NASA Marshall Space Flight Center
Project Lead, Science Discovery Engine
Project Scientist, Transform to Open Science (TOPS) Project Office

September 26, 2023

NASA

# Talk Outline

1. Open Science and AI/LLM Principles
2. NASA's Approach to LLM Enabled Search
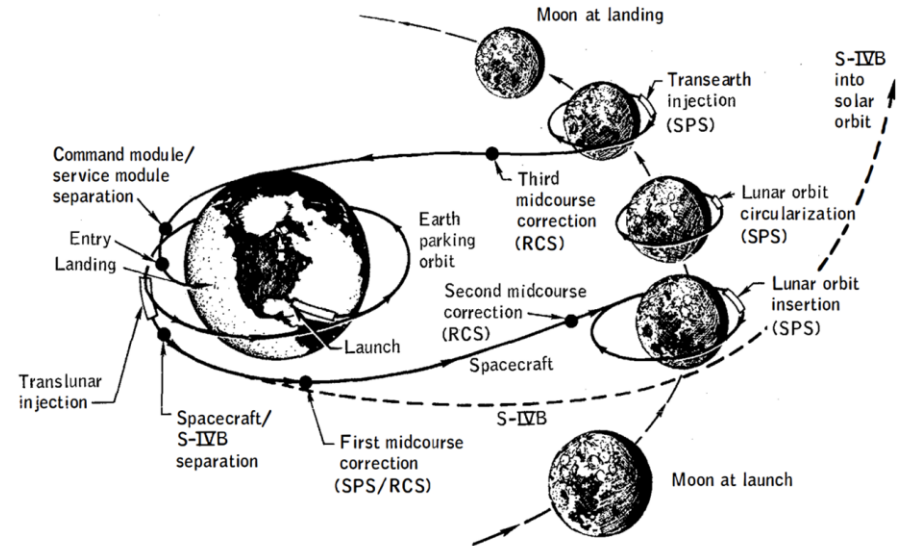3. Vision for the Future and Discussion



NASA-S-69-605

Figure 3-1.- Apollo 8 mission profile.
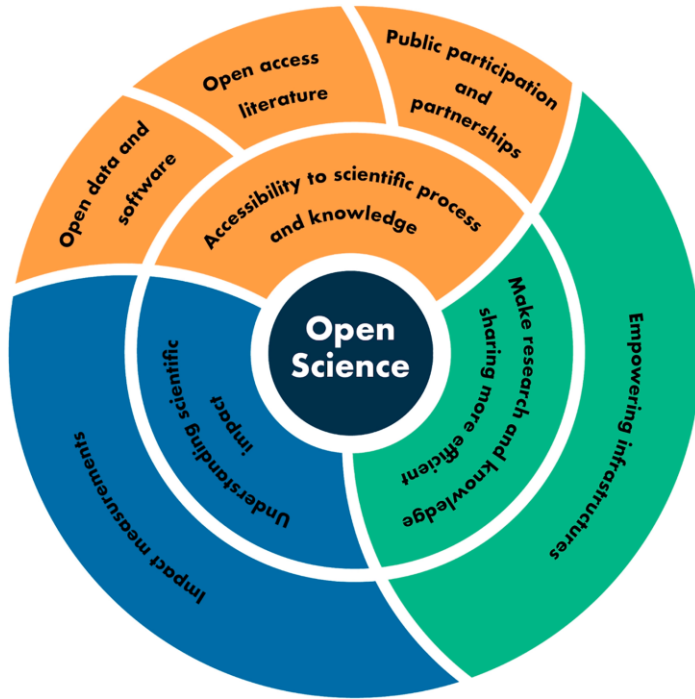
Image Credit: NASA

*"We are called to be architects of the future, not its victims." -R. Buckminster Fuller*

# Fears and Challenges Surrounding LLMs for Science

- Generative models are prone to hallucinations or making up answers

  - Problematic for agencies like NASA who have a responsibility to deliver trustworthy information to users

- LLMs are built on the shoulders of existing content

  - Propagation of bias

  - Lack of attribution

  - Gaps in content coverage

- Lack of clarity around the development process and data used for proprietary LLMs

  - Difficult to understand gaps, strengths, weaknesses

- ***These fears and challenges push up against open science values such as reproducibility, transparency, attribution and inclusion.***

# *Open Science and AI*

# What is Open Science?



"Open science is 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."

For data systems, open science has 3 core focus areas:

- Increasing accessibility to the scientific process & knowledge
- Making research & knowledge sharing more efficient
- Understanding scientific impact

***AI and Large Language Models will transform these 3 focus areas.***

Ramachandran, R., Bugbee, K., & Murphy, K. (2021). From open data to open science. *Earth and Space Science*, 8, e2020EA001562. https://doi.org/10.1029/2020EA001562

# AI Legislation, Guidance and Principles

- Legislation and government guidance exists to mitigate risks associated with AI
  - European Union's Artificial Intelligence Act
  - White House Office of Science and Technology Policy (OSTP)'s Blueprint for an AI Bill of Rights
- There are also a number of AI principles being discussed in the community including
  - HHH (Helpful, Honest, Harmless)
  - Google AI
  - Ethical and Responsible use of AI/ML in the Earth, Space, and Environmental Sciences
- There are currently no set of guiding principles for designing and implementing AI that were developed through the lens of open science.

**Safe and Effective Systems**

**Algorithmic Discrimination Protections**
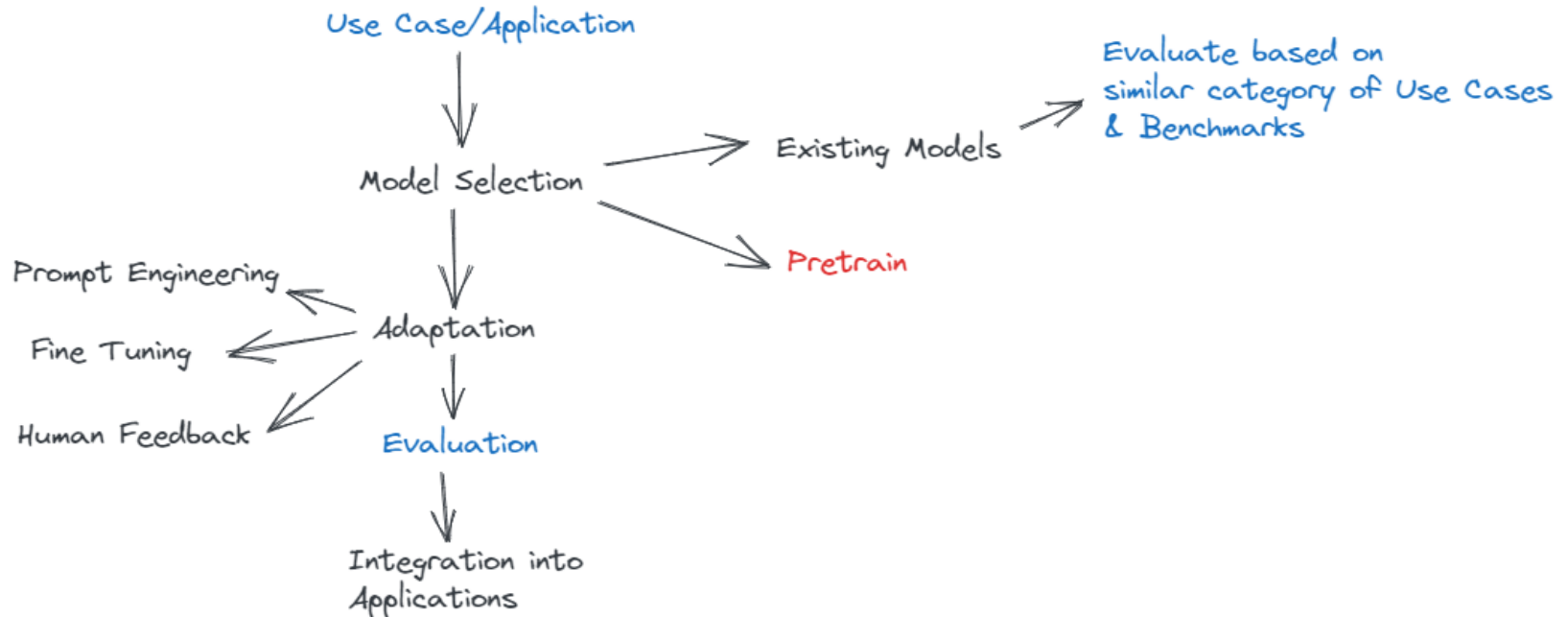
**Data Privacy**

**Notice and Explanation**
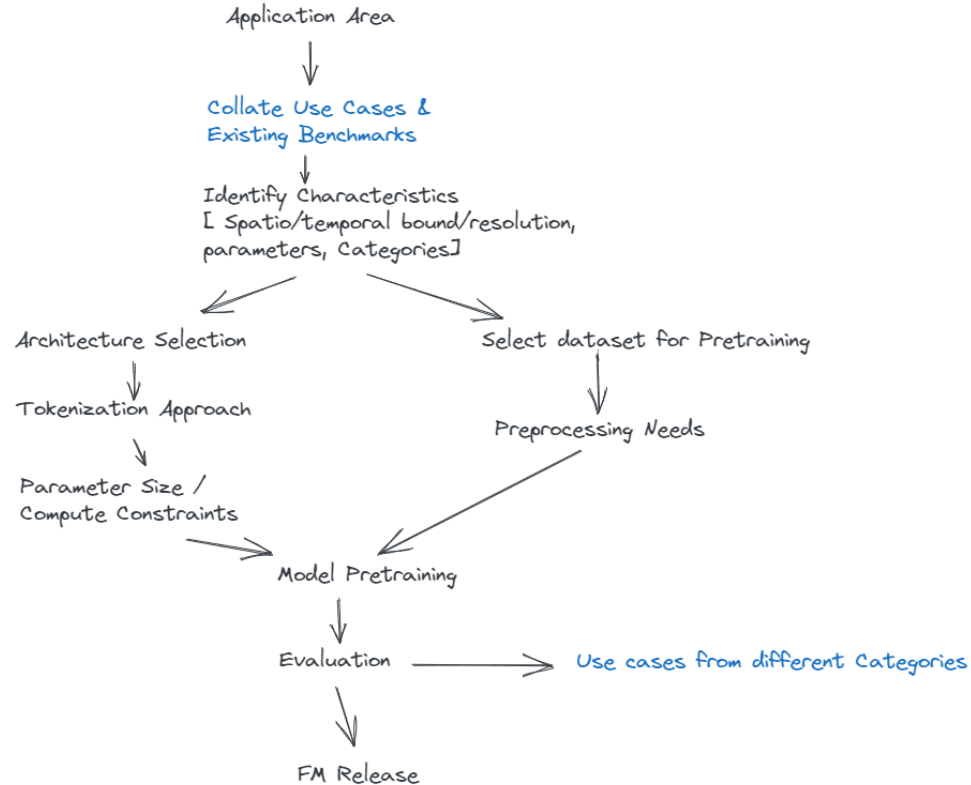
**Human Alternatives, Consideration, and Fallback**

Image Credit: OSTP

*For science to be open, we need to apply open science principles to the entire AI Life Cycle, not just the models themselves.*

# AI Project Lifecycle: Model Reuse

# AI Project Lifecycle: New Model Development



Image Credit: NASA IMPACT team

# Open Science Principles for the AI Lifecycle

### Transparency
Open Models
Open Workflows
Open Validation & Benchmarking
Open Data
Open Code

### Trust
Truthful
Attribution Given
Bias Reduced

### Teamwork
Collaborate
Community Participation
Share Resources

### Training
Workshops on AI Techniques
Easy-to-Understand Documentation
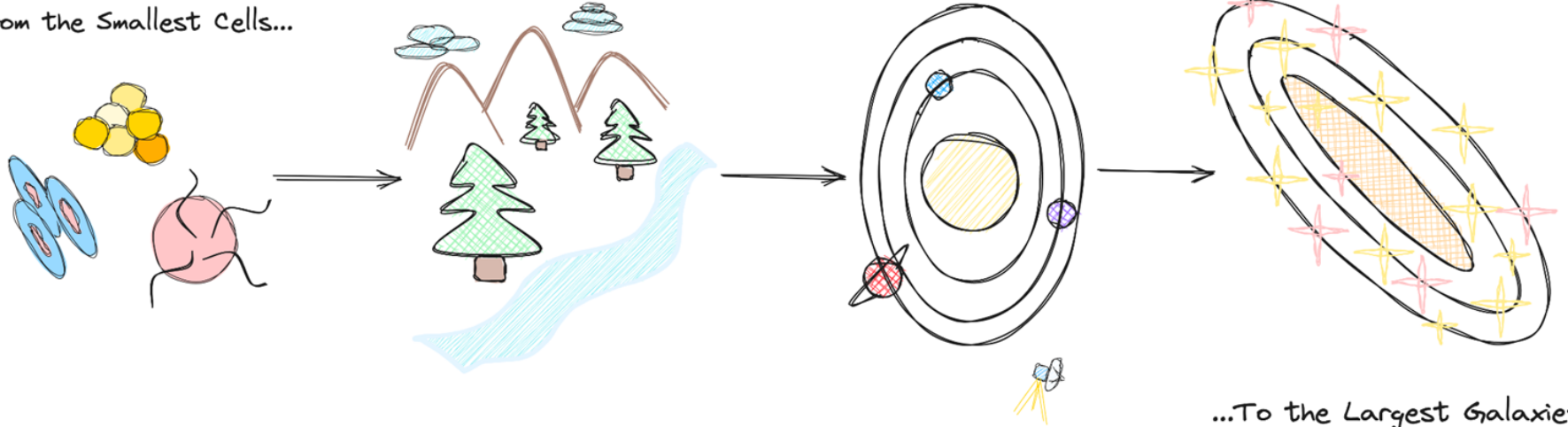Demonstration Notebooks

### Techniques
Retrieval Augmented Generation
Constitutional AI
Prompt Engineering

***The rapidly changing nature of AI requires emerging techniques be consistently monitored and assessed.

# *NASA's Approach to LLM Enabled Search*
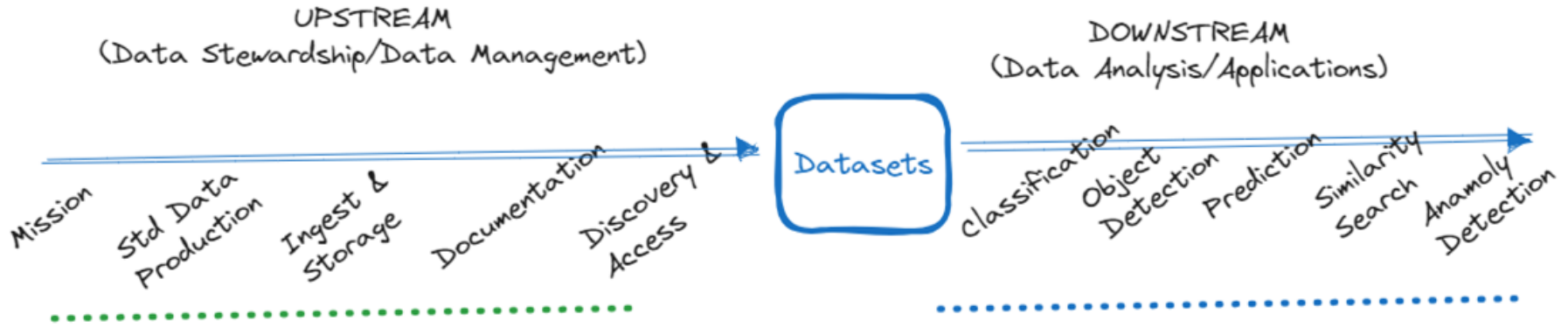
# Science at NASA
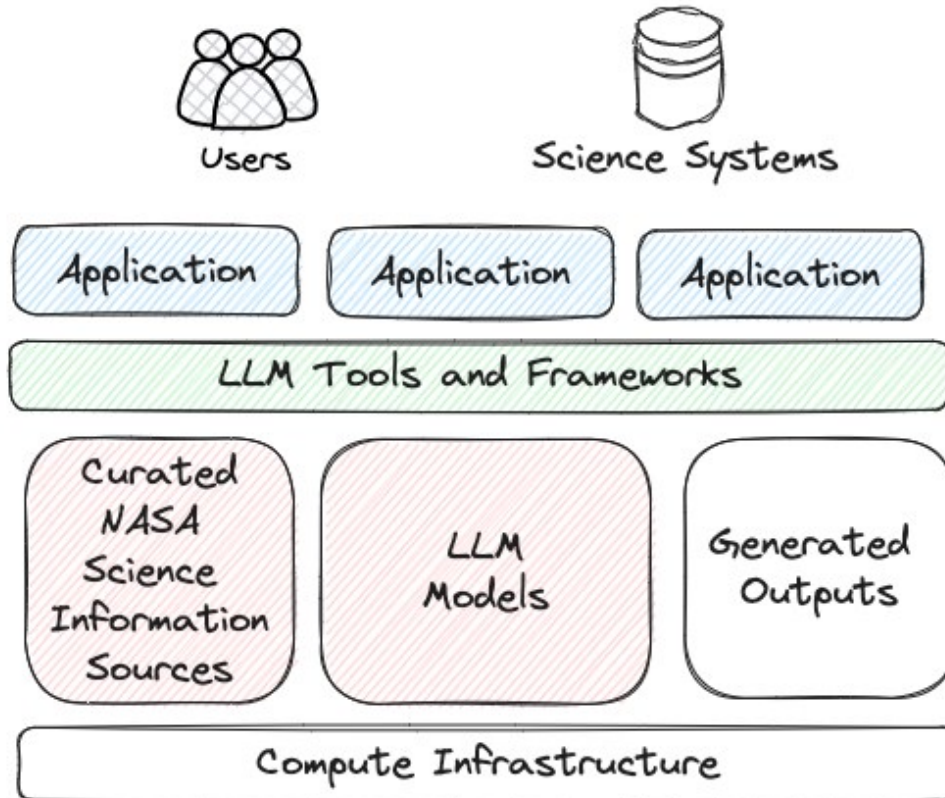


From the Smallest Cells...

...To the Largest Galaxies

"...the most powerful [ideas], the ones that are really transformative, they almost always come out of the collision between at least two or three or four different intellectual worlds or traditions or fields or disciplines. The border between disciplines is where the exciting stuff happens."
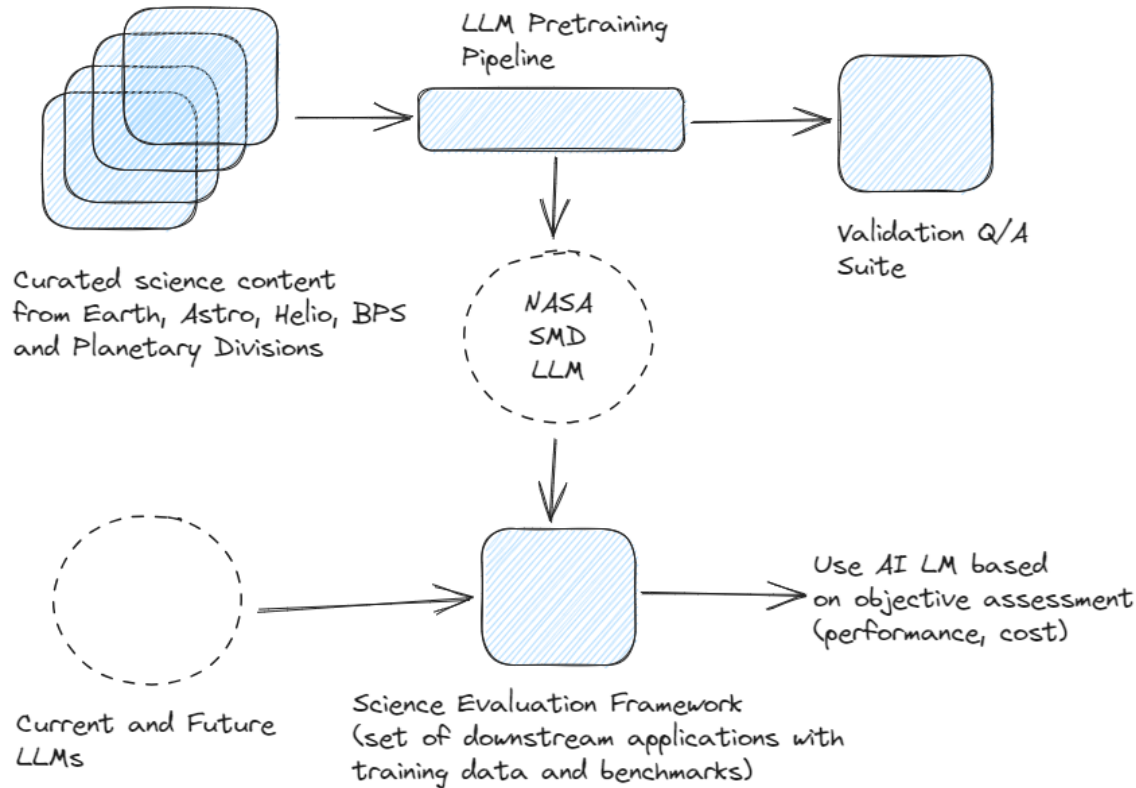
– Steven Johnson

# LLM Adoption in Science

# Platform for Building LLM-Based Applications

# Building a LLM for NASA Science
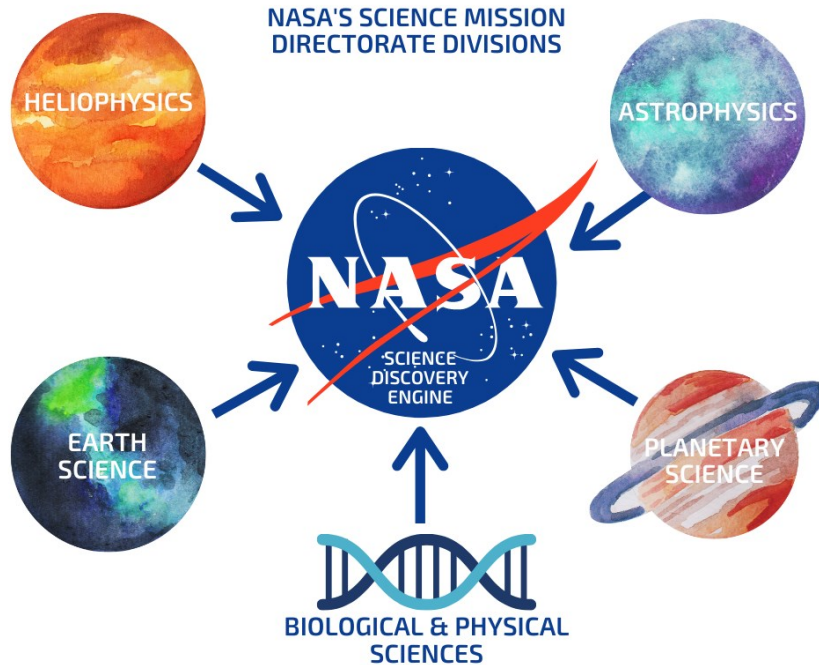


Image Credit: NASA IMPACT team

Collaboration between NASA science topical experts and IBM Research.

- Curating science resources (journal papers, reports, web pages) to train a NASA science specific domain model
- Utilize a common training pipeline to build the model
- Develop an evaluation suite to test the performance of this and *future* models

This LLM is in development now and we envision utilizing it for many different downstream applications.

# Application Example: Science Discovery Engine (SDE)



**NASA'S SCIENCE MISSION DIRECTORATE DIVISIONS**

HELIOPHYSICS

ASTROPHYSICS

NASA

SCIENCE DISCOVERY ENGINE

EARTH SCIENCE

PLANETARY SCIENCE

BIOLOGICAL & PHYSICAL SCIENCES

## <u>Goals:</u>

1. Enable rapid discovery of NASA science data, software and documentation

2. Support NASA's open science goals and infrastructure

3. Promote interdisciplinary science

4. Prototype emerging technologies and search techniques including Large Language Models

Image Credit: SDE team

# SDE: Challenges in Implementation

- Bringing heterogeneous topics into a single search environment

  - Domain-specific schemas

  - Different science concepts across the domains

    - Experiments versus Observations

- Identifying and curating dispersed content

  - Where does all this data and information live?

  - How do we know we have everything?

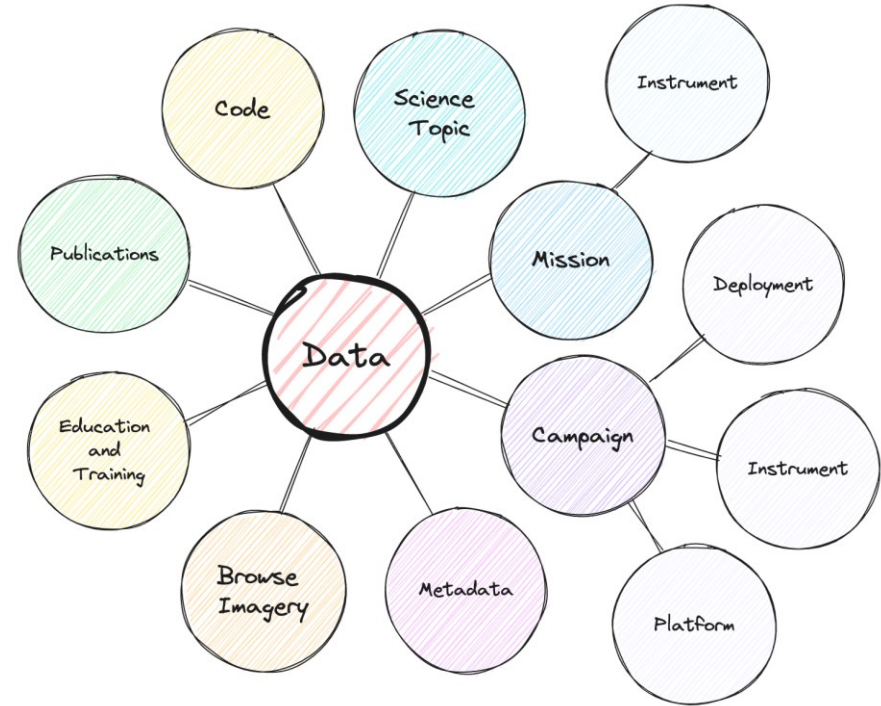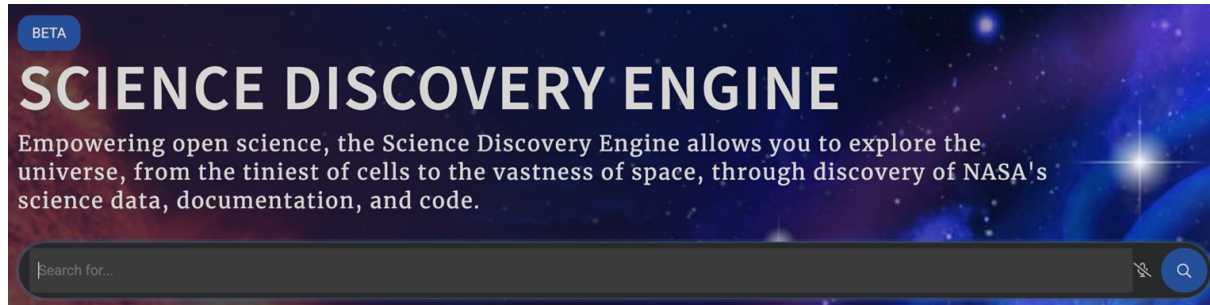- Enabling interdisciplinary search in a manner that is helpful



Image Credit: SDE team

# SDE Results

- The beta Science Discovery Engine (SDE) is **available now** and enables search of NASA's open science data and information.
- Includes over 84,000 metadata records about data and 600,000 documents.
- In the future, we plan to leverage the SMD LLM to refine and improve search in the SDE.



**sciencediscoveryengine.nasa.gov**

Image Credit: SDE team

# Application Example: Policy Compliance & Governance

- Scientists and data managers are required to follow both Federal and NASA policy on sharing scientific information
  - [NASA's Scientific Information Policy for the Science Mission Directorate (SPD-41a)](#)
- In addition, each science domain has community standards that need to be followed.
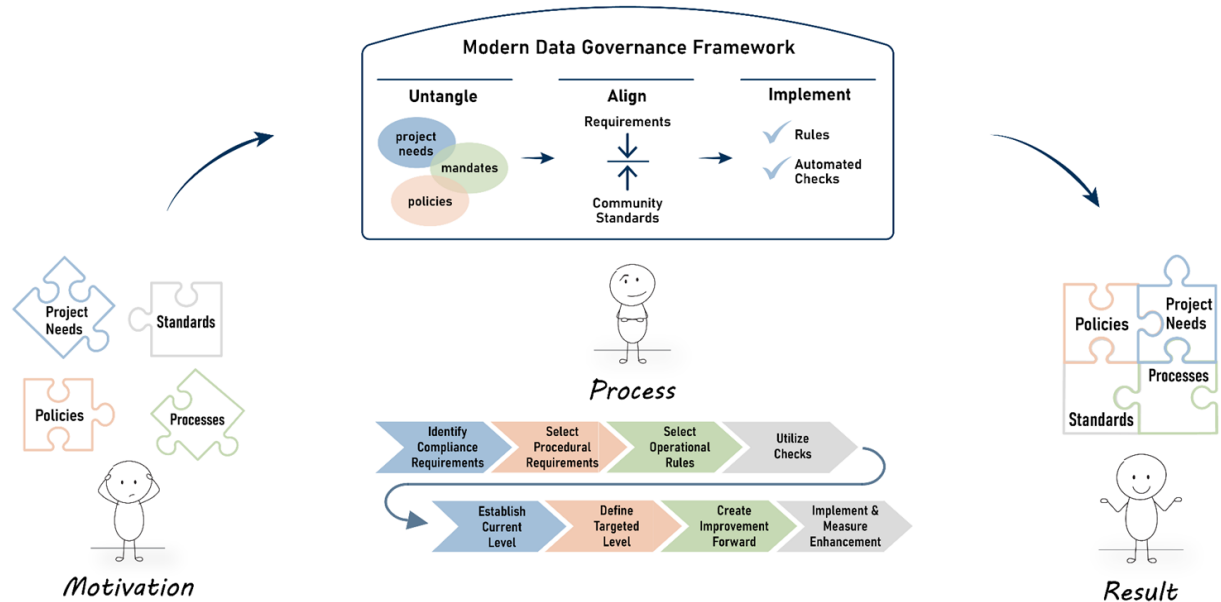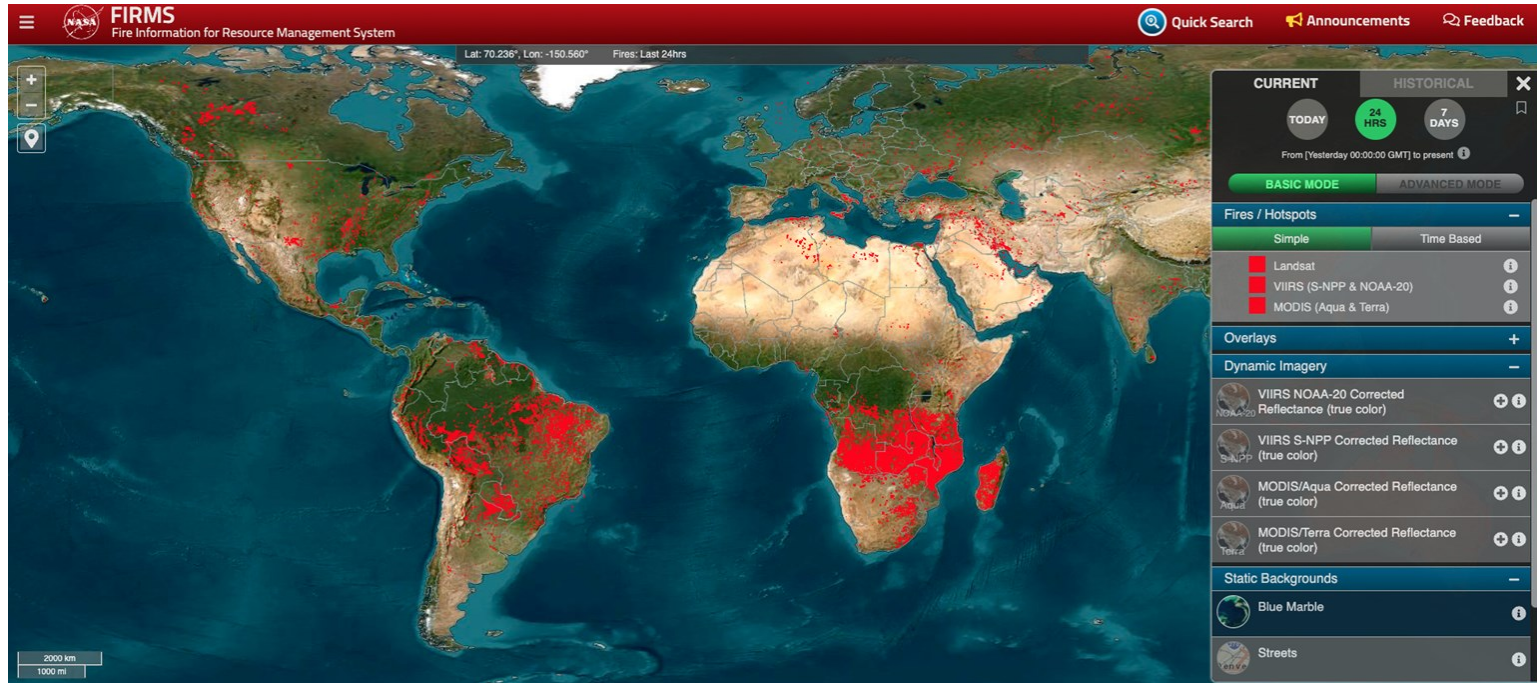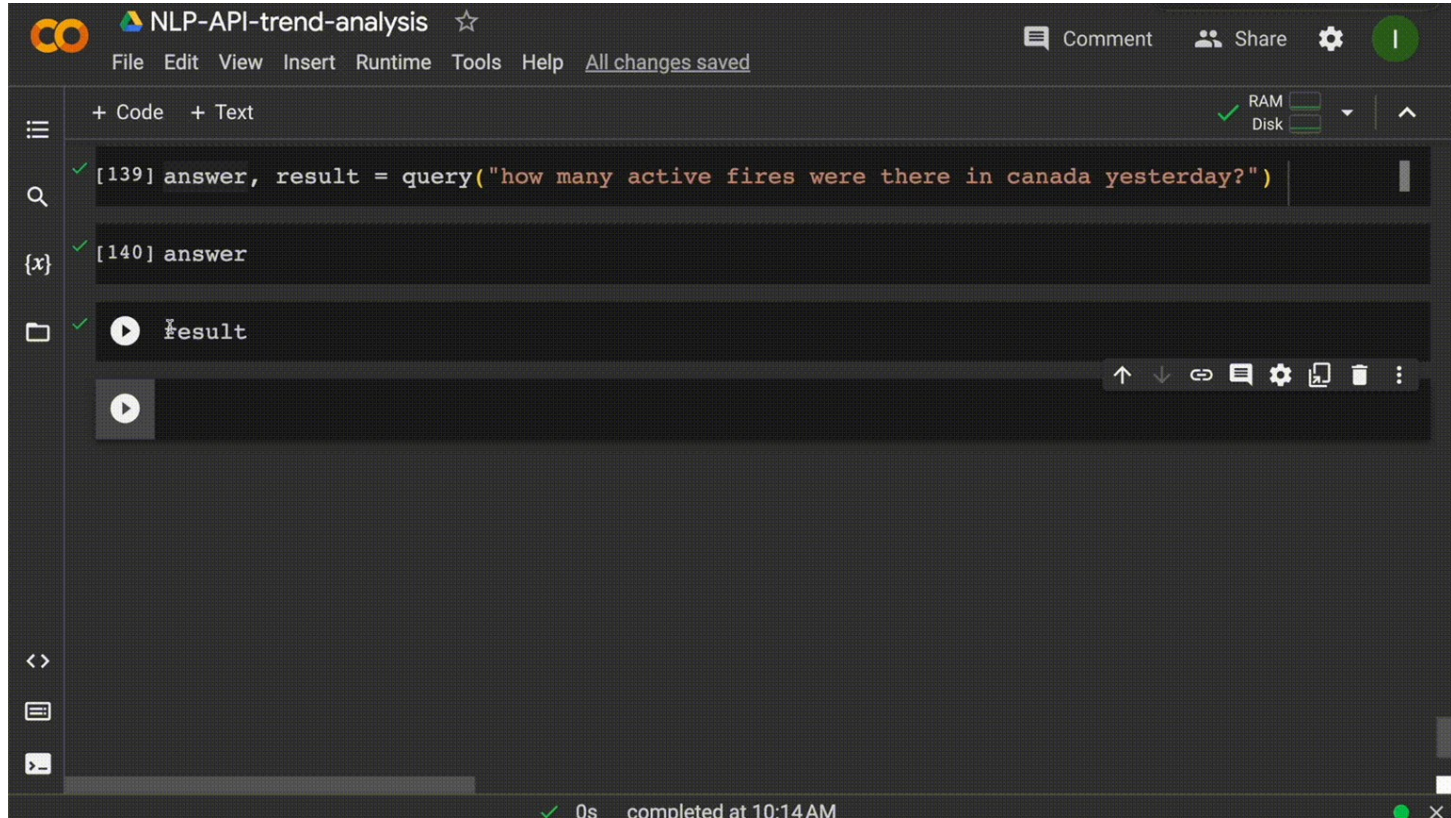- We are developing a Q&A service to allow a user to check for compliance



Image Credit: IMPACT team

# Example: Science Applications



Fire Information for Resource Management System (FIRMS) distributes Near Real-Time (NRT) active fire data

# Example: Science Applications FIRMS Q&A



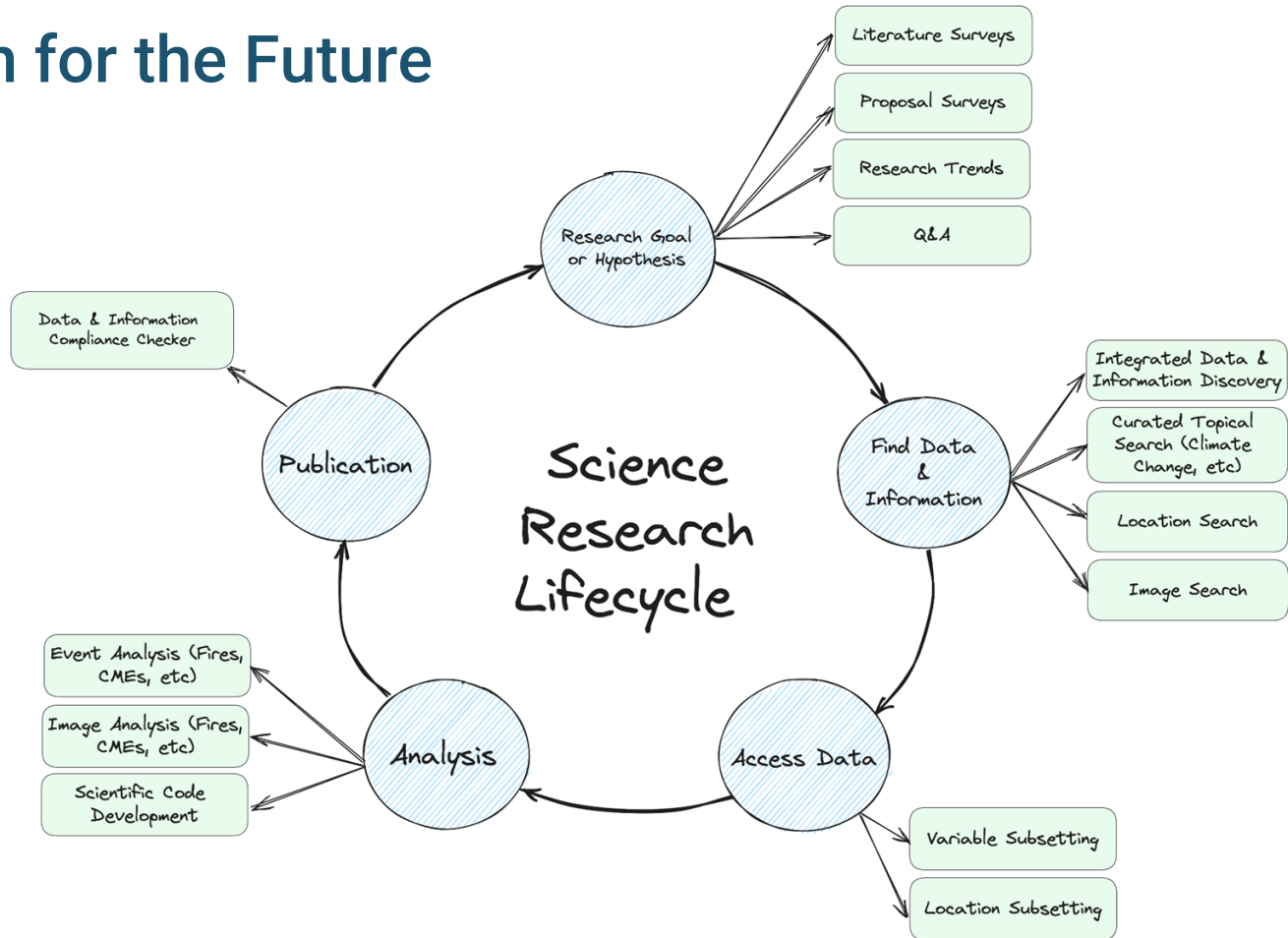Image Credit: NASA IMPACT team

# *Vision for the Future and Discussion*

# Discussion and Lessons Learned

- As a trusted scientific agency, NASA has a responsibility to provide knowledge with transparency and justification as to the source of that information.

    - Misinformation will always be a problem in the Internet age.

    - We need to have tools and techniques in place to both mitigate and detect misinformation.

    - We also need to systematically design and evaluate LLM workflows to ensure trustworthy results.

- Collaboration is essential for the successful development and implementation of LLMs

    - We shouldn't assume that all use cases are the same across diverse scientific disciplines - collaboration with scientific stakeholders is essential to ensure use cases and applications are developed for various needs.

    - Innovation is moving quickly in the LLM space - sustainable collaborations with partners like IBM help close the gap in understanding new techniques and architectures.

- Innovation is needed to allow for serendipitous discovery and to ensure epistemic diversity

    - New ideas happen when serendipitous discovery is possible. If only one answer is ever returned, this will limit the diversity of views into a topic. This poses problems for interdisciplinary science and non-traditional users such has decision makers, educators and the general public.

# Discussion and Lessons Learned

- We need to acknowledge existing biases in scientific systems.

  - There is a lot of discussion within the community about bias in AI systems. We should do everything in power to reduce bias.

  - However, we need to acknowledge that bias is already built into our existing scientific systems. Publications in and of themselves represent a biased sample of research conducted in that only research with successful results are published. Using things like citation counts or h-indexes for relevance is also another form of bias in that the frequently cited keep getting cited because they receive more exposure even though that publication may or may not be the best reference for a given search.

- We are committed to open science and open science principles for the AI lifecycle. However, not everyone will have the resources or collaboration opportunities to build open models. We should:

  - Make models open whenever possible so that others can benefit

  - Have a discussion within the science community as to whether there are some use cases, especially for implementation in applications, that do not require complete transparency and the use of open models.

# Vision for the Future

# NASA's Transform to Open Science (TOPS)

TOPS is a five-year initiative to promote open science through education and outreach efforts, by providing researchers and scientists with the tools, data, and resources they need to conduct advanced, inclusive, and impactful research.

## Objectives:

- Increase understanding and adoption of open science principles and techniques
- Broaden participation by historically excluded communities
- Accelerate major scientific discoveries

**Learn more at:**
**nasa.github.io/Transform-to-Open-Science**

## Open Science 101 - Coming Dec. 2023

Open Science 101: A **community-developed** introduction to open science with inclusivity, accessibility, and diversity at the forefront. Learn core open science skills:

- ✓ Writing open science and data management plans
- ✓ Open science tools and best practices (i.e. ORCID)
- ✓ Networking and connecting in the open science community

Pre-Enroll in OS101 and join our mailing list today!

Scan Here!

https://go.nasa.gov/3Y7MBsg

# *Thank You!*

**Contact:**
Kaylin Bugbee
kaylin.m.bugbee@nasa.gov