# FROM ARDS TO AODS: FUTURE OF ANALYTICS FOR EARTH OBSERVATIONS
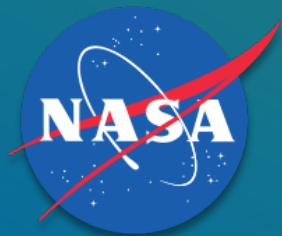
Dr. Rahul Ramachandran*, Kaylin Bugbee, Manil Maskey, Chris Lynnes

*Manager | Inter-Agency Implementation and Advanced Concepts Team (IMPACT)
Senior Research Scientist | Earth Science Branch (ST11)
Marshall Space Flight Center/NASA
Huntsville, Alabama 35812, USA

Acknowledgements: Hook Hua, George Change (JPL)
Aimee B (Dev Seed) + MAAP Team

IGARSS 2019, Yokohama, Japan

# Outline

- Concept Definitions
  - Existing ARD Definitions
  - Missing Component – Big Data Perspective
  - AODS Definition
- Notional Analytics Architecture for the Future
  - Background on NASA's Data and Information Systems
  - Drivers
  - Analytics Modes
  - Prototype: MAAP Project
  - Summary

Concept Definitions

# Analysis Ready Data

"Consistently processed to the highest scientific standards and level of processing required for direct use in monitoring and assessing landscape change" - U.S Landsat ARD

"Time-series stacks of overhead imagery that are prepared so that a user can analyze the data without having to pre-process the imagery themselves" - Planet

"Analysis Ready Data for Land (CARD4L) are satellite data that have been processed to a minimum set of requirements and *organized into a form* that allows immediate analysis with a minimum of additional user effort and interoperability both through time and with other datasets" - CEOS

- Lower barrier to data use
- Focus on (pre)processing needs of users (subsetting, atmospheric corrections, co-registration, and sensor calibration)

- Not new - maps directly to NASA's level 3 processing definition of providing "variables mapped on uniform space-time grid scales, usually with some completeness and consistency"

# Missing Component: Big Data Perspective

- **Data engineering**, or restructuring, is important from the big data analytics perspective because restructuring easily enables mapping data into analytics tools

- Characteristics of big data analytics tools :
  - Exploit parallelism
  - Minimize data movement between compute nodes
  - Utilize commodity computing to reduce cost
  - Scale to handle growth in data volumes.

# Analytics Optimized Data Stores (AODS)

- Addresses both the preprocessing requirements as well as the data engineering requirement.
  - Focus on restructuring and distributing the data in order to exploit parallelism and minimize data movement during computing.

- Characteristics:
  - Serve specific analytic goals
  - Preprocessed with the right steps to meet analytic goals
  - Engineered (data) to enable fast access patterns tuned to the expected analyses
  - Can be subsets of a total archive holding, preprocessed using consistent and well-known algorithms that cover a spectrum of steps

# Notional Analytics Architecture

# Earth Science – NASA's Strategic Goal

**This ability to *observe our planet comprehensively* matters to each of us, on a daily level. Earth information—for use in Internet maps, daily weather forecasts, land use planning, transportation efficiency, and agricultural productivity, to name a few—is central to our lives, providing substantial contributions to our economies, our national security, and our personal safety. It helps ensure we are a thriving society.  - NRC, 2018**

NASA's Strategic Goal 1.1:
"Understand The Sun, Earth, Solar System, And Universe."



National Aeronautics and Space Administration

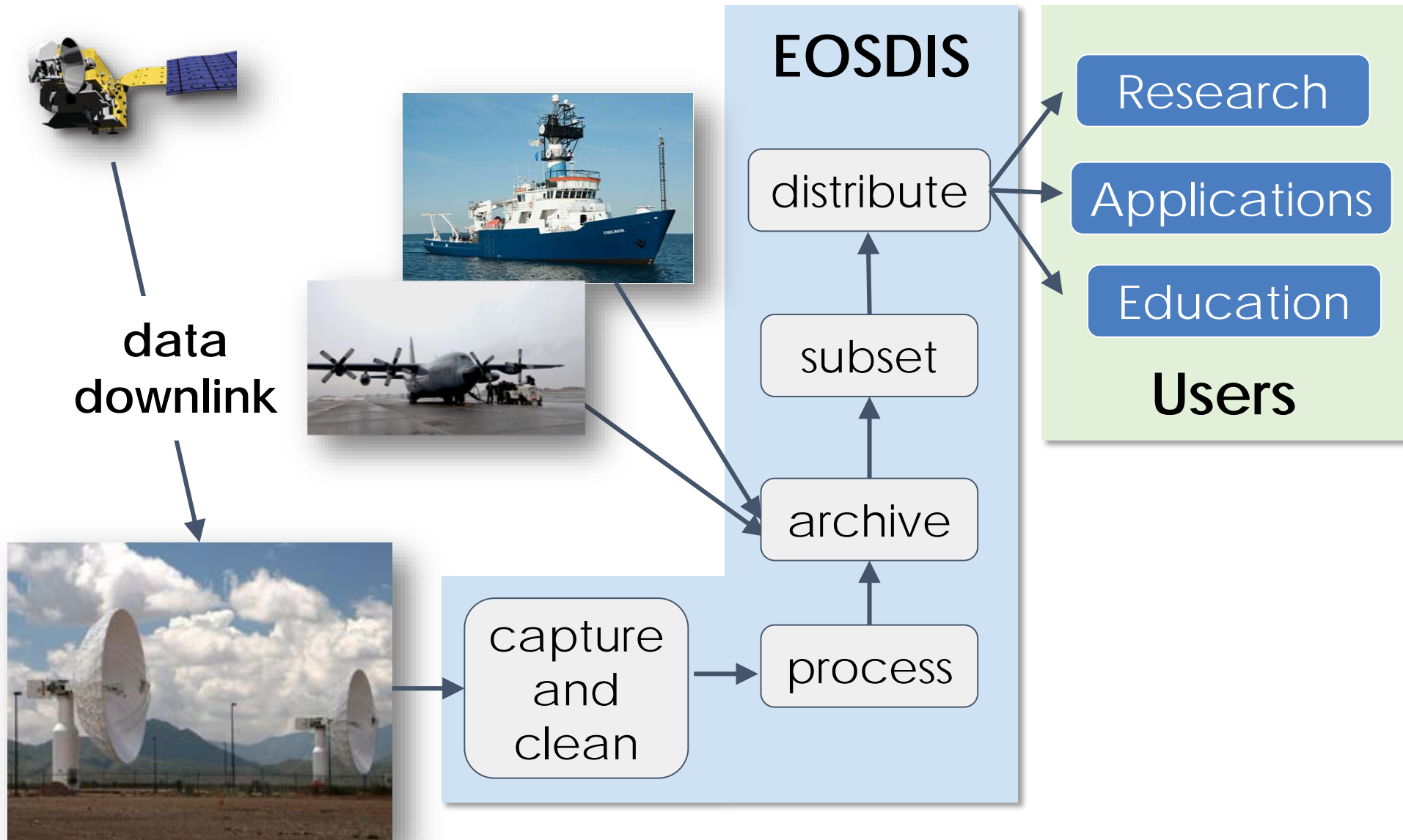NASA STRATEGIC PLAN 2018

www.nasa.gov

# Earth Science Data System Program

The Earth Science Data System Program is an essential *component of the Earth Science Division* and is responsible for:

- Actively *managing NASA's Earth science data* (Satellite, Airborne, and Field).
- *Developing unique* data system *capabilities* optimized to support *rigorous science investigations* and interdisciplinary research.
- *Processing* (and reprocessing) instrument data to create high quality long-term Earth science data records.
- Upholding NASA's policy of *full and open sharing of all data, tools, and ancillary information* for all users.
- Engaging members of the Earth science community in the *evolution of data systems*.

# Earth Observing System Data and Information System (EOSDIS)

EOSDIS is managed by the Earth Science Data and Information System (ESDIS) Project at GSFC and includes the following major core components:
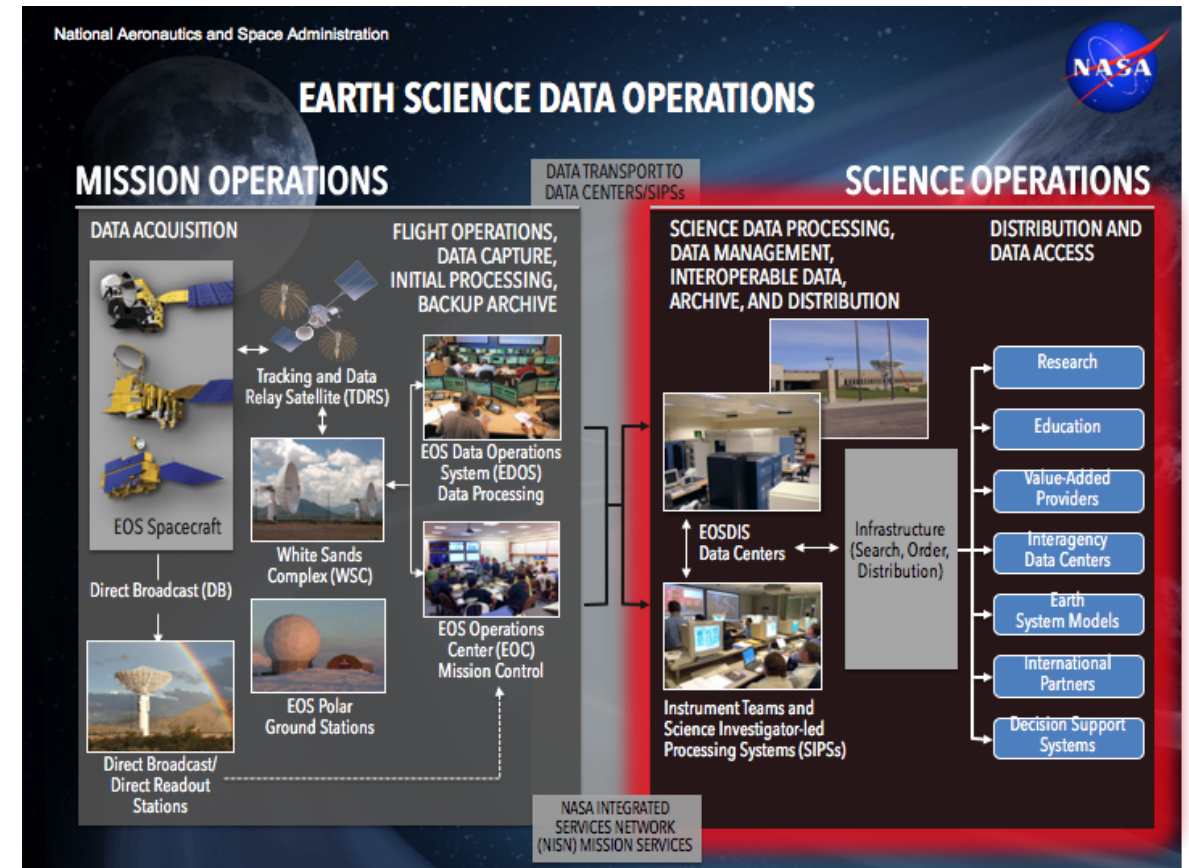
### Science Investigator-led Processing Systems (SIPS)

- Perform forward processing of standard data products and reprocess data to incorporate algorithm improvements

### Distributed Active Archive Centers (DAACs)

- Co-located with centers of science discipline expertise; archive and distribute standard data products produced by the SIPS and others
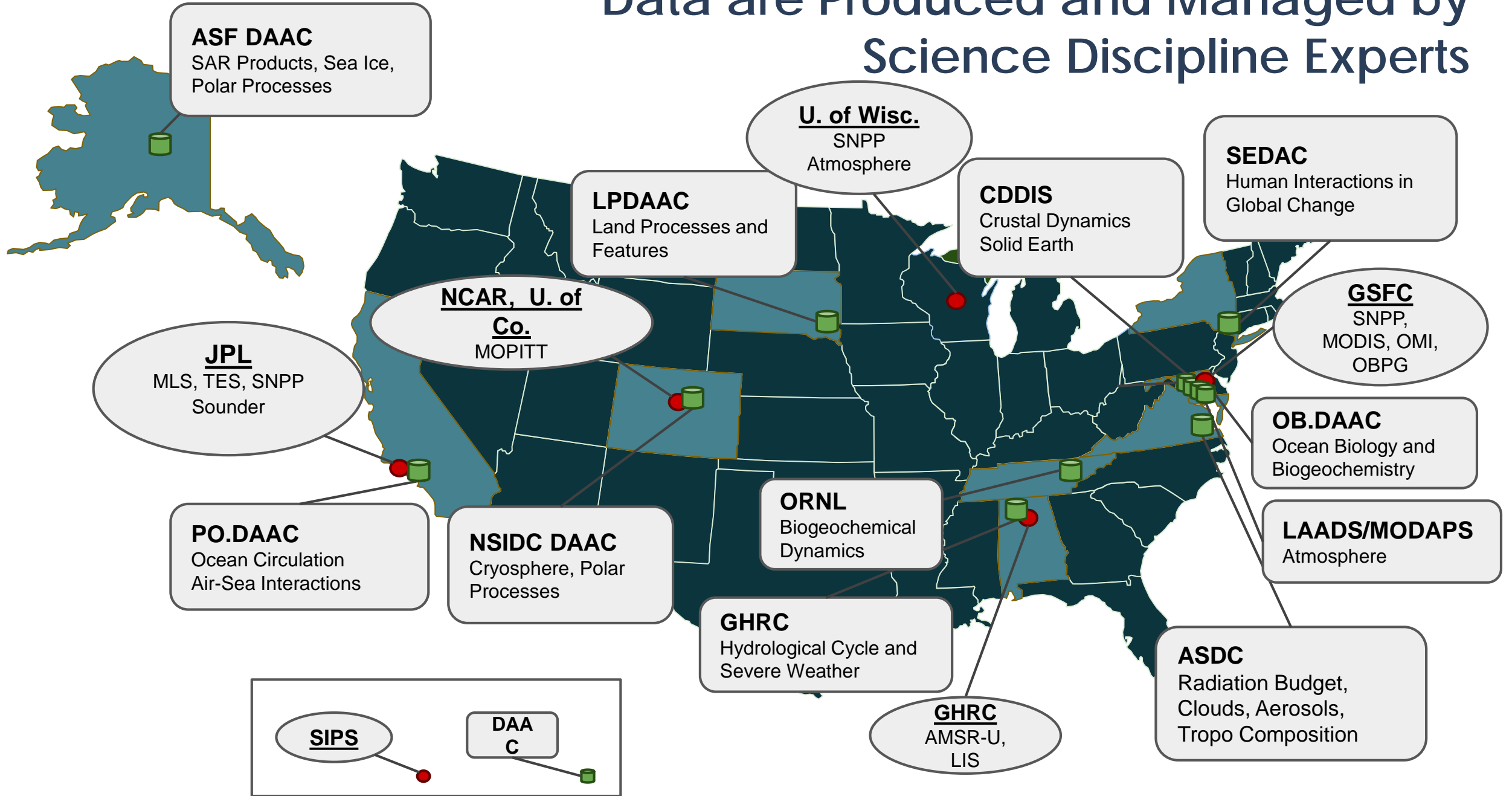
### Earthdata and Core Services

- Allows users to search, discover, visualize, refine, and access NASA Earth Observation data. Includes networking and security



*Red highlight indicates EOSDIS boundary.*

# Data are Produced and Managed by Science Discipline Experts

**ASF DAAC**
SAR Products, Sea Ice, Polar Processes

**U. of Wisc.**
SNPP Atmosphere

**LPDAAC**
Land Processes and Features

**CDDIS**
Crustal Dynamics Solid Earth

**SEDAC**
Human Interactions in Global Change

**GSFC**
SNPP, MODIS, OMI, OBPG

**NCAR, U. of Co.**
MOPITT

**JPL**
MLS, TES, SNPP Sounder

**OB.DAAC**
Ocean Biology and Biogeochemistry

**PO.DAAC**
Ocean Circulation Air-Sea Interactions

**NSIDC DAAC**
Cryosphere, Polar Processes

**ORNL**
Biogeochemical Dynamics

**LAADS/MODAPS**
Atmosphere

**GHRC**
Hydrological Cycle and Severe Weather

**GHRC**
AMSR-U, LIS

**ASDC**
Radiation Budget, Clouds, Aerosols, Tropo Composition

**SIPS**

**DAAC**

# Extensive Data Collection

## Started in the 1990s, EOSDIS today has 11,000+ data types (collections)

**Land**
- Cover & Usage
- Surface temperature
- Soil moisture
- Surface topography

**Ocean**
- Surface temperature
- Surface wind fields & heat flux
- Surface topography
- Ocean color

**Atmosphere**
- Winds & Precipitation
- Aerosols & Clouds
- Temperature & Humidity
- Solar radiation

**Human Dimensions**
- Population & Land Use
- Human & Environmental Health
- Ecosystems

**Cryosphere**
- Sea/Land Ice & Snow Cover

# NASA's Earth Science Data System in 2018

EOSDIS currently has over **27 Petabytes** of accessible Earth science data

Easy access and discovery of data to over **12,500 unique data products**

… of which 95% of granule searches complete in less than **1 Second**

EOSDIS delivered over **1.6 Billion** data products to over **4.1 Million** users from around the world

**33,000** Data Collections in the Common Metadata Repository (CMR)

EOSDIS also delivers near-real-time products in under **3 hours** from observation …

Over **330,000 users** have registered with EOSDIS to date

And Over **380 Million** data granules

American Customer Satisfaction Index (ACSI) survey scoring **79** from over **4,000** respondents

ACSI — American Customer Satisfaction Index®

*https://earthdata.nasa.gov/about/system-performance*

# Current EOSDIS Architecture



**Strengths**

- Discipline specific support and tools (DAAC data)
- Optimized for discipline specific archive, search and distribution
- Easily add new data products
- Supports millions of users

**Weakness**

- Uneven service and performance
- Significant interface coordination
- Limited on-demand product generation
- *Fragmentation – duplication of services, software and storage*
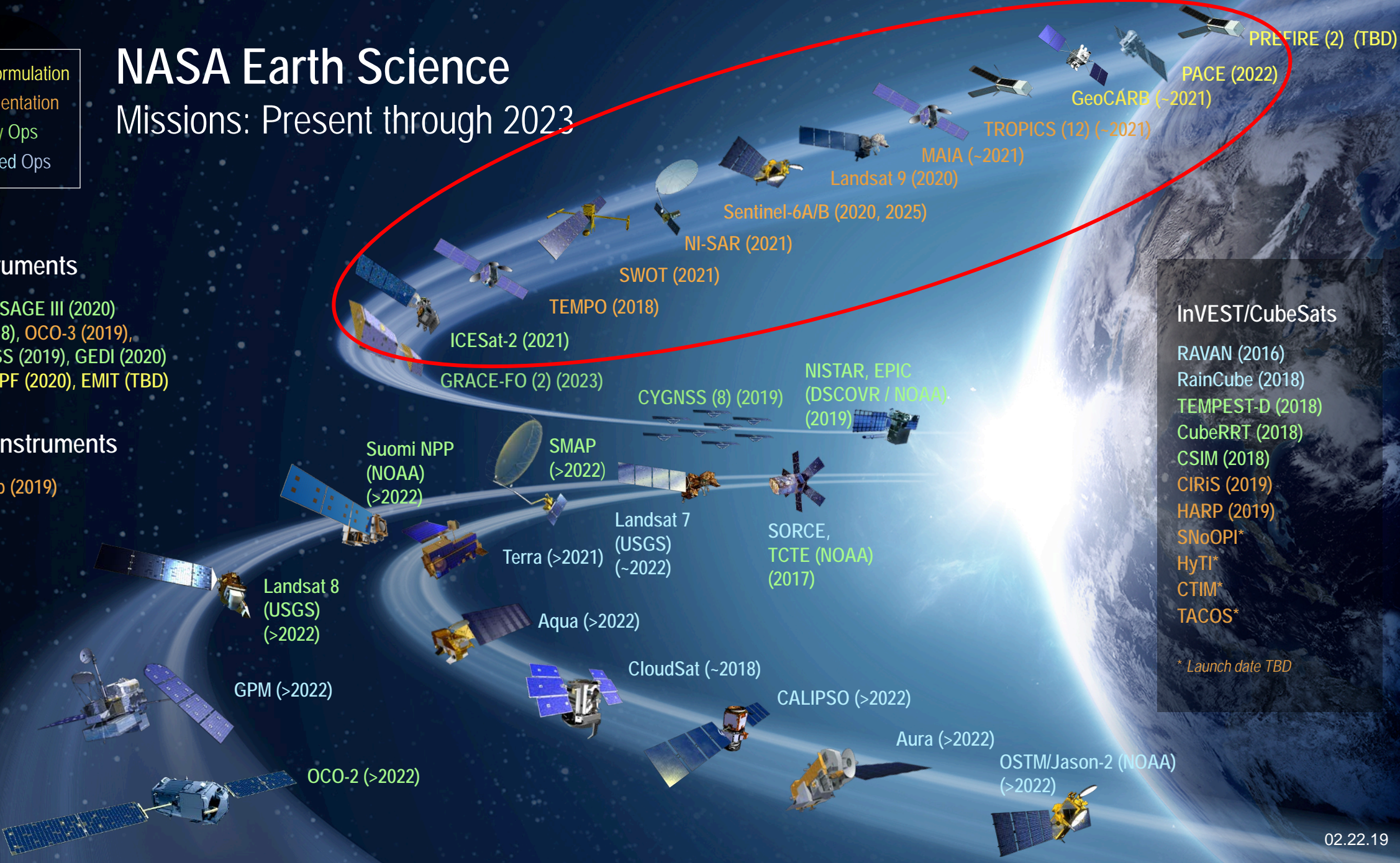
# NASA Earth Science
## Missions: Present through 2023

**Legend:**
- (Pre)Formulation
- Implementation
- Primary Ops
- Extended Ops

**ISS Instruments**

LIS (2020), SAGE III (2020)
TSIS-1 (2018), OCO-3 (2019),
ECOSTRESS (2019), GEDI (2020)
CLARREO-PF (2020), EMIT (TBD)

**JPSS-2 Instruments**

OMPS-Limb (2019)

PREFIRE (2) (TBD)
PACE (2022)
GeoCARB (~2021)
TROPICS (12) (~2021)
MAIA (~2021)
Landsat 9 (2020)
Sentinel-6A/B (2020, 2025)
NI-SAR (2021)
SWOT (2021)
TEMPO (2018)
ICESat-2 (2021)
GRACE-FO (2) (2023)
CYGNSS (8) (2019)
NISTAR, EPIC (DSCOVR / NOAA) (2019)
Suomi NPP (NOAA) (>2022)
SMAP (>2022)
Landsat 7 (USGS) (~2022)
Terra (>2021)
SORCE, TCTE (NOAA) (2017)
Landsat 8 (USGS) (>2022)
Aqua (>2022)
CloudSat (~2018)
CALIPSO (>2022)
GPM (>2022)
Aura (>2022)
OSTM/Jason-2 (NOAA) (>2022)
OCO-2 (>2022)

**InVEST/CubeSats**

RAVAN (2016)
RainCube (2018)
TEMPEST-D (2018)
CubeRRT (2018)
CSIM (2018)
CIRiS (2019)
HARP (2019)
SNoOPI*
HyTI*
CTIM*
TACOS*

* Launch date TBD

02.22.19

# EOSDIS Data System Evolution

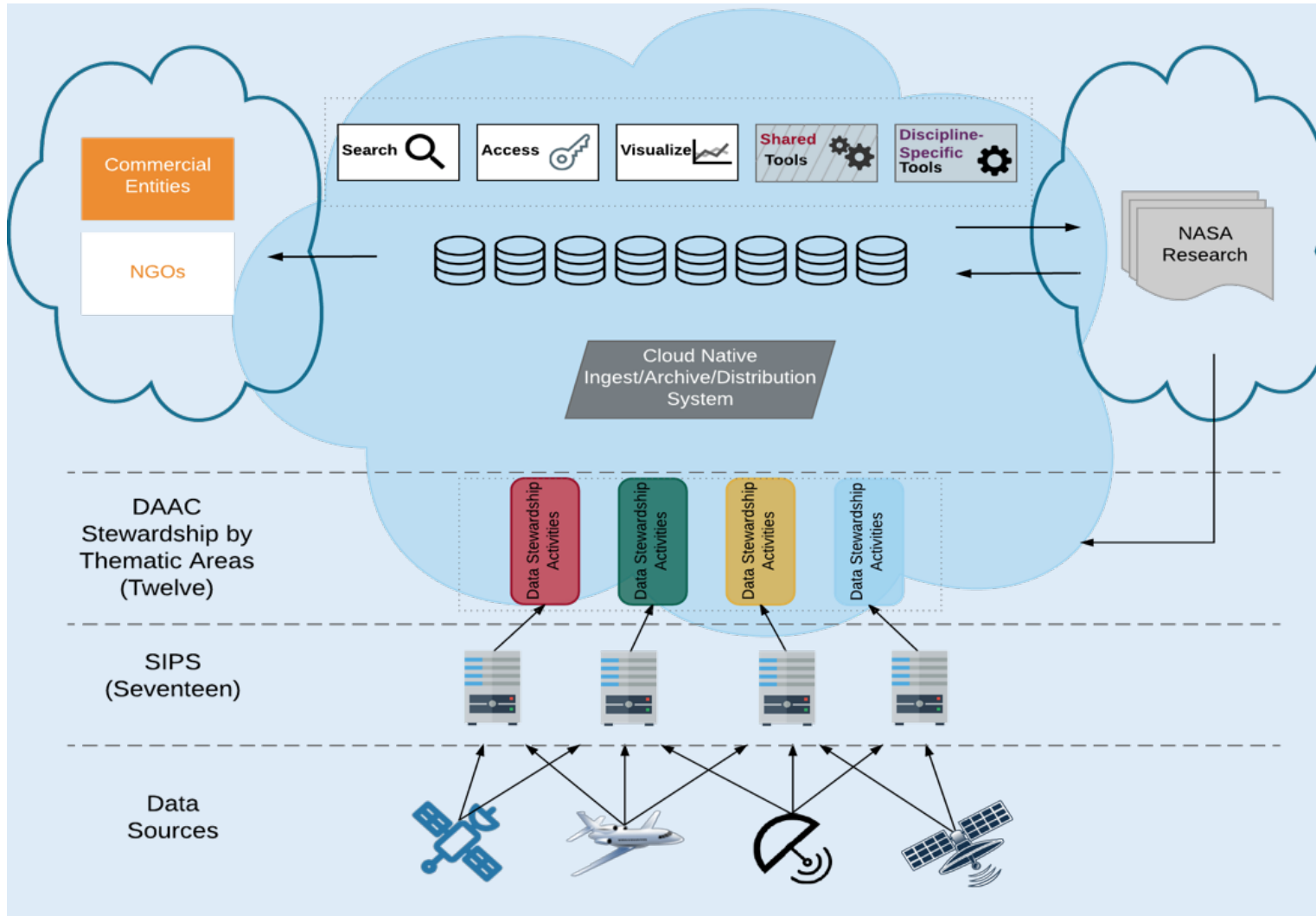The current architecture will not be cost effective as the annual ingest rate increases from 4 to 50PB/year

EOSDIS is developing open source cloud native software for reuse across the agency, throughout the government and for any other user.



**Cloud** *offers benefits like the ability to analyze data at scale, analyze multiple data sets together easily and avoid lengthy expensive moves of large data sets allowing scientists to work on data "in place"*

# Conceptual Cloud Architecture - 2021

## Open Science = Open Data + Open Source Software + Open Services



**Advantages**
- Scalability
- Processing next to data for anyone
- Optimized for multidisciplinary research

**Challenges**
- Develop coordination and documentation
- Cost management
- Business processes, security and skillsets
- Vendor lock-in

# Moving towards the cloud: Cumulus

*Lightweight, cloud-native framework for data ingest, archive, distribution and management*

A lightweight framework consisting of:

**Tasks** a discrete action in a workflow, invoked as a Lambda function or EC2 service, common protocol supports chaining

**Orchestration engine** (AWS Step Functions) that controls invocation of tasks in a workflow

**Database** store status, logs, and other system state information

**Workflows(s)** file(s) that define the ingest, processing, publication, and archive operations (json)

**Dashboard** create and execute workflows, monitor system

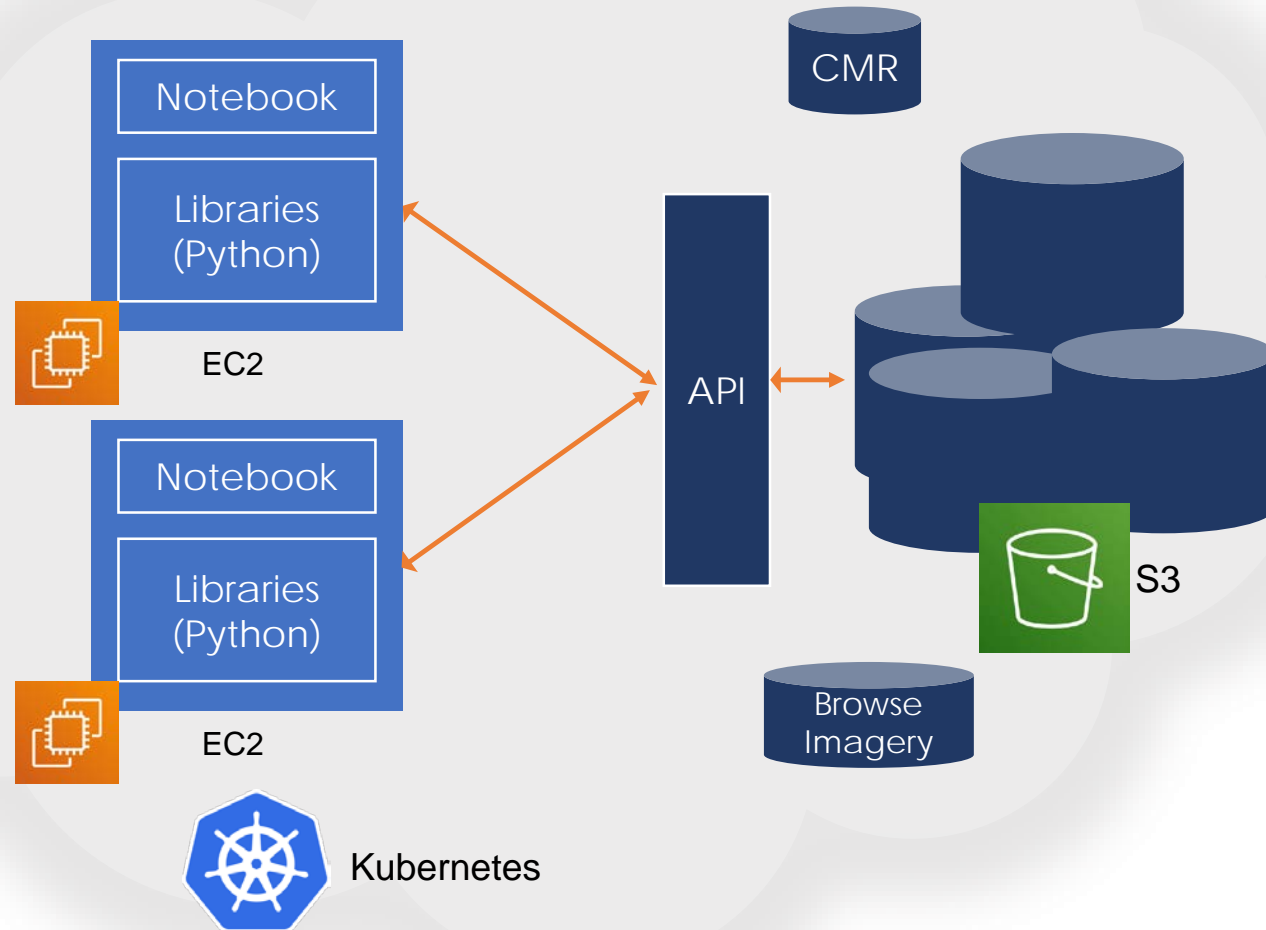How do we maximize the use of NASA data?

# Conceptual Architecture Supporting Analytics



**Three Modes of Data Interaction:**

- **Small Data/Interactive analysis**
- **Batch Processing**
- **Ephemeral AODS enabling big data analytics**

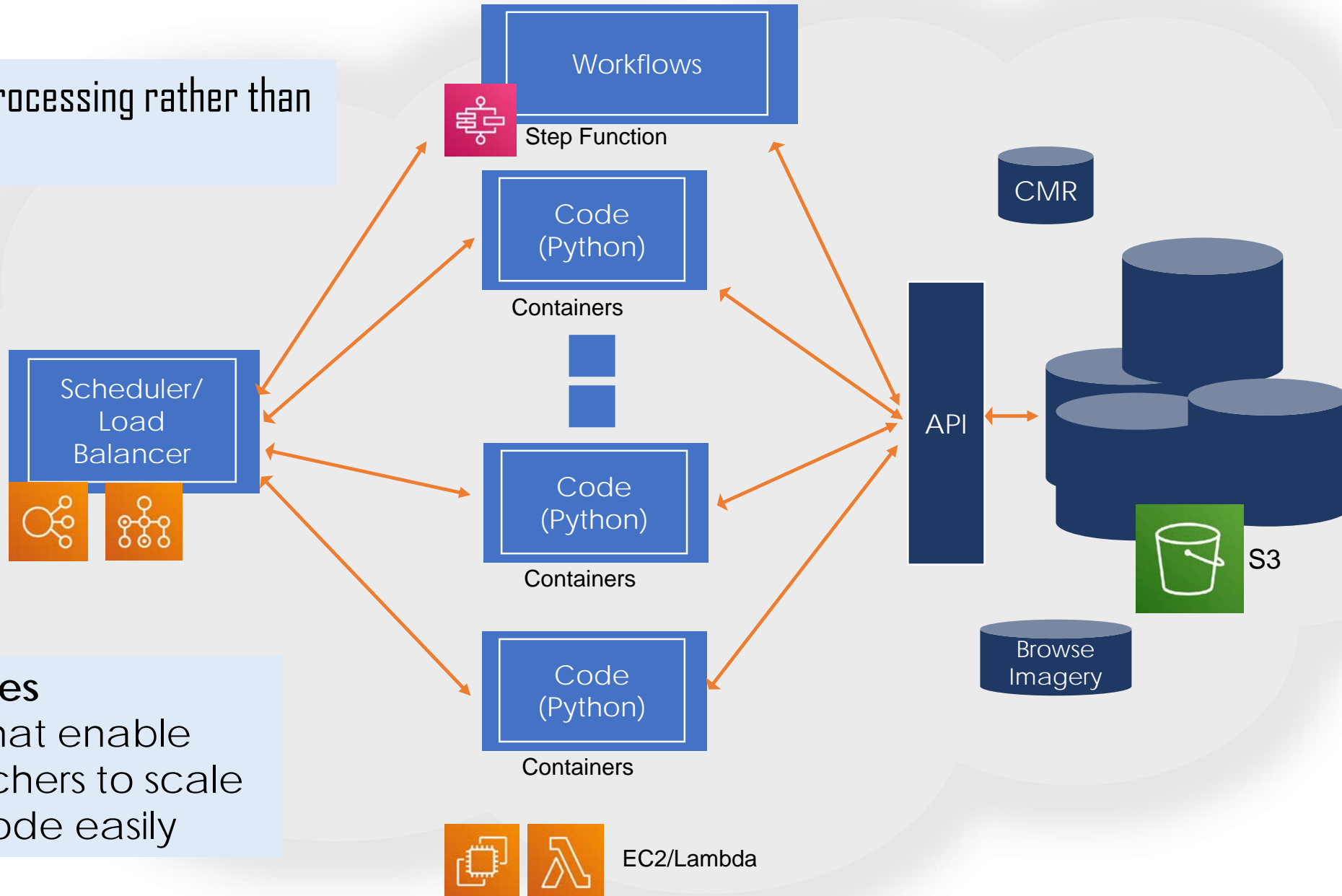# Mode 1: Small Amounts of Data/Interactive Analysis



- Moving individual researchers to the cloud
- Shifting the paradigm from downloading files for data analysis

**Challenges**
- Who provides and pays for the computing environment (eg. Pangeo)
- Researchers need to learn to instantiate their instances and manage their cloud accounts

# Mode 2: Scaling Out/Batch Processing

Focus is processing rather than analysis

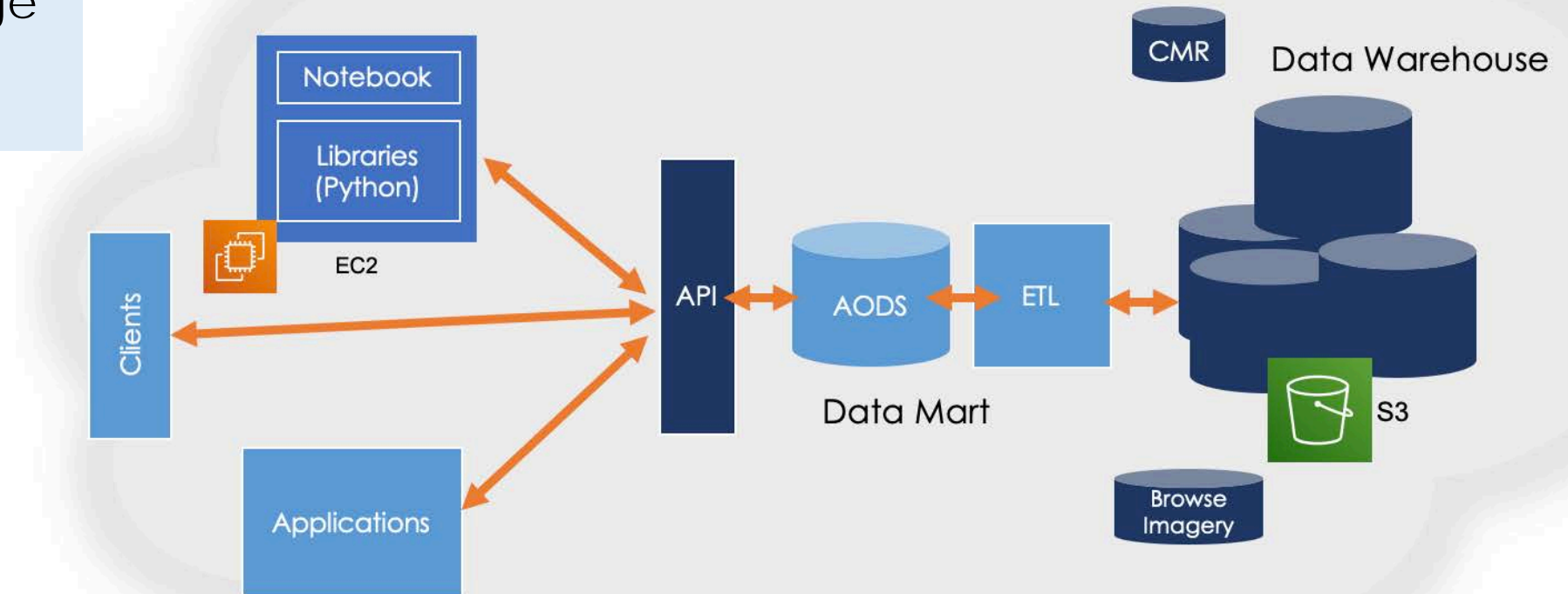Workflows

Step Function

Code (Python)

Containers

Code (Python)

Containers

Code (Python)

Containers

Scheduler/ Load Balancer

API

CMR

S3

Browse Imagery

EC2/Lambda

**Challenges**
- Tools that enable researchers to scale their code easily
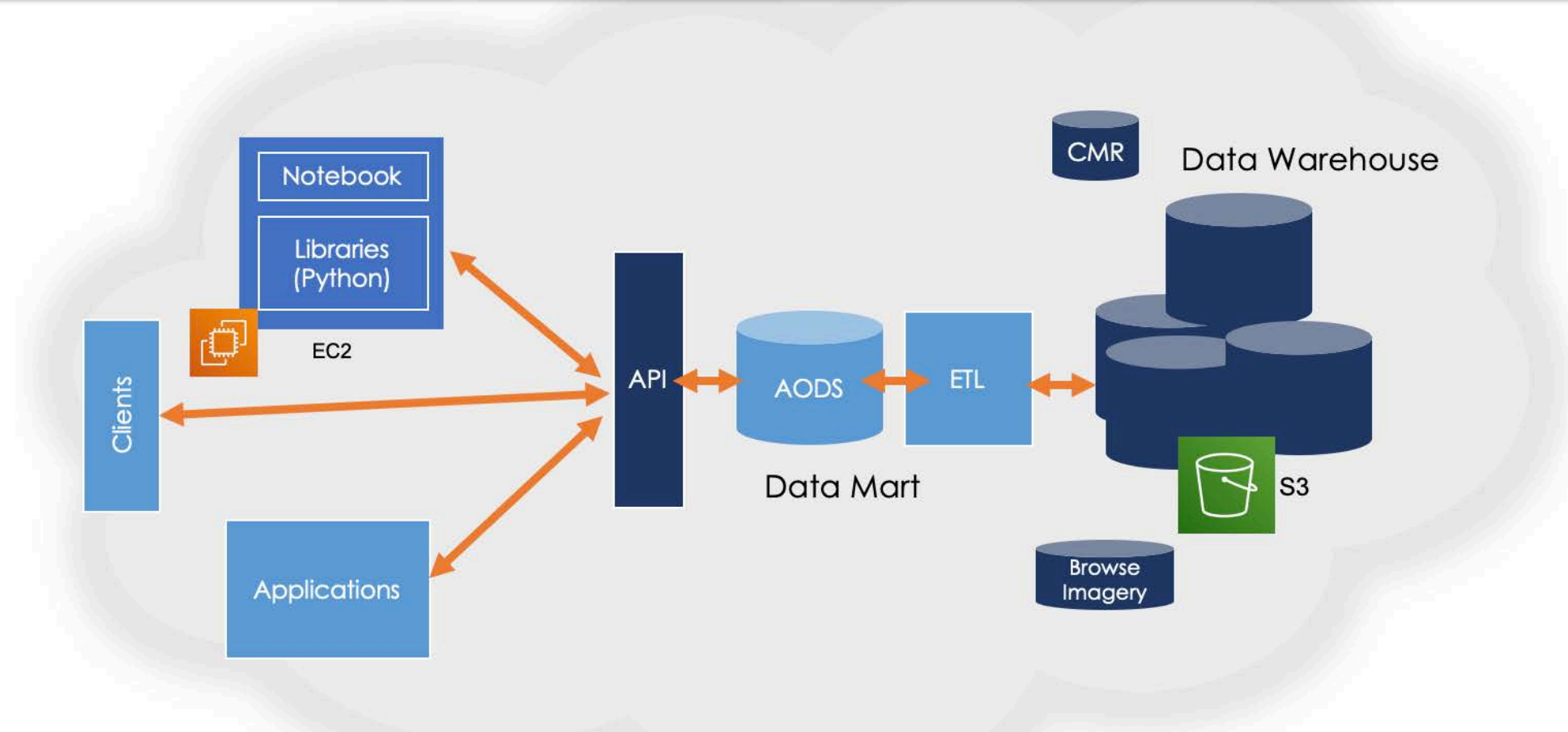
# Mode 3: Big Data/Interactive Analysis

Enable researchers to iteratively query large amounts of data to perform an analysis.



- Akin to data marts but for the Earth science domain
- AODS can be constructed and an analytics software stack spun up to support the researchers' analysis (can be ephemeral)
- Data expeditions - "time-bound" analytics projects on answering specific science questions

# Mode 3: Big Data/Interactive Analysis



**Additional Extract Transform Load (ETL) costs:**
- Data curation to meet data expedition needs
- Preprocessed using consistent and well-known algorithms that cover a spectrum of steps
- Data restructuring to enable data parallel computation analysis

# Multi-Mission Algorithm and Analysis Platform (MAAP)

- The MAAP is a virtual environment dedicated to the unique needs of sharing and processing data from relevant field, airborne and satellite measurements related to ESA and NASA missions
  - Jointly managed by ESA and NASA and accessible to designated ESA and NASA scientists.
  - Initially populated with pre-launch and complimentary data from other projects.
- Science focus is to improve the understanding of global terrestrial carbon dynamics & to support algorithm development
- Addresses a need expressed by the science community to more easily share and process data collected by NASA and ESA activities

# Multi-Mission Algorithm and Analysis Platform (MAAP)

**Mode1 Stack:**

- Eclipse Che

- Docker/Kubernetes

**Mode 2 Stack:**

- HySDS

**Mode 3 Stack:**

- Dask/Zarr

- AWS Athena

# Summary

- Current ARD definitions are incomplete as they do not encompass the data engineering component needed to exploit parallelism and minimize data movement during compute.
- AODS extends ARD to encompass big data analytics to enable data-driven science.
- Adoption of cloud platforms by the science community provides an opportunity to provide different modes of analytics
- Ephemeral AODSs can be created to support data expeditions focused on addressing specific science goals by analyzing large data volumes.
- Need exists to **develop a common set of ETL tools that enforces uniformity and transparency** on preprocessing/data structuring is paramount for scientific analysis

# Thank you!

Contact info: rahul.ramachandran@nasa.gov

Slides Attribution:
Kevin Murphy ESDS Program PE, NASA HQ
Andy Mitchell ESDIS Project Manager, NASA GSFC
IMPACT Team, NASA MSFC
MAAP Team, NASA JPL/NASA MSFC [UAH, Development Seed]