# ARM Data and Computing Capabilities

GIRI PRAKASH

**ARM Data Center, Oak Ridge National Laboratory**
**palanisamyg@ornl.gov**

**BERAC, April 21, 2023**

# Comprehensive Sets of Measurements Deployed in Diverse Climate Regimes



- Background atmospheric state
- Surface energy balance
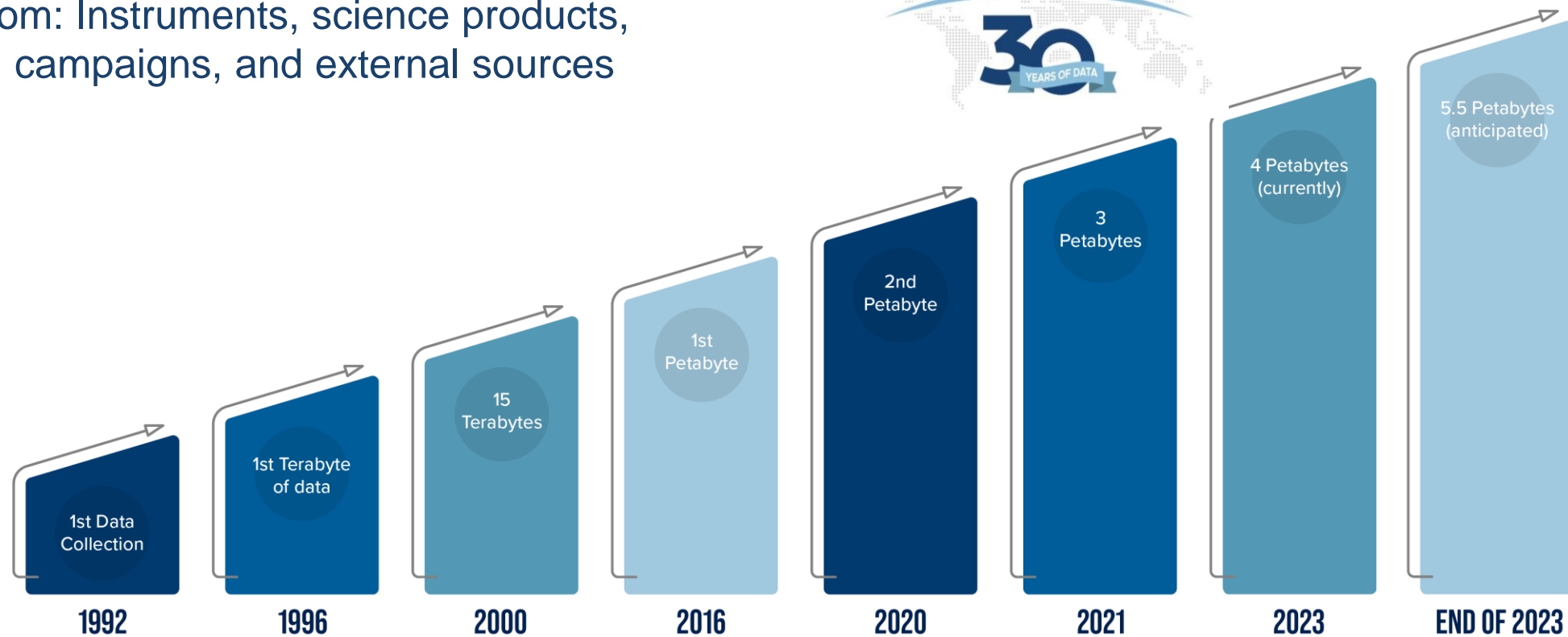- Aerosol and hydrometeor profiles
- Near-surface aerosol properties
- Aerial measurements

Legend:
- Fixed Site
- Past Fixed Site
- Deployments
- Aerial Deployment
- Ship Deployment

# Data and Users At a Glance

► Data From: Instruments, science products, models, campaigns, and external sources

# Data and Users At a Glance
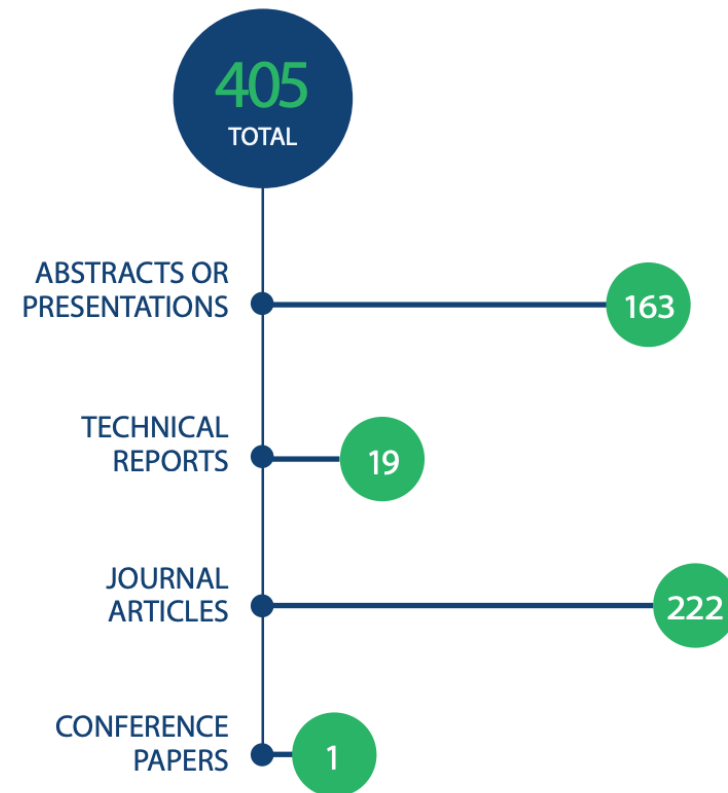
**ARM**

## USERS BY COUNTRY
**37 COUNTRIES**



## PUBLICATIONS USING ARM*

**405** TOTAL

ABSTRACTS OR PRESENTATIONS — 163

TECHNICAL REPORTS — 19
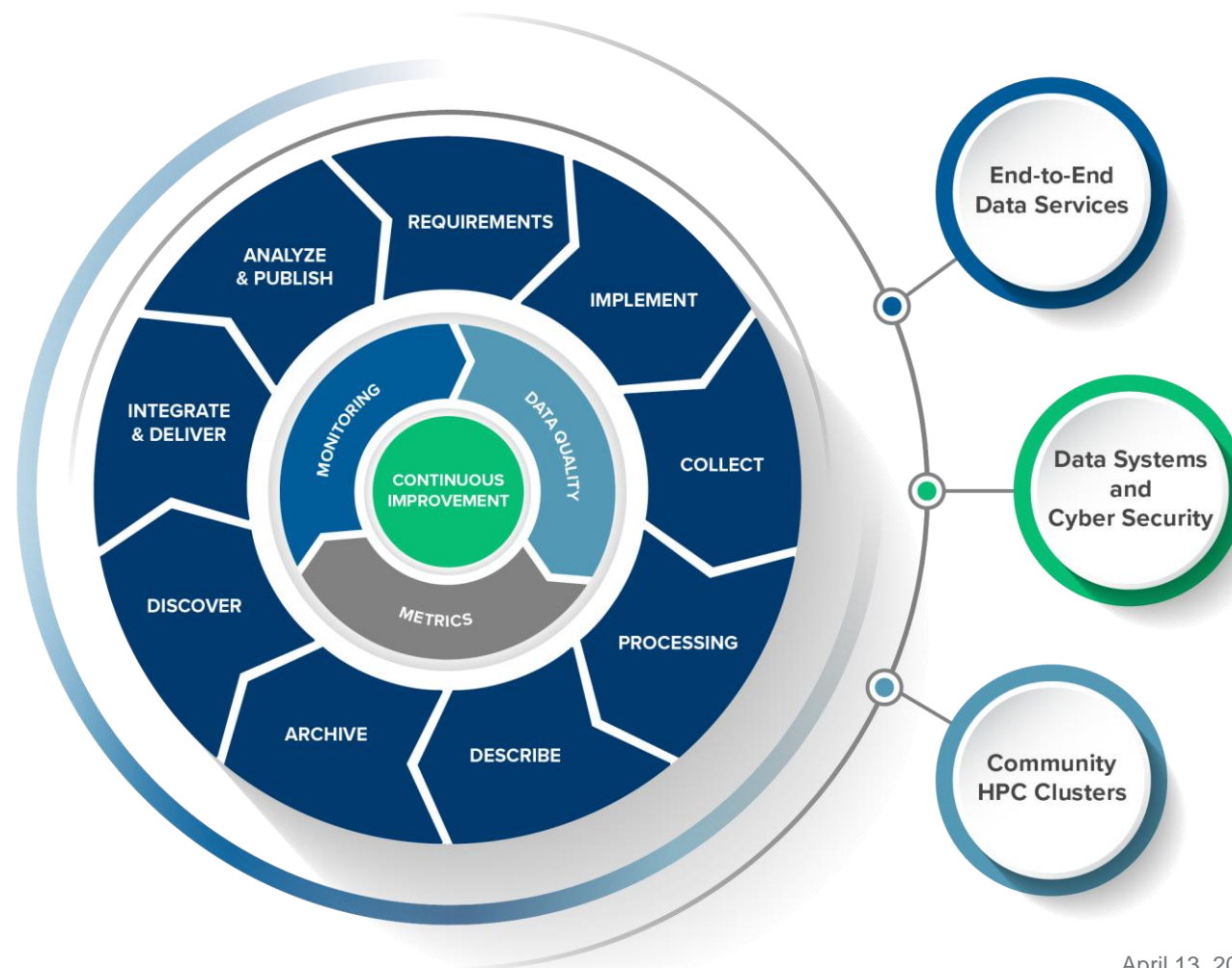
JOURNAL ARTICLES — 222

CONFERENCE PAPERS — 1

*Publication statistics were collected as of December 2022. Journal article numbers will continue to increase over time.
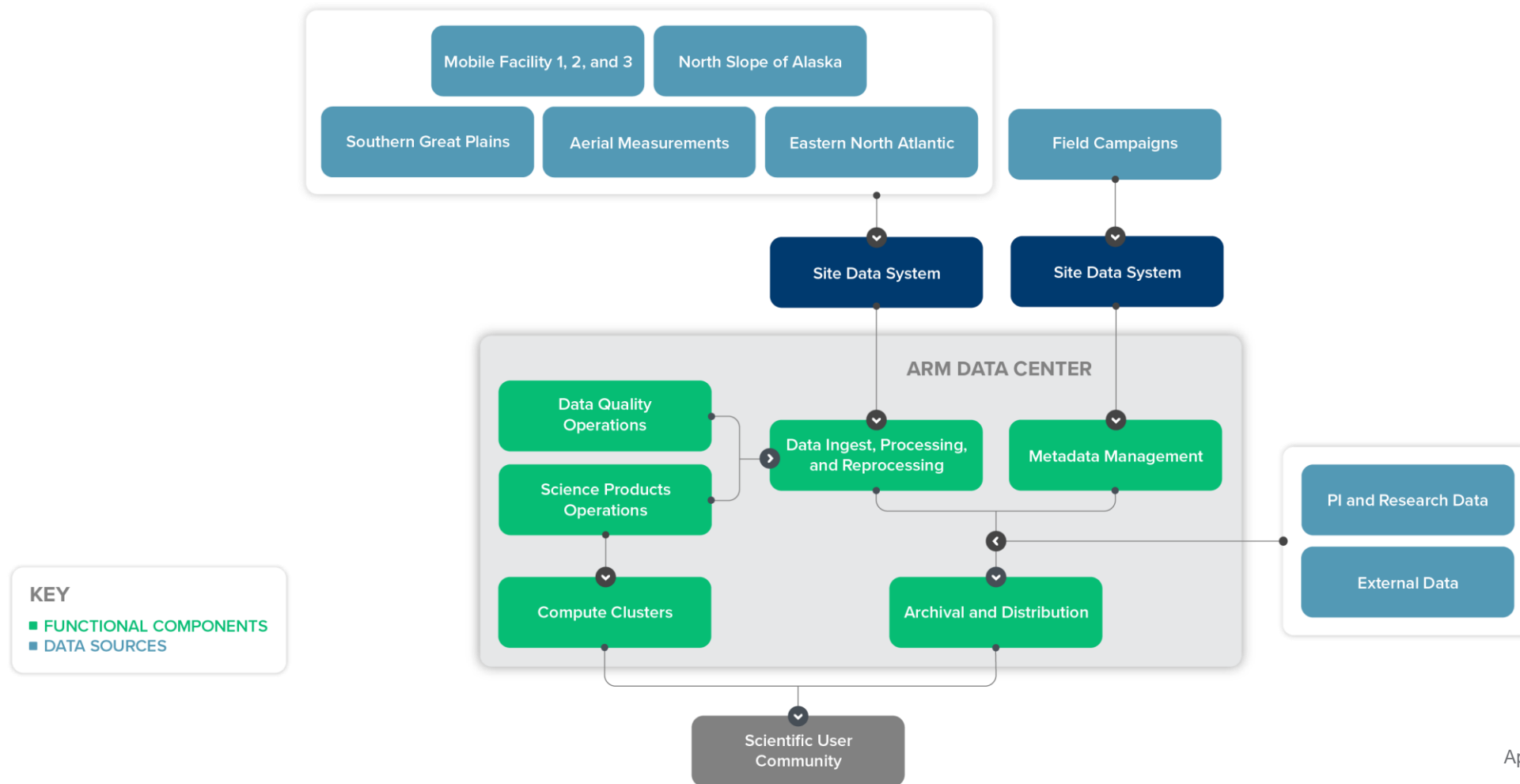
# About The ARM Data Services

**Provides a robust integrated data and computing ecosystem to advance understanding of atmospheric measurements**

► Data flow operations and monitoring
► Advanced data collection systems
► High-performance computing (HPC)
► Comprehensive Data Processing
► Data Interoperability:
  ■ Advanced strategies for utilizing metadata
  ■ Data Discovery
  ■ Data workbench
  ■ FAIR, Standards, and Protocols
► User Management and Citations
► AI-based approach in data management
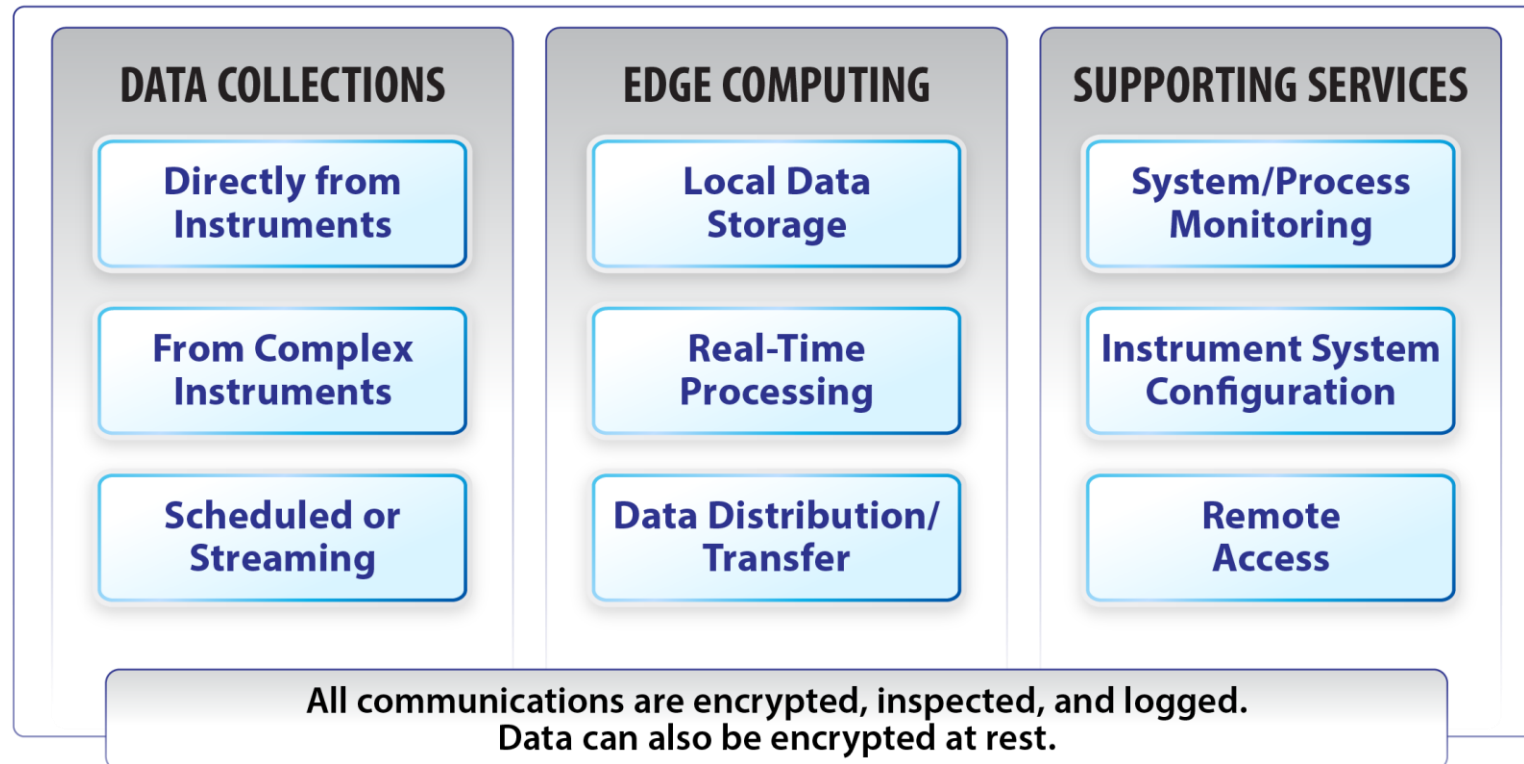
# ARM Data Flow: From Collection to Distribution

► Offers powerful and adaptable infrastructure capabilities to support a wide range of data pipeline requirements, enabling efficient and streamlined processing of data from various sources.

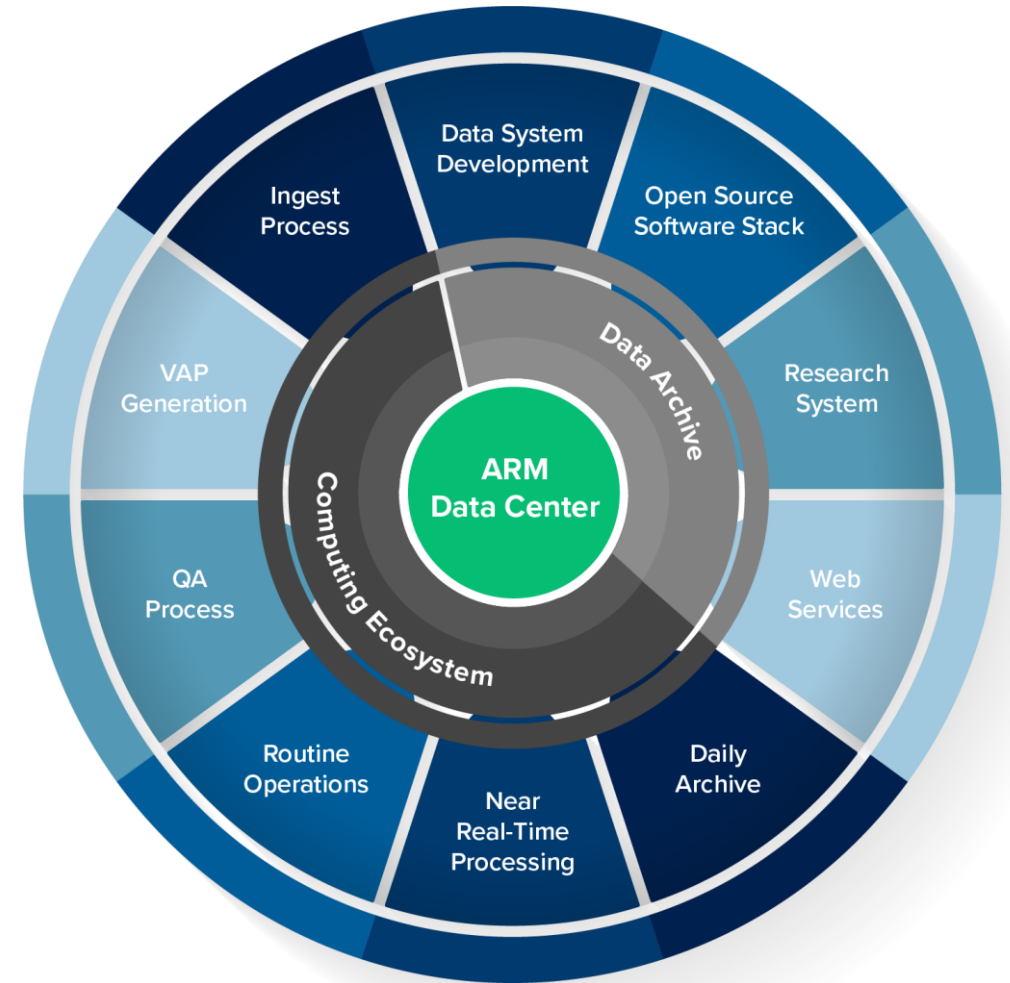# Advanced-Data Collection Systems
# For Next-Gen Sensor Networks

**ARM**

► Scalable data systems with proven hardware and software solutions
► Real-time data access to enable data reduction and edge computing (e.g., Supervised Learning)
► Future development of next-generation instrument computing with Machine Learning

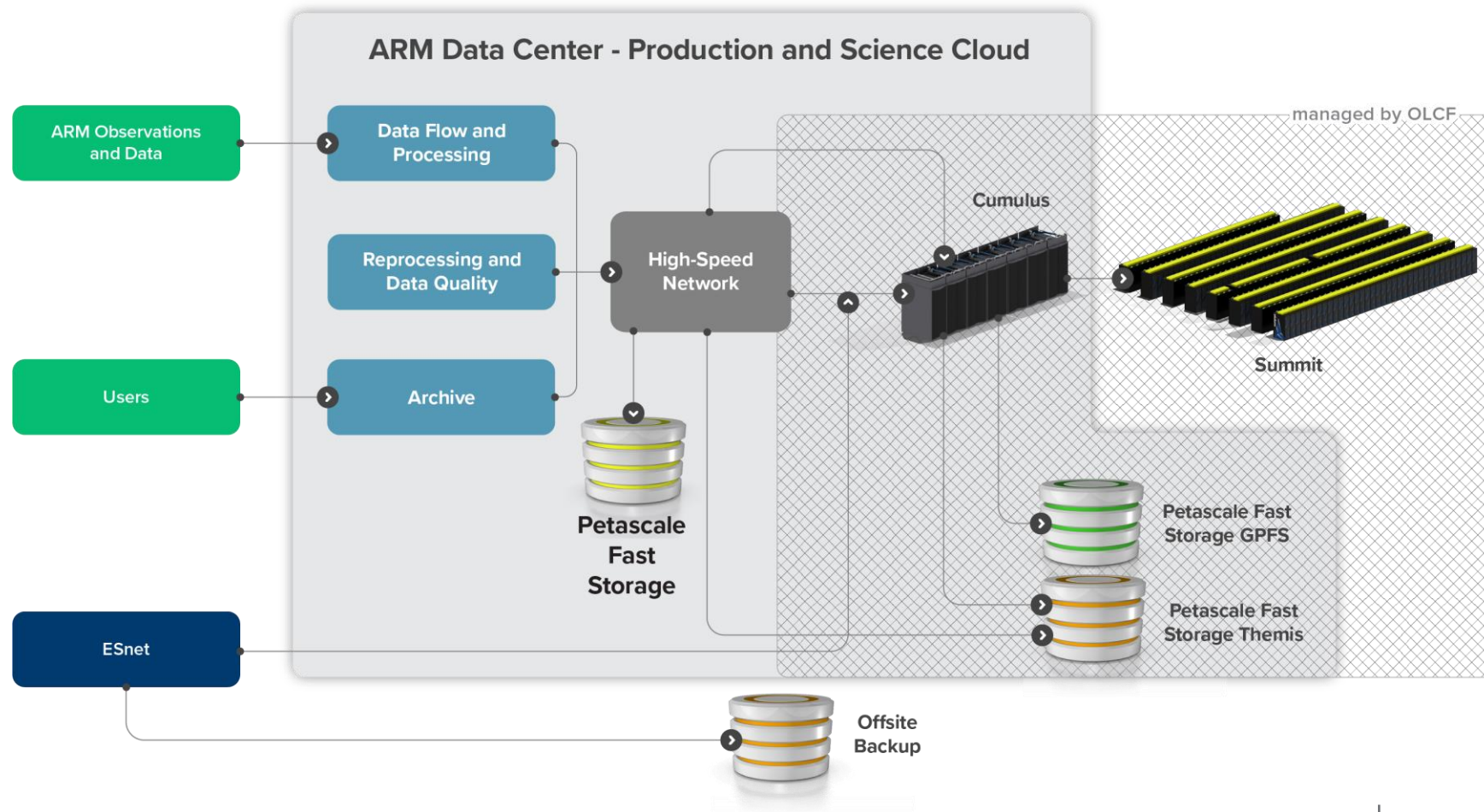# Comprehensive Data Processing Capabilities For Effective Data Management

**ARM**

▶ Efficient flow control for smooth processing

▶ Thorough data quality review for accuracy and reliability

▶ Immediate online access to high-demand data streams

▶ Near-term and long-term reprocessing capabilities

▶ Comprehensive monitoring of thousands of processing data streams.

U.S. DEPARTMENT OF ENERGY

# Computing Capabilities

## ARM Data Center Cyberinfrastructure: Enhancing synergy across DOE computing facilities

► Enables fast parallel processing of:

- Data ingest operations
- Complex ARM datastreams (e.g., Radars, Value Added Products)
- Large Eddy Model Simulations of ARM cases (LASSO)
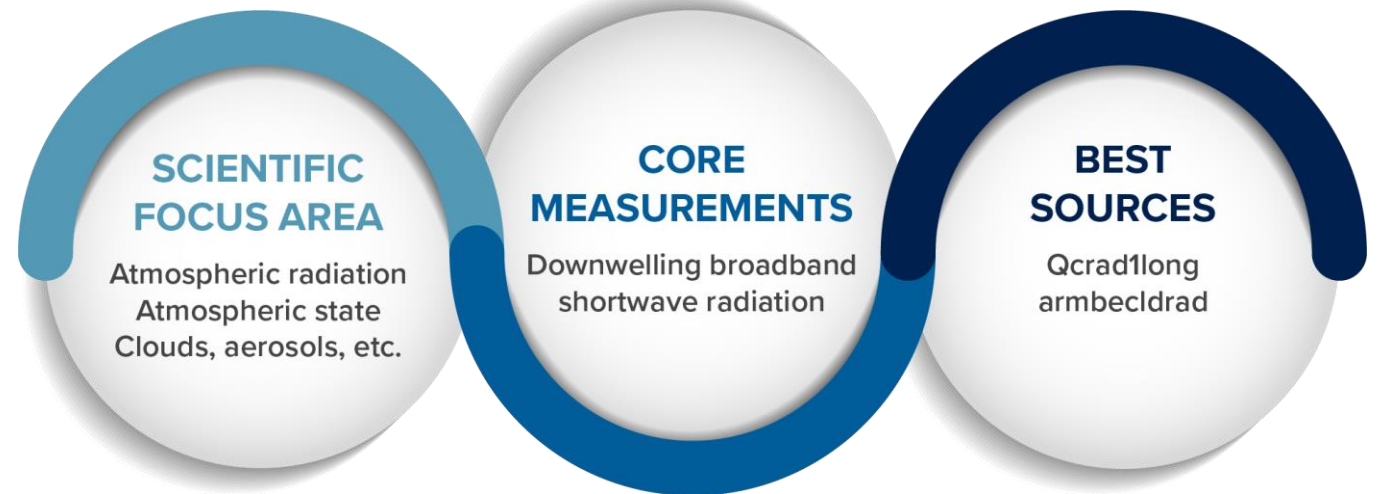- Data analysis using large volumes of ARM data

# Selecting Quality Data Sources: Harnessing the Power of Rich Metadata

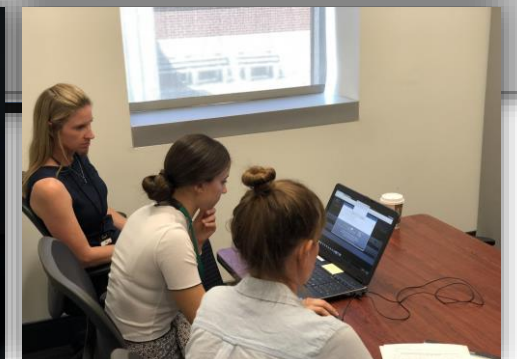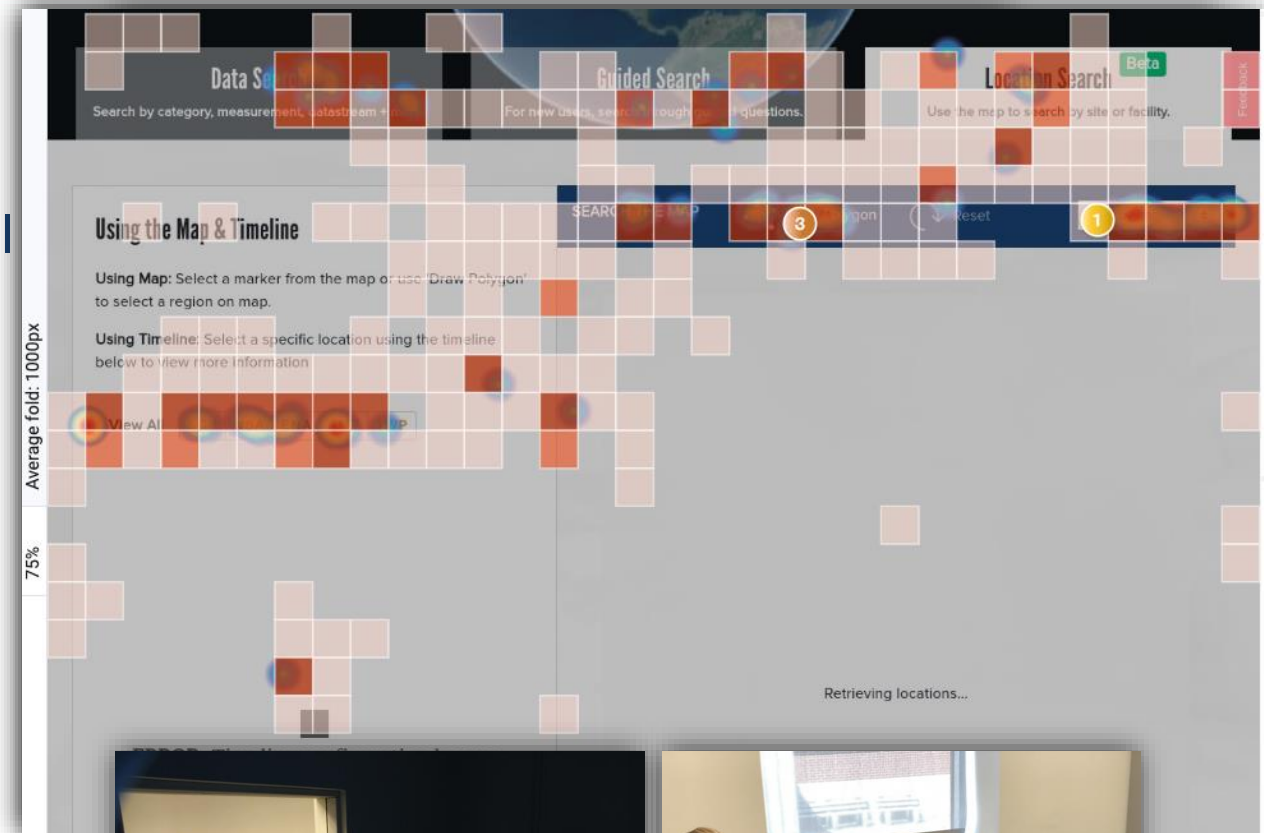**Over 11,000 Data products from 450+ instruments, science products, and model simulations**

▶ Robust metadata workflow system effectively used for operations, discovery, and data interoperability

▶ Recommends best data sources for the core measurements (i.e., Data Epoch)

▶ Semi-automated process includes input from subject matter experts

**SCIENTIFIC FOCUS AREA**
Atmospheric radiation
Atmospheric state
Clouds, aerosols, etc.

**CORE MEASUREMENTS**
Downwelling broadband
shortwave radiation

**BEST SOURCES**
Qcrad1long
armbecldrad

# Advanced Data Discovery: Leveraging Modern Architecture and Search Capabilities
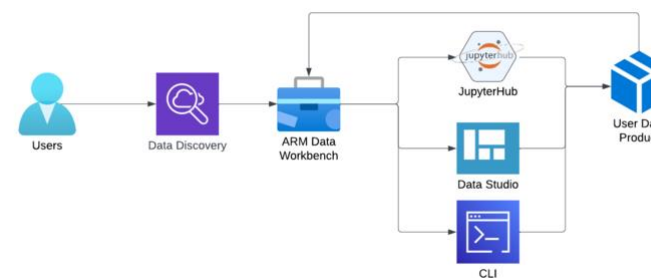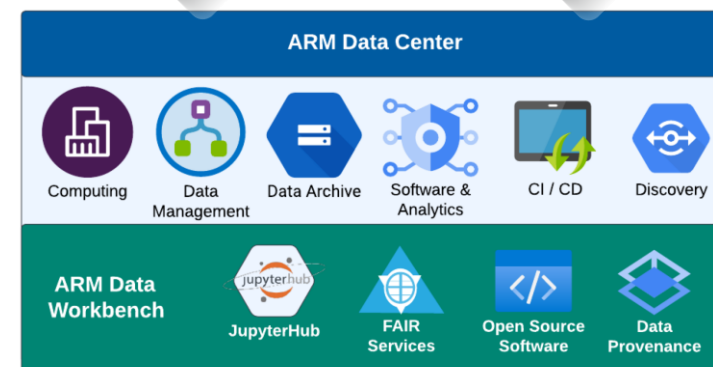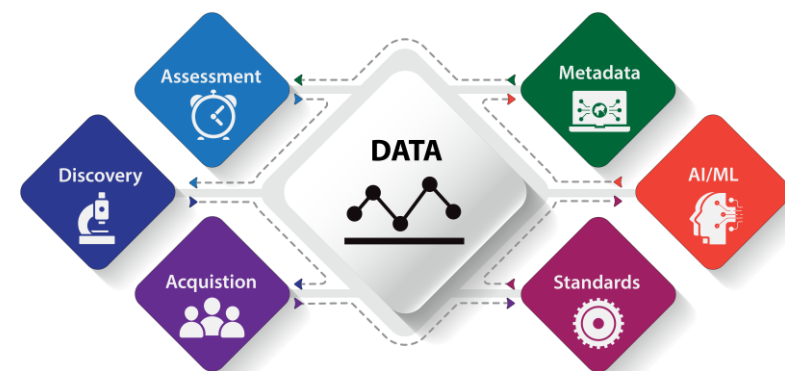
- User-centric design and improvements using modern software architecture with Continuous Integration and Deployment (CI/CD)

- Intelligent search capabilities based on the actual data, guided search based on user experience

- Recommendations, data tagging based on epochs or golden periods

- Near real-time access via secured webservices (API access)

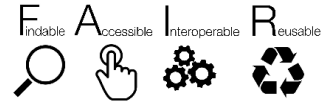- Customized interface for ARM high-resolution model simulations

# Data Workbench: Enabling Data Interoperability

▶ Aims to achieve transformative knowledge discovery by providing modular computing, data, and software capabilities

▶ Facilitate easier interaction with ARM data and enable interoperability with other data sources

■ Provide a collaborative and dynamic computation environment for data analysis, scientific computing, and machine learning (e.g., JupyterHub)

■ Facilitate data access to external datasets (e.g., weather radar, satellite, model data, etc.)

▶ Enable FAIR-based access to ARM data and computing for initiatives such as AI4ESP

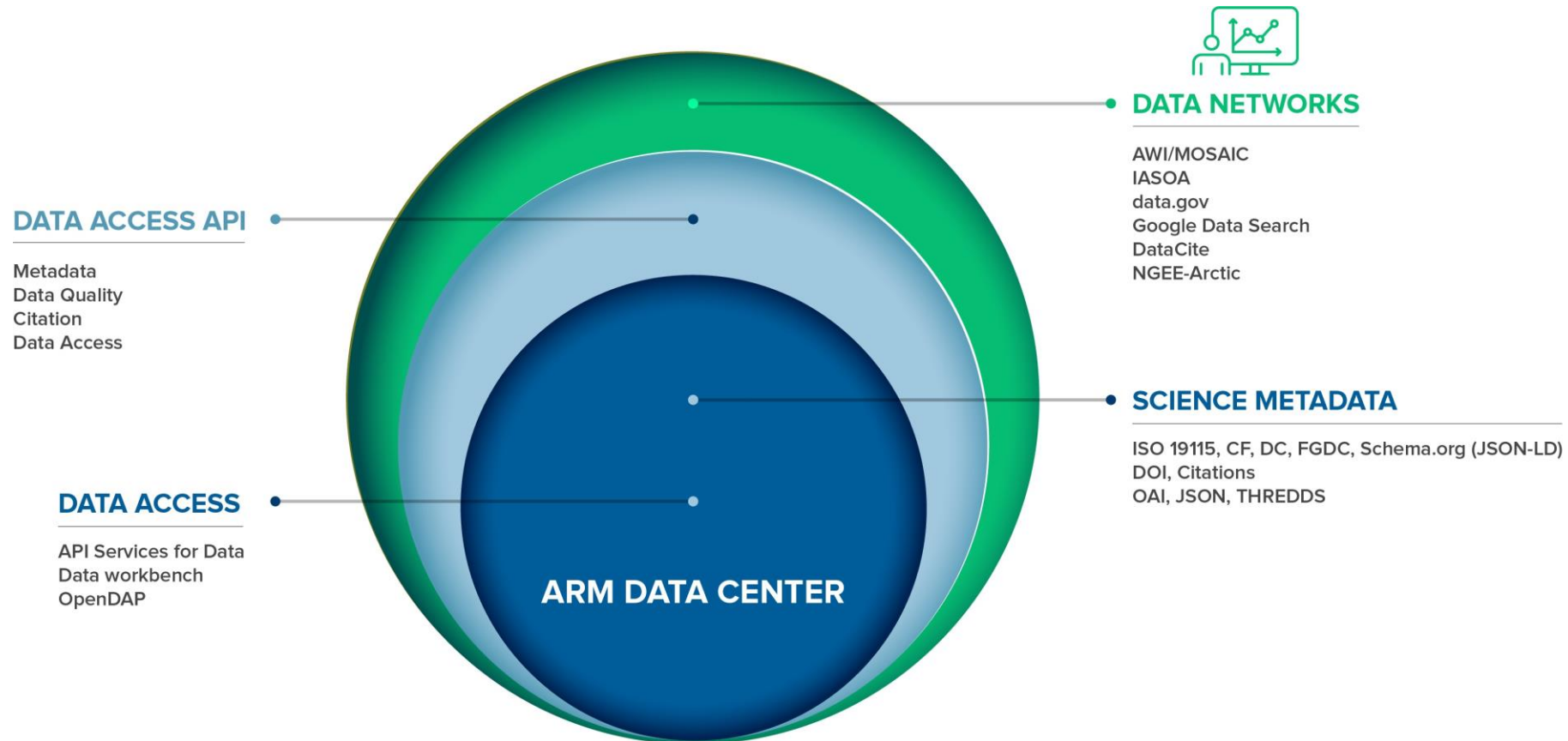# FAIRness Assessment and Community Engagement

► Review of data management capabilities and obtaining certifications

► Continuous collaboration with broader data networks

► Active contribution to national and international working groups

# Putting FAIR Principles into Practice: Standards and Protocols in Data Interoperability



**DATA NETWORKS**

AWI/MOSAIC
IASOA
data.gov
Google Data Search
DataCite
NGEE-Arctic

**DATA ACCESS API**

Metadata
Data Quality
Citation
Data Access

**DATA ACCESS**

API Services for Data
Data workbench
OpenDAP

**ARM DATA CENTER**

**SCIENCE METADATA**

ISO 19115, CF, DC, FGDC, Schema.org (JSON-LD)
DOI, Citations
OAI, JSON, THREDDS

# Expanding the Reach of ARM Data: Data sharing Examples

▶ Data Access:
  - ■ Ensure the latest version of data are available for users
  - ■ Data endpoints are provided in the metadata
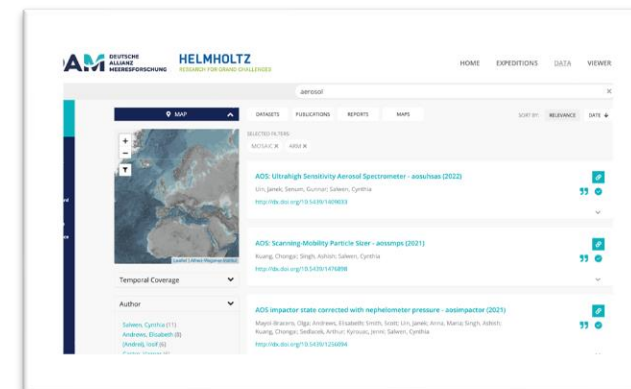  - ■ Direct access via API-based services (live data service, Globus, OpenDAP, and JupyterHub)

▶ Provide access to data quality, plots, and other ancillary details

▶ Options for users to get notified of any data quality changes or new data versions

▶ Interoperability:
  - ■ ARM Data is currently discoverable in partner portals
  - ■ External data are shared through the ARM discovery interface
  - ■ Currently in discussion with BER data centers such as EMSL, ESGF, and Ameriflux
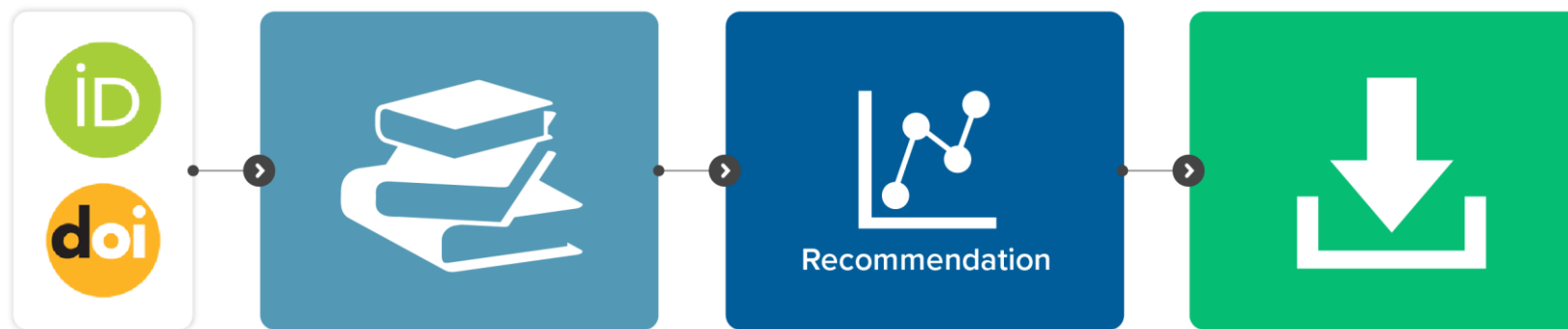
MOSAIC Data Portal



NGEE-Arctic Data Portal
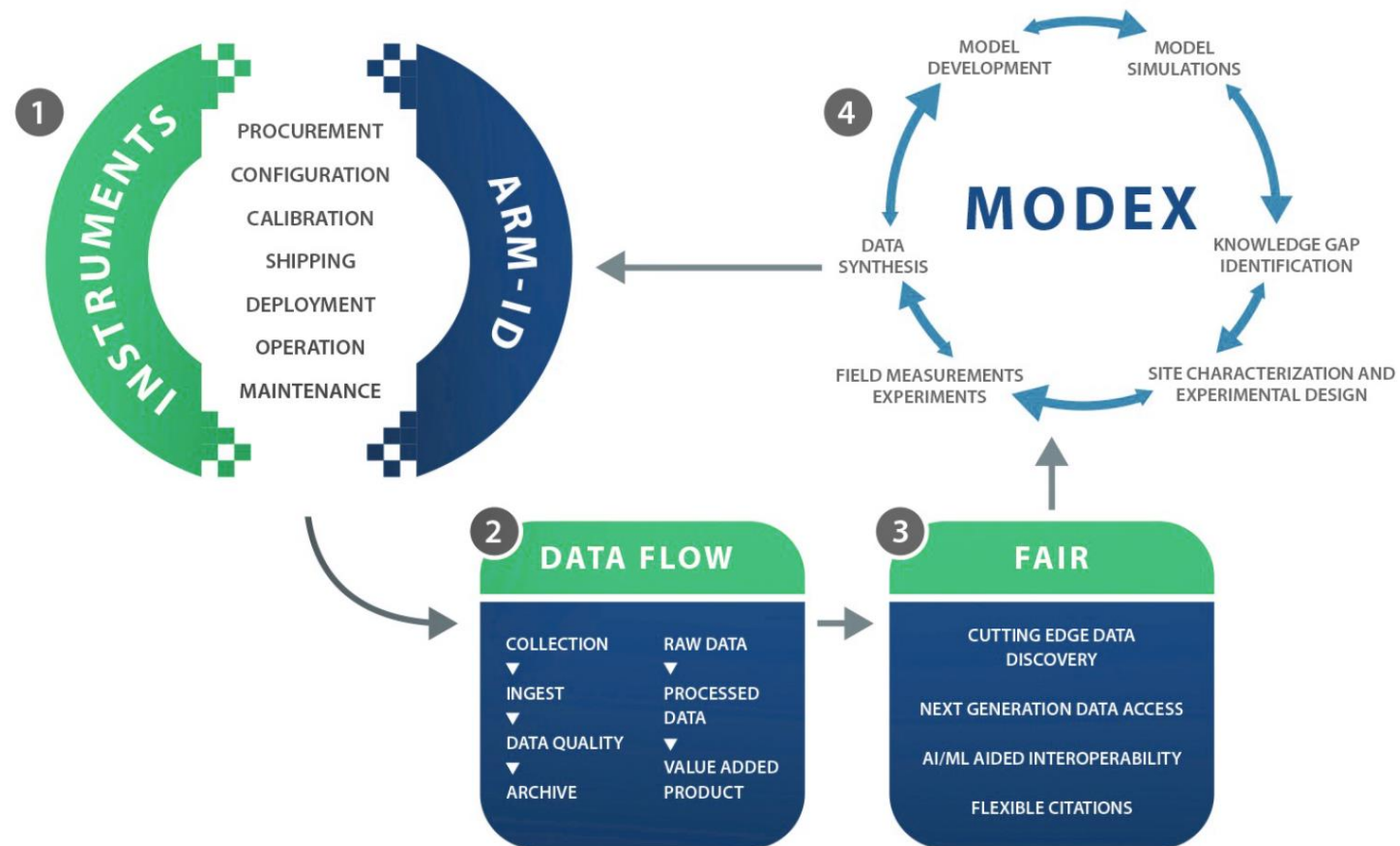
# Maximizing User Management and Data Citations

► Integrating ORCID with other user metrics improves the program's ability to manage the quality of user details and metrics preparation

► Opportunities exist to improve user experience using AI/ML techniques
  ■ Discover relationships between ORCID identifiers, users, publications, data, metadata etc. Then use these relationships to improve the user experience with finding and using ARM Data

# Looking Ahead: Unlocking the Power of Data. The Role of AI in Enhancing Observational Data Centers

▶ Enabling interdisciplinary research through modernization of data pipelines from collection to distribution using AI-based approaches

▶ Near real-time data analysis and data collection configurations using edge computing

▶ Developing and extending community-based standards between data repositories and AI models

▶ Data tagging to identify benchmarking/training datasets

# Summary

## ARM Data Services

▶ Provides robust data collection, processing, archival, and distribution capabilities

▶ Enables unified data, computing, and software ecosystem for scientists and facility operations

▶ Empowers data interoperability with broader scientific data networks by putting FAIR principles into practice.