# ACME – Accelerated Climate Modeling For Energy

Marcia Branstetter

Katherine Evans

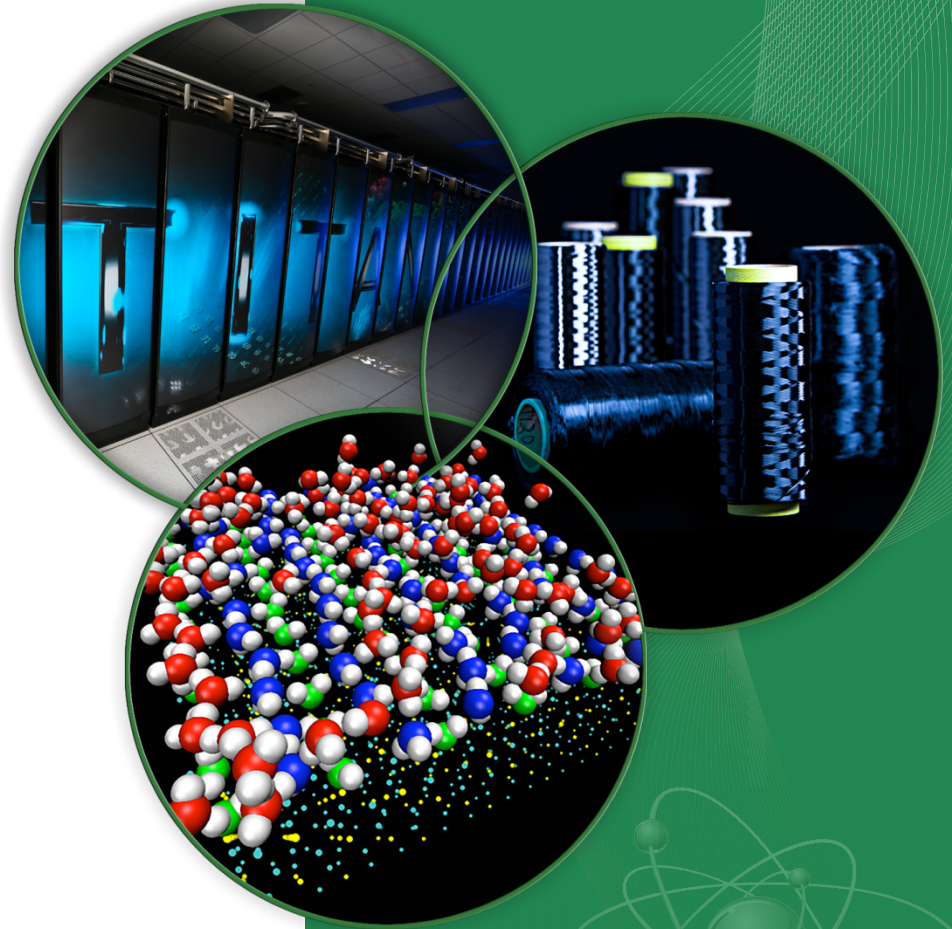**John Harney**

Benjamin Mayer

Daniel Ricciutto

Galen Shipman

Brian Smith

Chad Steed

Peter Thornton

# What is ACME?

- *Multi-institutional, multi-disciplined, BER funded climate science effort*

- *Proposal accepted June 1 for 3 years*

- *Mission:*

   *Build and test a next-generation earth modeling system that can be run on future generations of exascale computing systems at Office of Science computing facilities*

# What is ACME?

- ***Science Drivers:***
  - How do the hydrological cycle and water resources interact with the climate system on local to global scales?
  - How do biogeochemical cycles interact with global climate change?
  - How do rapid changes in cryospheric systems interact with the climate system?

- ***Science Needs:***
  - Support of wide range of model runs and workflow types
  - Capture of model runs, settings and data during development
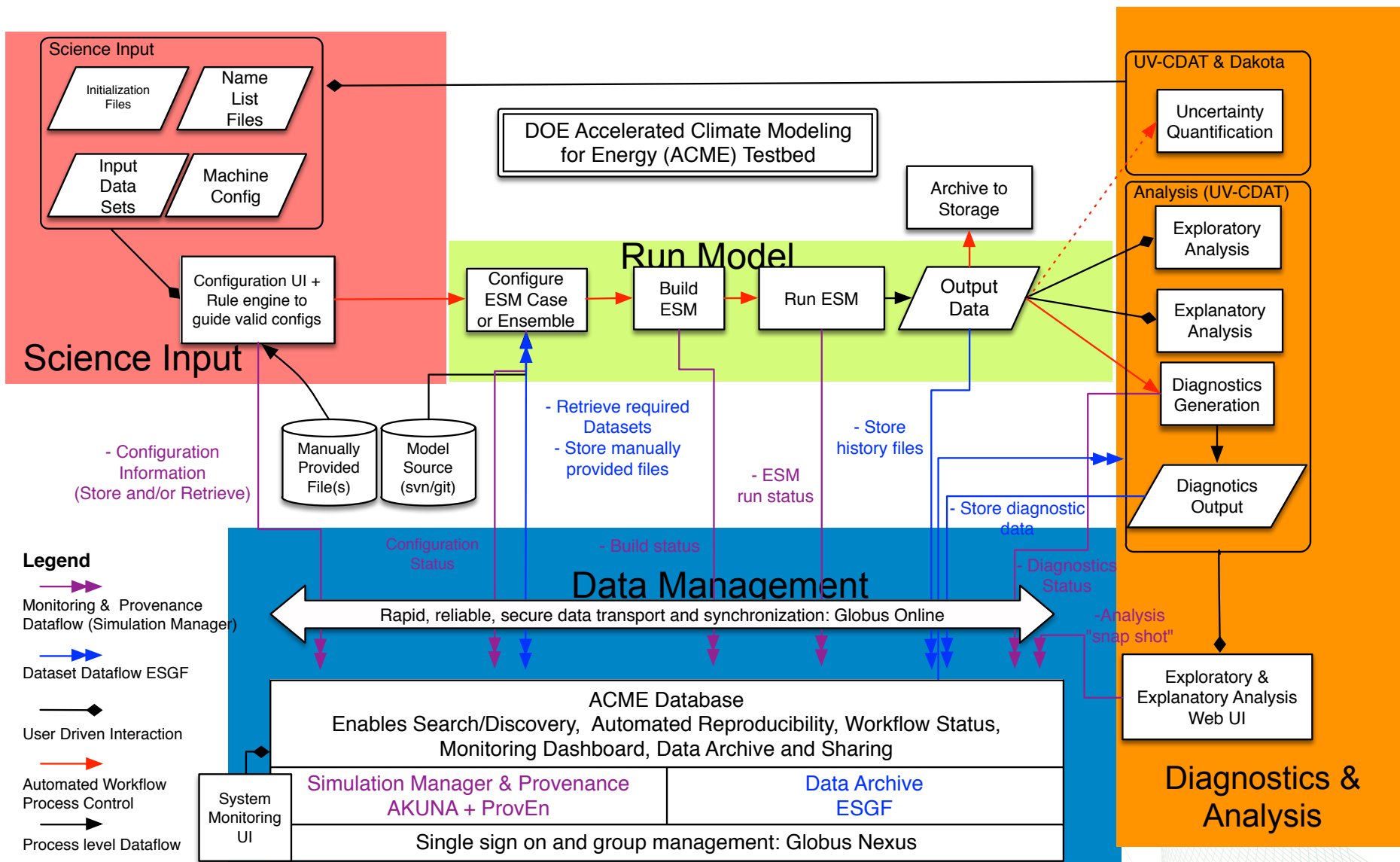  - Quickly evaluate and validate model behavior

**OAK RIDGE** | OAK RIDGE
National Laboratory | LEADERSHIP
COMPUTING FACILITY

# ACME Solutions

| Science Needs | Solution |
|---|---|
| **Support of wide range of model runs and workflow types across the project** | Define use cases: 1) Model Development (DEV); (2) Exploratory Model Runs (EXP); (3) Production Model Runs (PROD)<br><br>Hi-res and low-res |
| **Capture and record suites of runs and their settings during model development** | Automated provenance and archiving throughout the life cycle of the model |
| **Quickly evaluate coupled model behavior** | Advanced diagnostics of the coupled system within one software system |

OAK RIDGE
National Laboratory | OAK RIDGE
LEADERSHIP
COMPUTING FACILITY

# Major challenges for the ACME end-to-end system

| Challenges | Description |
|---|---|
| **Installation** | Software must adapted to multiple hardware platforms and operating systems located throughout the ACME |
| **Heterogeneous Data Sets** | The same infrastructure must also allow scientists to access and compare data sets from multiple sources, including from observational satellite and instrument sources |
| **Analysis, Diagnostics, and Visualizations** | The generation of new and improved analysis, diagnostics, and visualization techniques for the better model development and evaluation |
| **Server-side and In Situ Computing** | Server-side and in situ computation is necessary as the increase in data size and complexity of algorithms lead to data-intensive, compute-intensive challenges for ACME diagnostics, UQ, analysis, model metrics, and visualization |

OAK RIDGE
National Laboratory

OAK RIDGE
LEADERSHIP
COMPUTING FACILITY

# End-to-end workflow

# Science Input & Run Model – Streamlining CESM

- Model development tasks (configure, interpolation, run, build) are tedious and time-consuming

- ACME partners developers and scientists to define and iteratively implement in (semi-)automated end-to-end science workflow. Ideally reduces:
  - Scientist effort to setup/run an experiment
  - Time to solution by automating running of CESM
  - Queue wait times
  - Delays of human intervention

- Implement so automation "just happens" or can be used as helper scripts by hand

# Science Input & Run Model – Workflow automation

- Workflow progress on compute resources is reported back by each workflow component via AMQP messaging broker
  - Leverage existing DOE technology to consume messages (e.g. Akuna from PNNL)

- Create central place to see progress of simulations
  - Current progress (e.g. 23 / 50 years)
  - Current status of data (simulation progress, HPSS, etc)

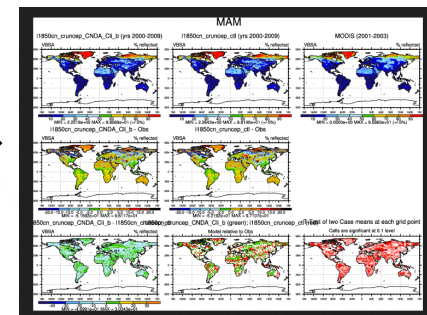- Direct links to archives and diags framework

# Diagnostics and Analysis

- Diagnostics
  - Used to quickly evaluate models and validate their results
  - Traditional CESM NCL scripts from NCAR produce static HTML and plots (gif, jpg, etc)
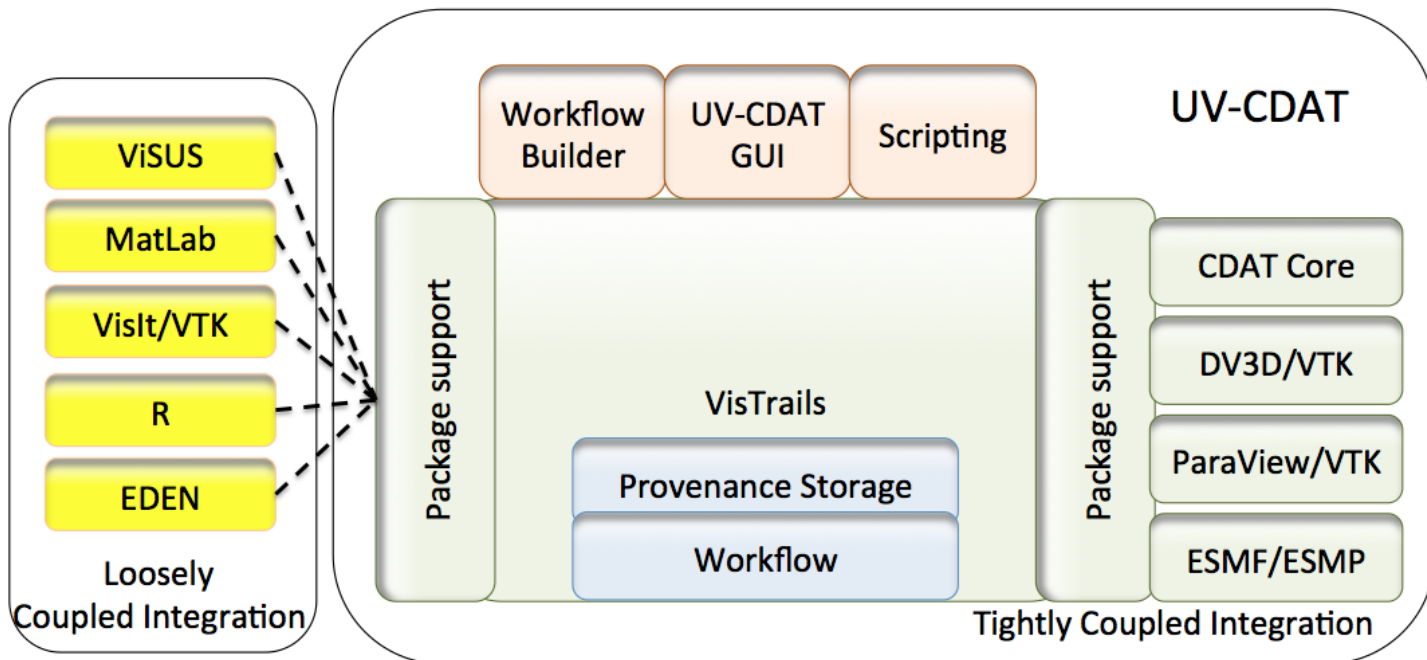
**Diagnostics Home**
**(Plot types of realms)**

**Variable List**

**Plot**



  - Need a more scalable, dynamic and intuitive diags package so results can be ascertained and validated quickly

OAK RIDGE National Laboratory | OAK RIDGE LEADERSHIP COMPUTING FACILITY

# Diagnostics and Analysis - UVCDAT

- ## What is UV-CDAT:
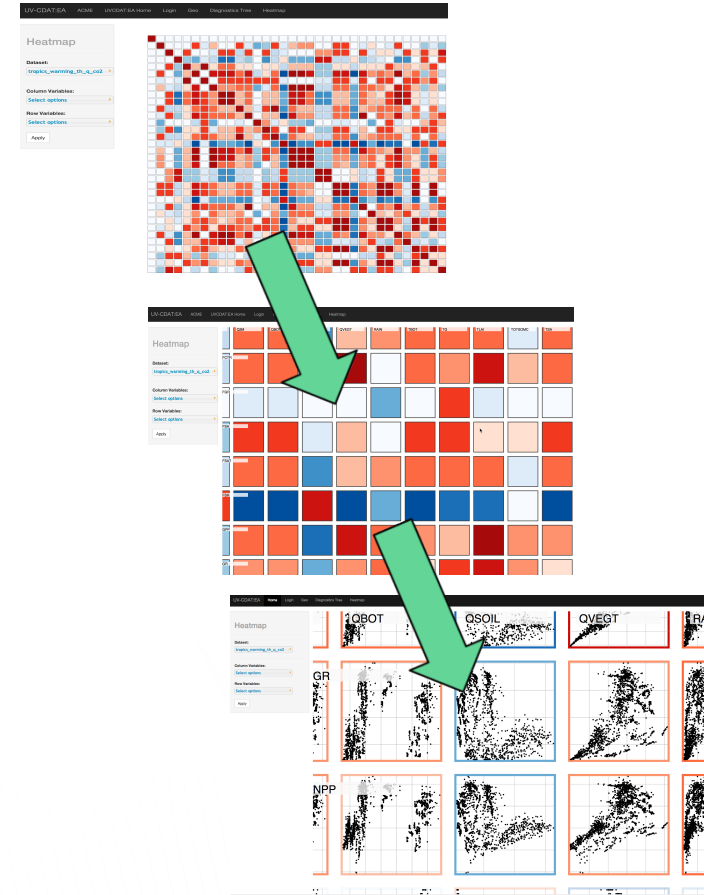  - ❑ A seamless environment for open-source data analysis and visualization packages



- ## What is UV-CDAT purpose:
  - ❑ Bring together robust tools for climate data processing and reproducibility
  - ❑ Integration heterogeneous data sources (e.g., simulations, observation, re-analysis)
  - ❑ Local and remote data access and visualization

OAK RIDGE
National Laboratory | OAK RIDGE LEADERSHIP COMPUTING FACILITY

# Diagnostics and Analysis – Web-based Visual Analytics

- ## Attractive alternative to standalone, thick client

  - Democratizes advanced visual analytics capacities

  - Greater efficiencies through distributed architecture

  - Runs in a web browser (D3, django, jquery)

  - Supports social collaboration

  - Compatibility with UVCDAT

  - Multiple views developed so far:
    - Interactive geospatial-temporal view
    - Diagnostics Tree Viewer
    - Heatmap variable correlation

**E.g. Web-based Heatmap Correlation**

**OAK RIDGE** National Laboratory | OAK RIDGE LEADERSHIP COMPUTING FACILITY

# Diagnostics and Analysis – Examples

OAK RIDGE | OAK RIDGE LEADERSHIP COMPUTING FACILITY
National Laboratory
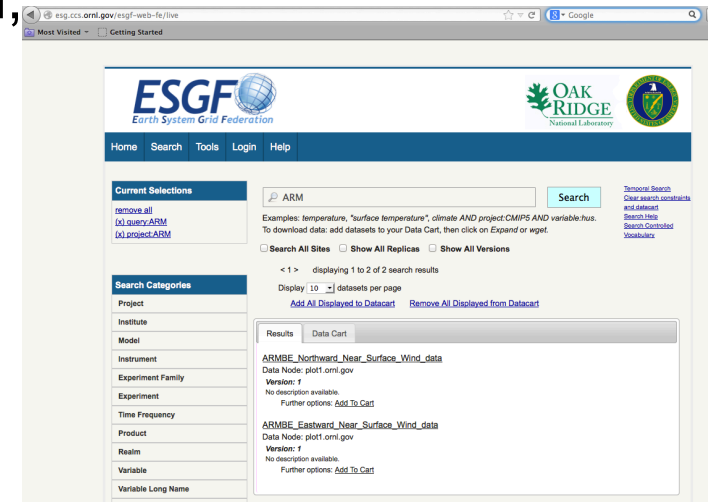
# Data Management

- Output from completed models enter the ACME data management life cycle which supports:
  - Data Collection and Publication
  - Archiving and stewardship
  - Search and Discovery
  - Provenance Capture (i.e. Data/Algorithm reproducibility)
  - Fast/secure acquisition of data (downloads, transfers, etc)

**Data Lifecycle**

Plan → Collect → Assure → Describe → Preserve → Discover → Integrate → Analyze → Plan

OAK RIDGE National Laboratory | OAK RIDGE LEADERSHIP COMPUTING FACILITY

# Data Management – Data Services
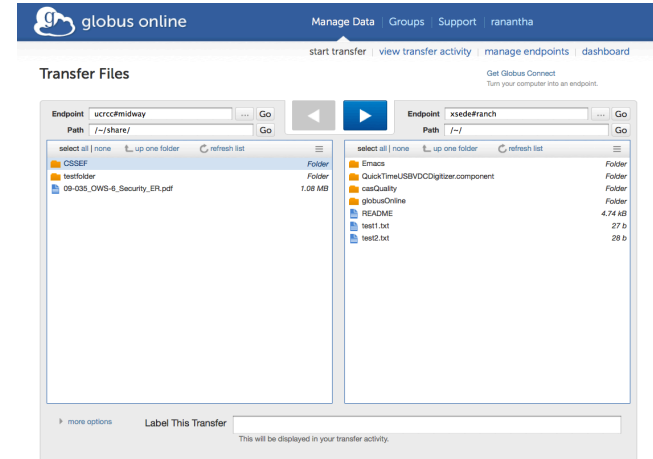
***Earth System Grid Federation (ESGF)***

- Distributed enterprise system that deploys software infrastructure for the management, collaboration, dissemination, and analysis of climate model output and observational data
- Over 30k published federation-wide datasets (over 1PB) in over 20 climate institutions worldwide
- Peer-to-Peer design for increased scalability and fault tolerance
- Solr search engine designed for optimal metadata cataloging and dataset retrieval
- Secure file access and group-based authorization services

OAK RIDGE National Laboratory | OAK RIDGE LEADERSHIP COMPUTING FACILITY
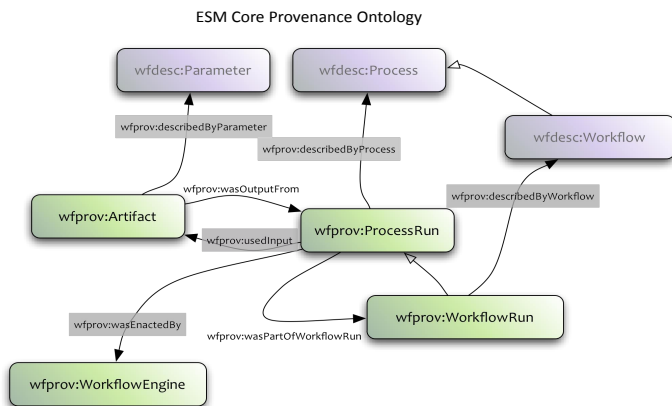
# Data Management – Data Services

## *Globus Online*

- Reliable, secure, high-performance file transfer and synchronization
- "Fire-and-forget" transfers
- Automatic fault recovery
- Fully integrated with ESGF services
- Key technology in transferring large datasets (e.g. high resolution)
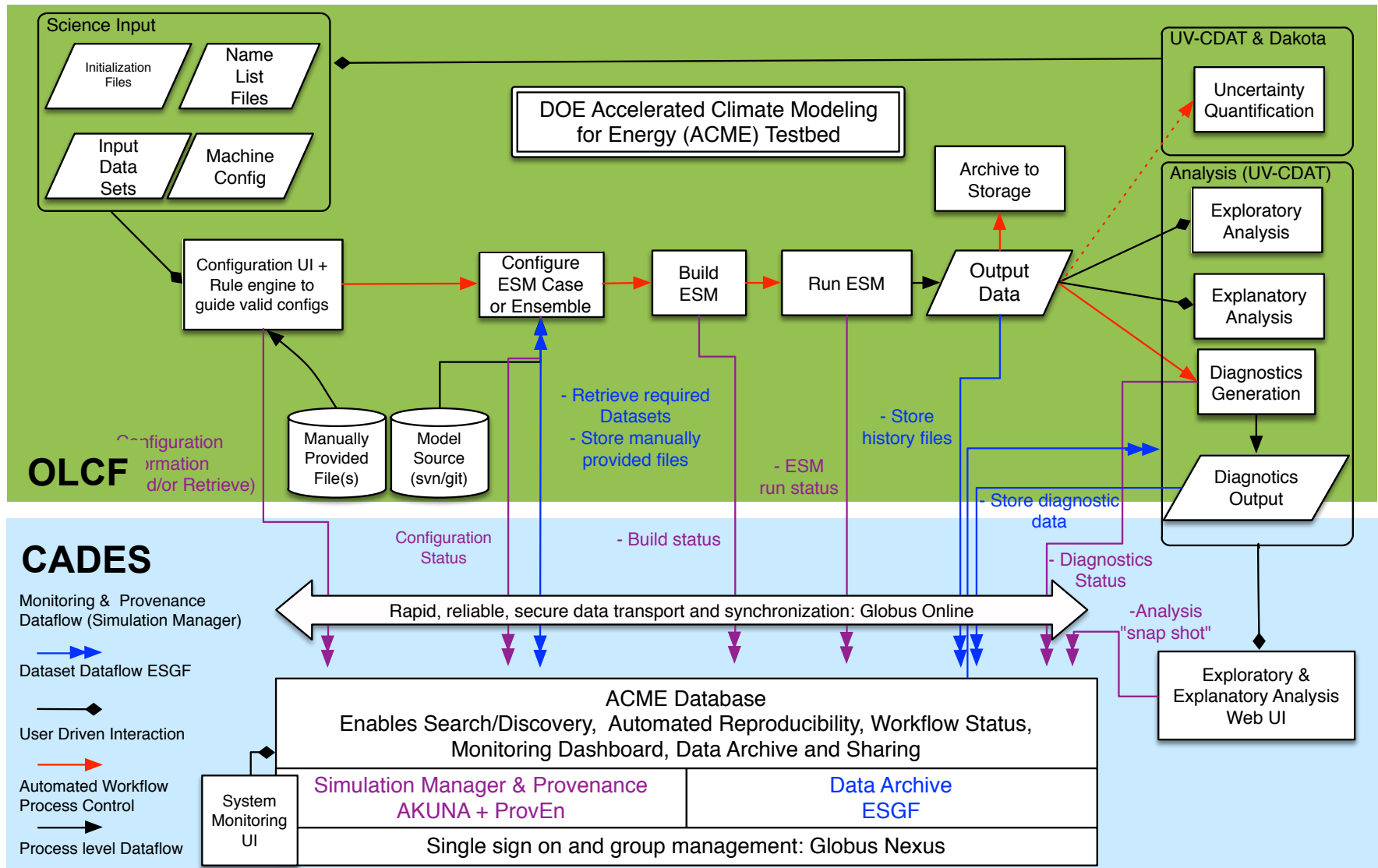




ESM Core Provenance Ontology

## *PROVEN*

- Capture historical information from any native source necessary to describe the origin of the dataset
- Store this information in a cross-referenced form through the use of internationally recognized standards W3C, Dublin Core, CF
- Use this cross-referenced information to provide finished products to different kinds of consumers
- Will be integrated with ESGF in future release

OAK RIDGE National Laboratory | OAK RIDGE LEADERSHIP COMPUTING FACILITY

# End-to-end workflow and data infrastructure architecture

# CADES
## Compute & Data Environment for Science

CADES provides core compute and data services required by major science facilities, large projects, small teams, and single principal investigators
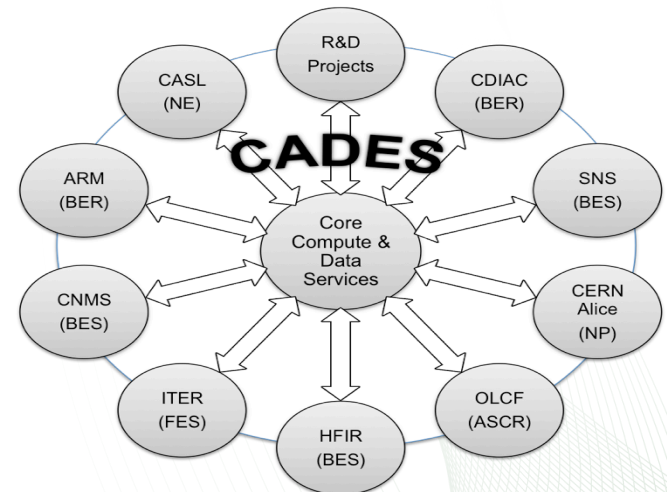
CADES is a cross-cutting center: it shares both data infrastructure and compute & data science expertise with and among many projects

A rich set of flexibly composable services coupled with experts in data science partnering with domain scientists on their challenges

# Experiences and Lessons Learned (so far)

- Communication and Relationships (scientists-developers, lab-lab, component-component, etc)

- Unstructured grids and interpolation

- Messaging across heterogeneous security enclaves

- High Resolution data issues
  - Storage (HPSS)
  - Movement (Globus and security policy)

**OAK RIDGE** National Laboratory | OAK RIDGE LEADERSHIP COMPUTING FACILITY
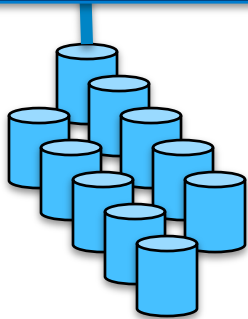
# Discussion

Questions?

# CADES integration with OLCF

- The CADES cloud infrastructure is a new resource that has been added to the existing services of CADES

- This platform will enable new CADES services

ESNet, Internet2, XSEDE

OpenStack IaaS & Parallel Compute/Data Environment

Scalable I/O Backplane

Shared home

### Scalable storage
Lustre, Cinder, Swift, RADOS

### High performance
Data Mining, Fusion & Analytics
Modeling & Simulation

### Utility compute
Databases
Web Servers
Workflow Engines
Other Persistent Services

### Integration with extreme scale
**27 Petaflops** – Titan
**200 Petabytes** – HPSS
**30 Petabytes** - Lustre

**OAK RIDGE** National Laboratory | OAK RIDGE LEADERSHIP COMPUTING FACILITY

# Testing and execution framework tailored for ACME

| Features | Description | Impact |
| --- | --- | --- |
| **Infrastructure** | Creates a flexible, extensible infrastructure for future ACME efforts and related DOE projects, automates laborious, repetitive simulation data tasks to improve productivity | Heightens productivity and user experience |
| **Data Sharing** | Supports broad data sharing within ACME project teams and with scientific collaborations; including NGEE, ARM, CDIAC, etc. | Accelerates model development and result dissemination |
| **Provenance** | Enables reproducibility, archiving and reuse of high-volume simulation data, provenance captures set up, execution and analysis details coupled with standard metadata creation, annotation, and forums for group discussions and sharing of any part of a workflow | Increases reproducibility, productivity and credibility of collaboration |
| **Model set up and execution** | Rule based support for model setup, specialized collaborative portal with a checklist for approvers of new model setups before job launching, Links to User interface and infrastructure for job submission | Enables new users to be effective quickly, allows control over model set ups in distributed, collaborative teams |
| **User Interface** | Specialized software when needed to enable web job submission, running, monitoring, and debugging capabilities on several HPC centers | One stop shop to all needed capabilities, increases productivity, reproducibility |
| **Ensemble and Automated Runs** | Job launching interface for submission of hundreds of production runs and enabling specialized monitoring of multiple ensemble, automated runs | Increases productivity |