# DataNet Federation Consortium & iRODS

# Interoperability
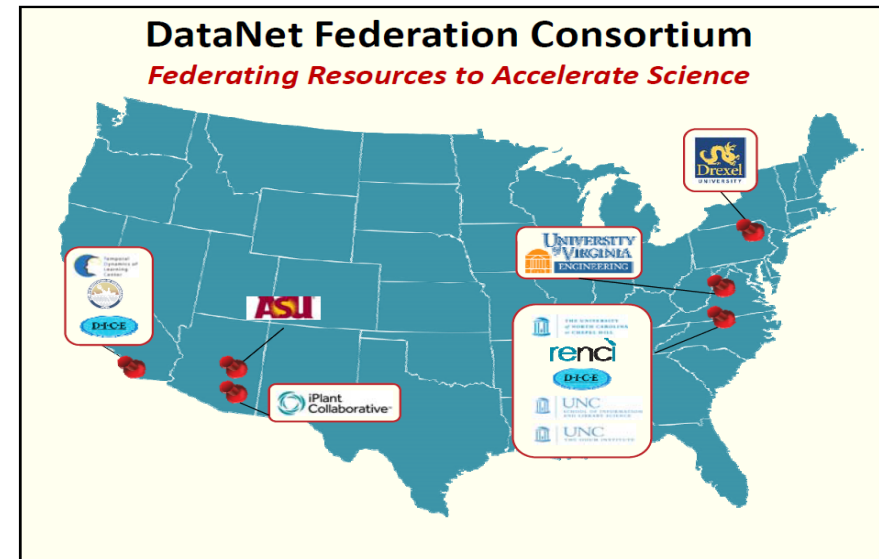
Charles Schmitt
CTO & Director of Informatics
RENCI

# DataNet Federation Consortium (DFC)

**NSF Sustainable Digital Data Preservation and Access Network Partners (DataNet)**
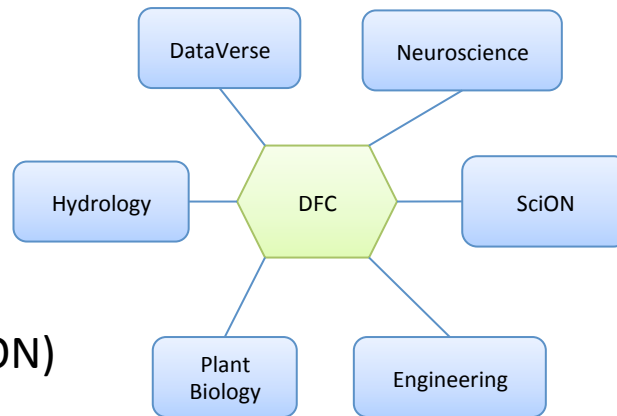
**Goals**

- Enable scientific collaboration

- Support access to live research data
  - Data Sharing

- Support reproducible data-driven research
  - Workflow Sharing

- Establish national data cyber-infrastructure
  - Persistent and Extensible Architecture



DataNet Federation Consortium
*Federating Resources to Accelerate Science*

# DFC Communities

**Research Communities**
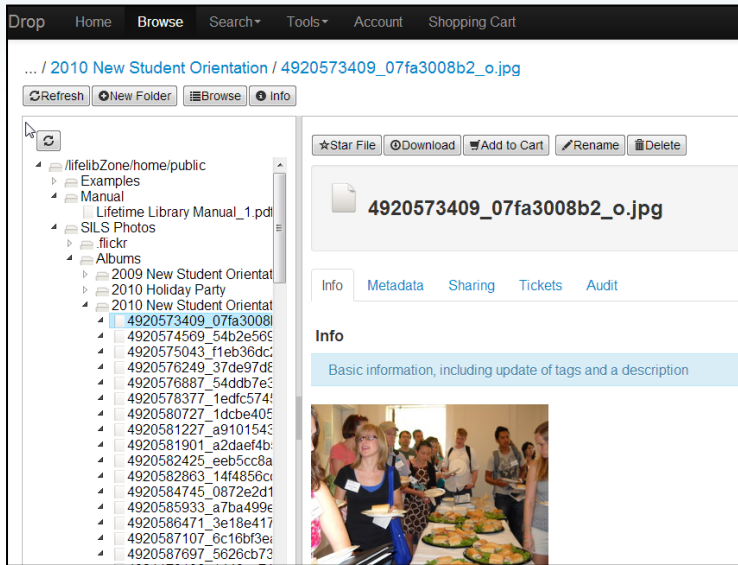
- Hydrology (UVA)

- Social Science (DataVerse)

- Neuroscience (TDLC)

- Ocean & Environment (SciON)

- Engineering (Drexel)

- Plant Biology (iPlant)
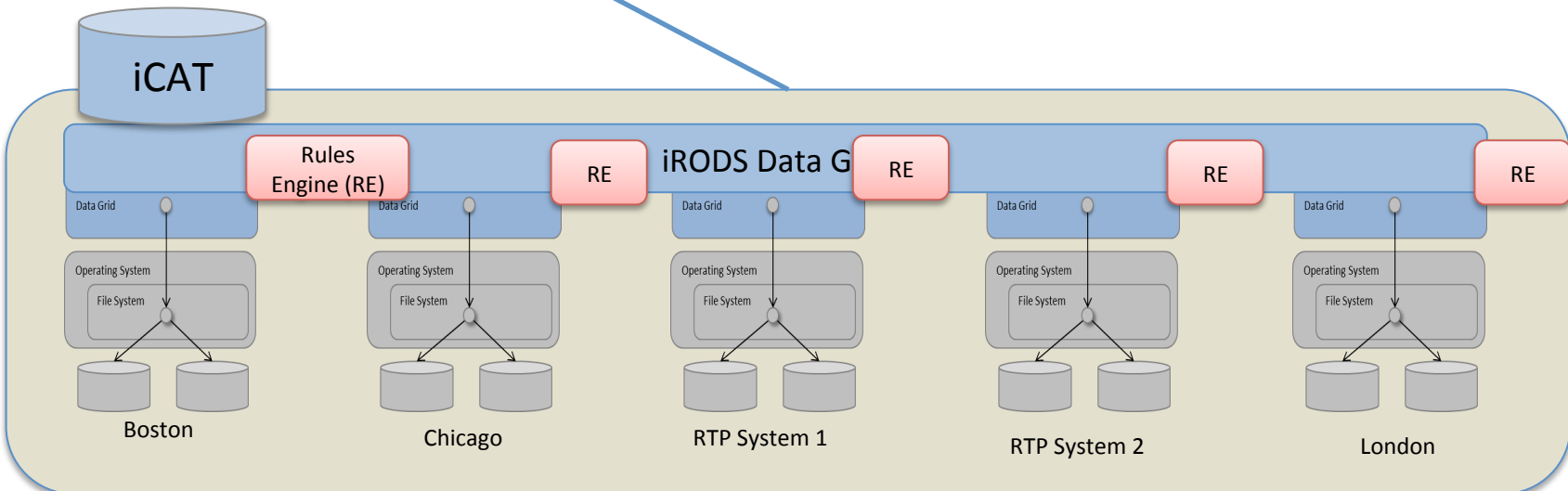
**Communities of Practice**

- Science and Engineering

- Technology and Research

- Facilities and Operations

- Policies and Standards

- Education and Outreach

- Sustainability

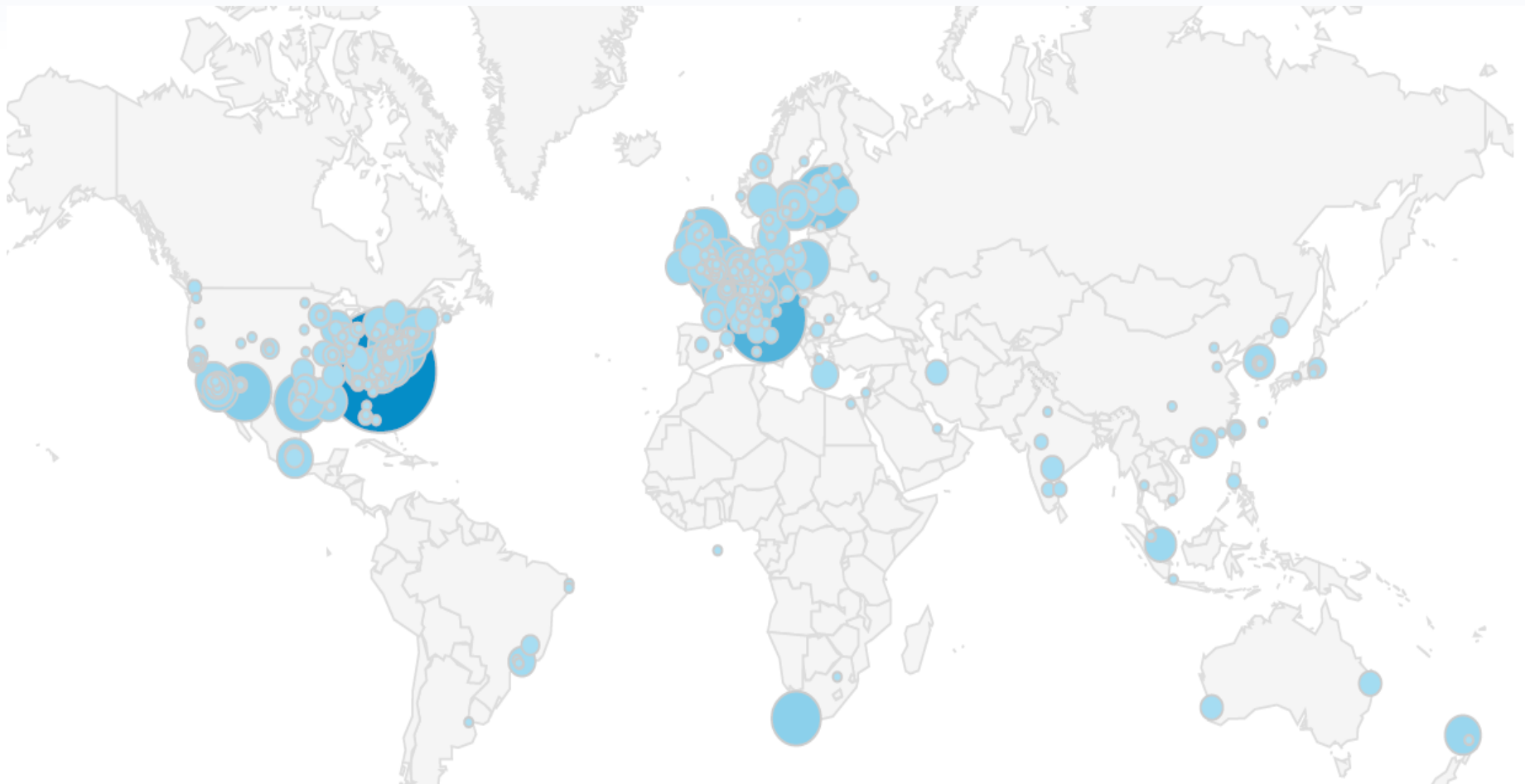# iRODS – Integrated Rules Oriented Data System



**Policy-based Data Management Platform**
- Virtual, Distributed Object Space
- Centralized Metadata
- Distributed Rules Engine allows:
  - Enforcement of policies
  - Automation of tasks
  - Orchestration of workflows
  - Controlled ingest, processing, and dissemination of digital objects

# iRODS Download Base



iRODS 4.0 downloads:  04/01/14-10/16/2014

renci

# iRODS EcoSystem

- Federal Users
  - National Aeronautics and Space Administration (NASA)
  - National Oceanic and Atmospheric Administration (NOAA)
  - National Optical Astronomy Observatory (NOAO)
  - US Geological Survey (USGS)

- Resellers/Redeployers/Partners
  - DataDirect Networks
  - Distributed Bio
  - Computer Sciences Corporation (CSC)
  - EMC
  - Seagate

- Commercial Users
  - DOW Chemical
  - Beijing Genome Institute

- Research Institutions
  - Broad Institute
  - International Neuroinformatics Coordinating Facilities (INCF)
  - Wellcome Trust Sanger Institute
  - Computer Center of the French National Institute of Nuclear and Particle Physics (CC-IN2P3)
  - CineGRID

- Hundreds of academic institutions worldwide

- Supported by the iRODS Consortium

renci

# DFC Models of Federation

**Strong Federation**

    Full and complete protocol-level federation across iRODS data grids

    Seamlessly Move from one grid to another

    "Mu Casa Su Casa"

    Used in DFC to federate Science and Engineering domain grids

**Weak Federation**

    DFC to External Resources

    Micro-services and Workflows

    DFC needs to 'know' the external protocol - plug-ins & wrappers

    Still seamless - external access problems hidden from user

    Used in DFC

        To access THREDDS (netCDF), Sensor system, federal data resources
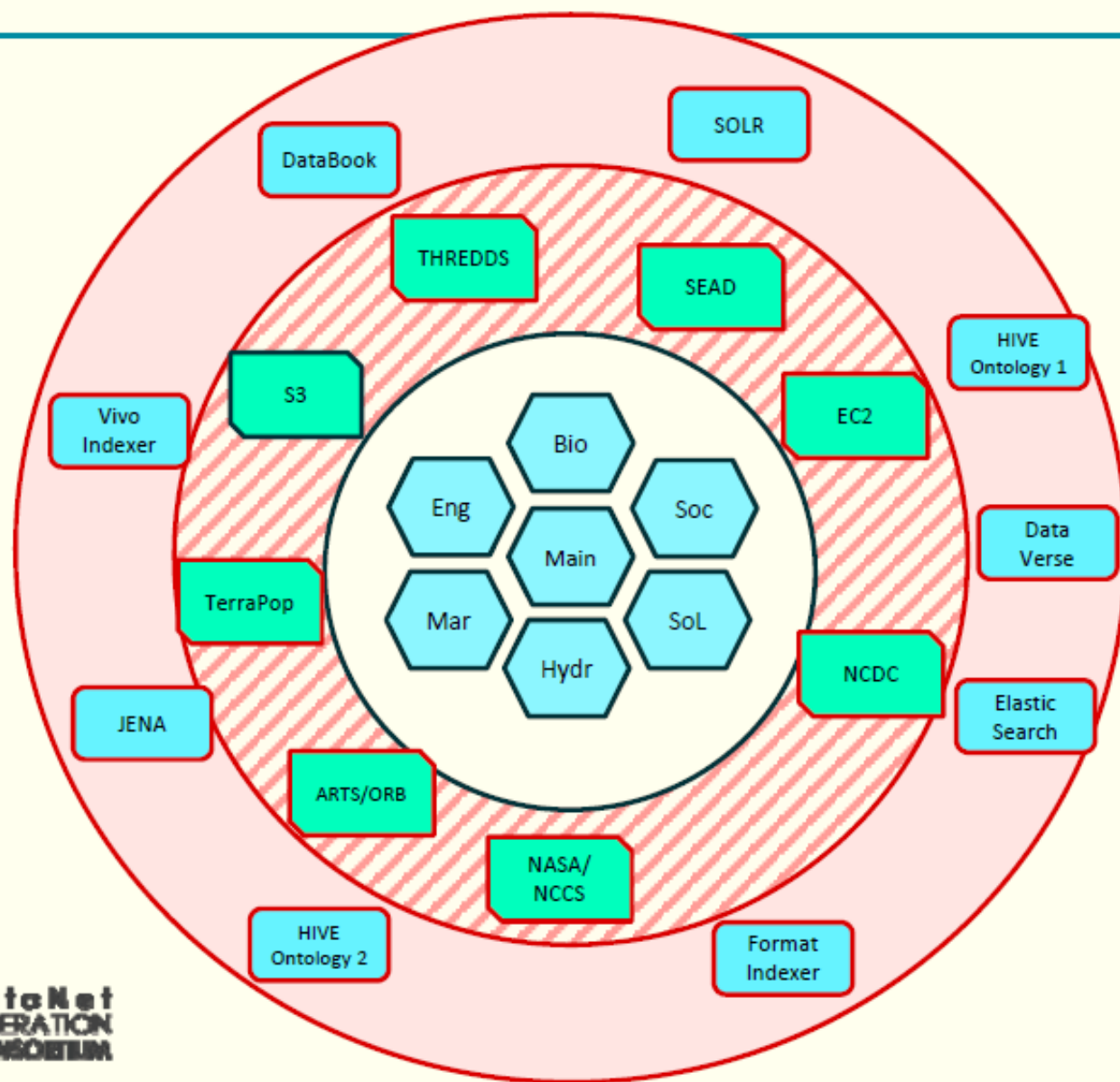
        To connect to SEAD, TerraPop, and DataONE

        To access Amazon Web Services (EC2, S3) and External workflows
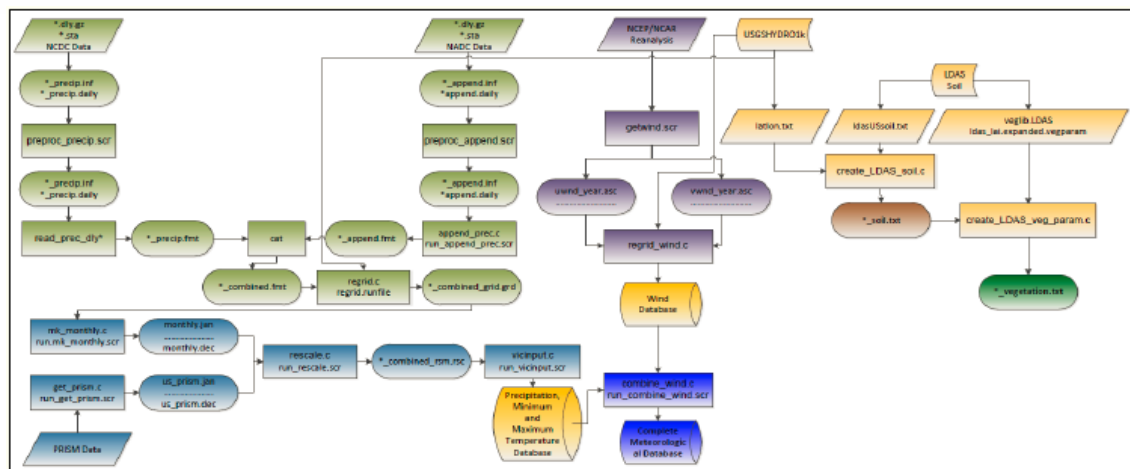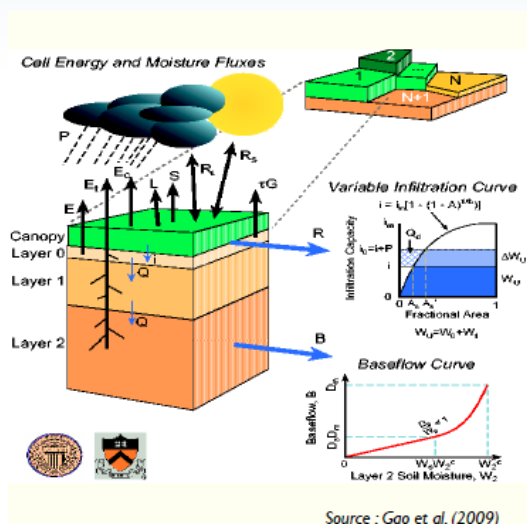
renci

# DFC Models of Federation

**Asynchronous Federation**

- Defined & Developed to meet needs of DFC
- Motivation Example: Provision access to indexing services
- Multiple indexing services – each with its own protocols
- DFC does not want to federate each separately – access is not seamless
  - Promote a common connectivity - based on message bus
  - Any indexer who can 'talk' this common connectivity can play
- Other examples: Metadata services, Ontologies, Formatting Services
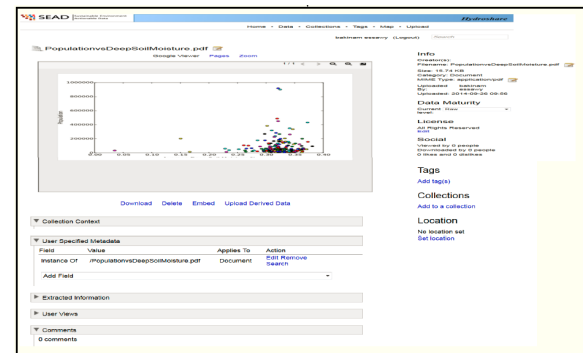
# Three Federations in DFC

# Interoperability Example



Source : Gao et al. (2009)



Source : Billah et al. (2014)

DFC Hydrology workflow (UVA)

- Input data from TerraPop data systems  (weak federation)
- Data processing across differenet iRODS data systems (strong federation)
- Output data to SEAD hydrology system (weak federation)
- Data generation/usage events can be monitored through DataBook (asynch federation)



SEAD Hydroshare Project Space

# Data Systems

- Purpose
  - Reason a collection is assembled
- Properties
  - Attributes needed to ensure the purpose
- Policies
  - Controls for enforcing desired properties
- Procedures
  - Functions that implement the policies
- Persistent state information
  - Results of applying procedures
- Property verification
  - Validation that state information conforms to purpose

- *These change throughout the lifecycle of the collection*

- *These are defined by the user/community of interest*

# Multiple aspects of interoperability in Data Systems

- Protocols & Interfaces
  - Services, APIs
- Authentication & Authorization
- Data
  - Discovery, Link, Retrieval, Submission, Manipulation, Syntax, Semantics
- Metadata
  - Discovery, Link, Retrieval, Submission, Manipulation, Syntax, Semantics
- Workflows and Computation
  - Automated control and data flows between systems
- Policies
  - Usage rights, audits, embargoes, verifications, encryption, FISMA, leakage protection