**An Earth Science Ontology Dialog**

# Value Proposition of Ontology and Semantic Technology for the Earth Science Community

## How Can Semantics Change Data Practices of the EarthCube Geoscience Community?

**Krishna Sinha**
**Geological Sciences**
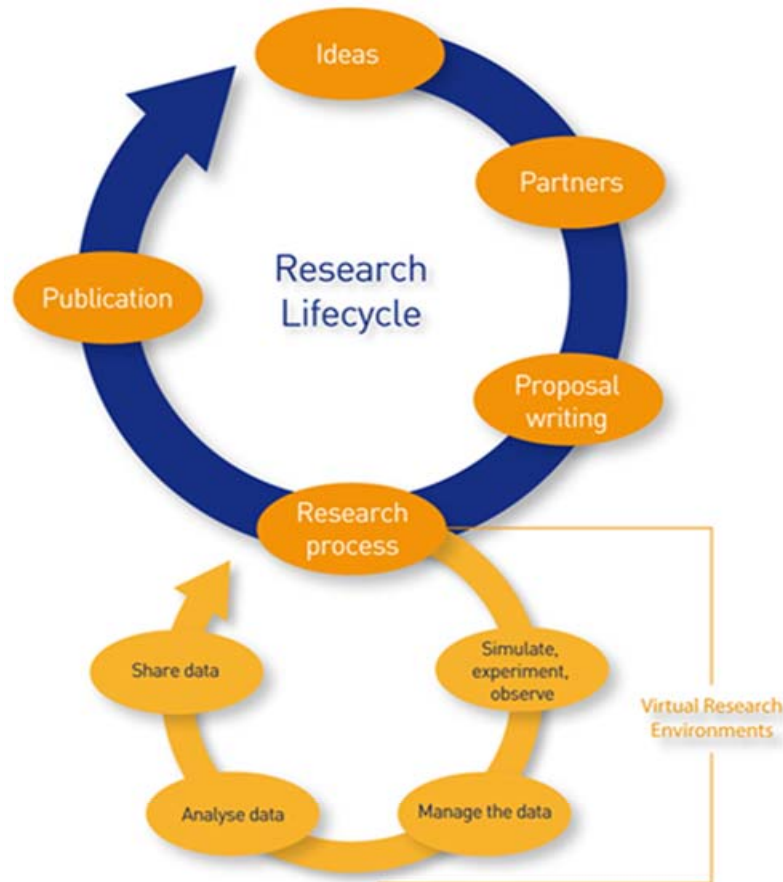**Virginia Tech**
**pitlab@vt.edu**

# Current Status of Geoscience community data practices

- Data and tool sharing practices face many barriers, and are especially <u>dominant in the world of individual researchers (long tail of science community)</u>

- However, these data are required to understand how natural systems (***earth systems***) change over time through physical, chemical, and biological processes.

- Reasons: lack of time, future publishing opportunities, ownership of data, concerns related to misuse of data, credit for professional advancement, lack of institutional support, and opportunities for commercial applications.

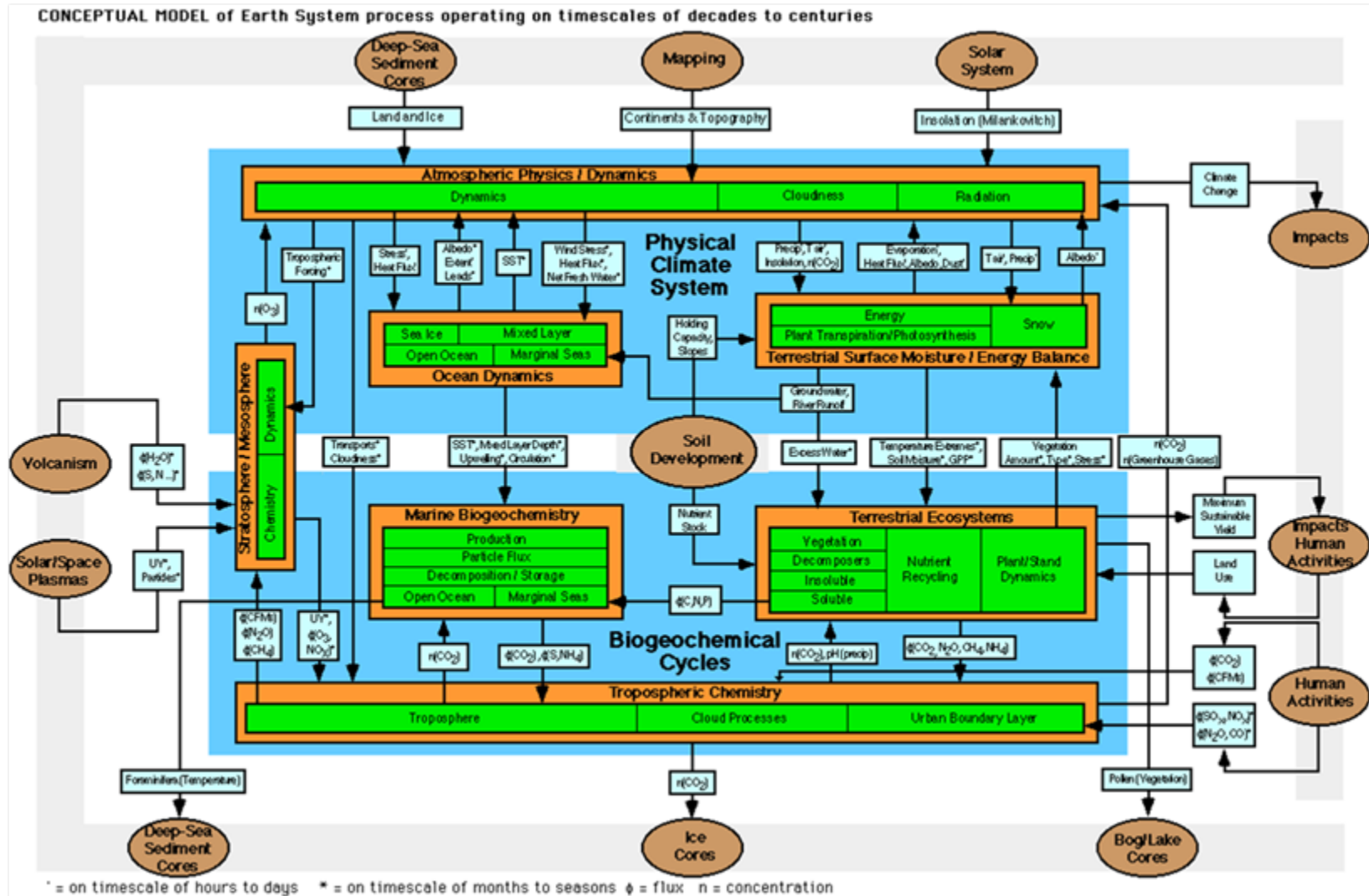**Krishna Sinha, Geological Sciences, Virginia Tech, pitlab@vt.edu**

# INTRODUCTION

- Data Life Cycle : involvement of individuals
- Complexity of data: earth as a system
- Data Environments : long tail of science
- Sharing data: Credit and motivation
- Motivational Integration scenario: multiple data resources and increased efficiency
- Sharing data at many levels of  domain ontology
- Individuals have more than data to share: ontology classes
- Stages in building the semantic infrastructure

3

# Research and Data Lifecycles

**Value Proposition of Ontology** ne.org/article/info:doi/10.1371/journal.pone.0021101    **Krishna Sinha, Geological Sciences, Virginia Tech, pitlab @vt.edu**

PLoS one

# Bretherton Diagram demonstrates the dynamic interaction of systems: the real world complexity of data and sources



CONCEPTUAL MODEL of Earth System process operating on timescales of decades to centuries

' = on timescale of hours to days    * = on timescale of months to seasons    φ = flux    n = concentration

**Krishna Sinha, Geological Sciences, Virginia Tech, pitlab@vt.edu**

# Data Environments



Data environments and properties

**Discipline specific vocabulary**

At established data centers
well planned
curated
Highly visible

Structured and homogeneous data
sensor based
large volumes

Resides in DATA CENTERS

Small volume of structured and homogeneous data

Resides in

Libraries ► Controlled vocabulary

Distributed
Poorly curated
Not visible
Dark data
Small projects

Unstructured
Heterogeneous
Small volume

Individual desktop computers

Personal notations, terms and acronyms

Long tail of Science Individuals

Multi-disciplinary Controlled vocabulary

Page 3 EarthScienceOntolog A.K.Sinha

6

# Discovery, Reuse, and Credit
# Mooney and Newton, 2012

- **scientists agree that data sharing is a desirable practice in theory, the fear of receiving no credit and losing funding or publishing opportunities is a serious deterrent to actual practice**

- **fear that sharing data could result in someone else publishing with no reward given to the sharer since there is no system of acknowledgement,**

- **some fairly famous cell lines were generated by obscure people**


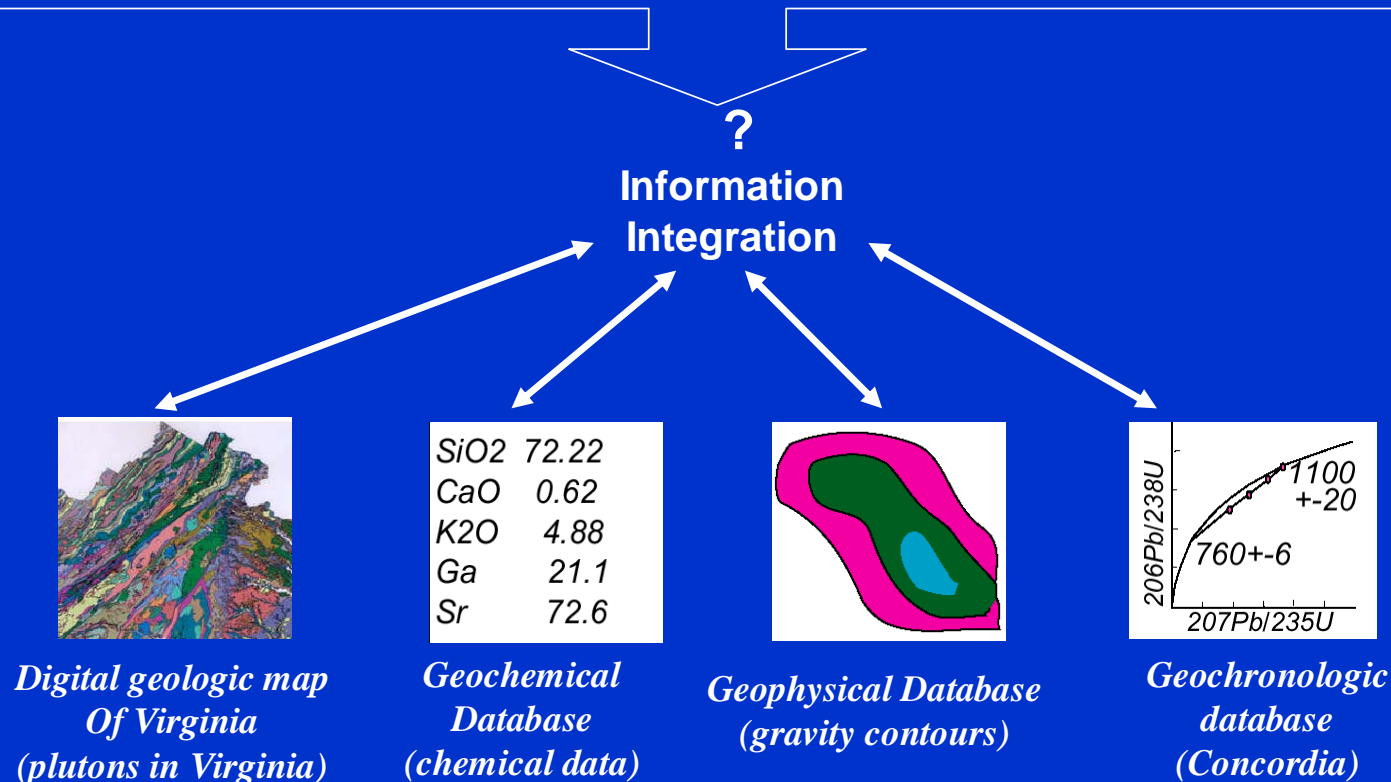- **What constitutes an adequate data citation**

# The Anatomy of a Data Citation: Discovery, Reuse, and Credit

- **Reward structures must be in place to encourage data publication, and citation is the appropriate tool for scholarly acknowledgment. Data citation also allows for the identification, retrieval, replication, and verification of data underlying published studies**

- **Promotion of data citation will foster a scholarly communication system that allows for identification, retrieval,**

- **Repositories publishing data should include appropriate metadata and mandate citations as a condition of reuse.**

- **Normalizing expectations for dataset citation will incentivize data sharing and promote secondary research,**

Krishna Sinha, Geological Sciences, Virginia Tech, pitlab@vt.edu

# Incentive through providing easy problem solving environment : use case exemplars resolved through ontologies
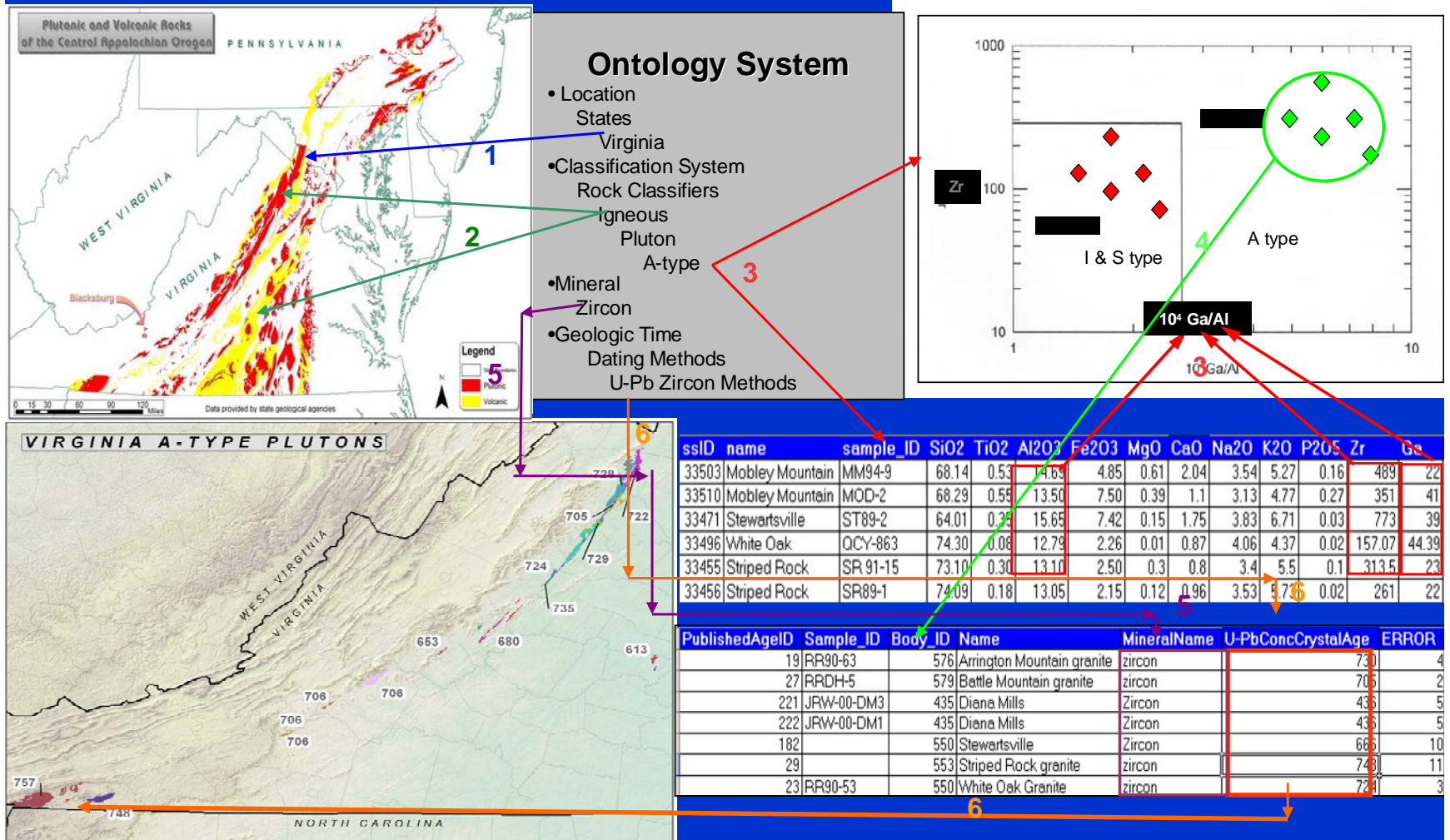
## A Geoscientist's Information Integration Scenario

*What is the distribution and U/ Pb zircon ages of A-type plutons in VA? How about their 3-D geometry using gravity data ?*

**?**
**Information Integration**

| SiO2 | 72.22 |
|------|-------|
| CaO | 0.62 |
| K2O | 4.88 |
| Ga | 21.1 |
| Sr | 72.6 |

*206Pb/238U*
1100 +-20
760+-6
*207Pb/235U*

*Digital geologic map Of Virginia (plutons in Virginia)*

*Geochemical Database (chemical data)*

*Geophysical Database (gravity contours)*

*Geochronologic database (Concordia)*

9

Krishna Sinha, Geological Sciences, Virginia Tech, pitlab@vt.edu

Integration Scenario: Stages for access to data and tools in a workflow environment

# Providing the ontology framework

**Sharing Geoscience Knowledge in a semantic world: Requires three classes of ontology frameworks**



- **Objects** represent our understanding of the state of the system when the data were acquired, while **processes** capture the physical and chemical forcings on objects that may lead to changes in state and condition over time. **Service** provides tools (e.g., simulation models and analysis algorithms) to assess multiple hypotheses, including inference or prediction.

- These three classes of ontologies within the semantic layer of the scientific cyberinfrastructure are thus required to enable automated discovery, analysis, utilization, and understanding of data through both induction and deduction

11

Krishna Sinha, Geological Sciences, Virginia Tech, pitlab@vt.edu

# Semantic infrastructure development stages: data to knowledge pathway



Thousands of providers for tools and services (APPS)

Tens of thousands of Data Providers

Types of Data

Text
Graphics
Numerical
Models

Heterogeneity
Volume

Data Pool

Data Centers
Agencies
Individuals
Libraries

Ontologies

**Object**
Upper Level
Disciplinary
level

Mid Level
Foundation
Level

**Service**
Service
Process

| Object | Service | Process |
|--------|---------|---------|
| What, where, when | Tools | How/Why |

Data Discovery          Integration

Disciplinary
Sub-discipline

Data Level

Linked data
Warehousing
Aggregation

Fusion

Knowledge Discovery

Education and training of the next generation workforce

**Tasks/Technology needed**

SITE registration

APPS registration

Credit
Motivation
Trust
Security
Ownership
Access
Data for sale

Key words
AGU / AGI /
Geoscience
World

Data level
registration
"Smart Search"

Registration
Technology

Registration of
tools and services

Integration
Technology

Access object
and service
(including
process)
ontologies

Center for
Geoinformatics

Legacy data Services  |  Standards Services  |  Open Access Services  |  Ontology Services  |  Registration services  |  Integration Services  |  Preservation Services

**Krishna Sinha, Geological Sciences, Virginia Tech, pitlab@vt.edu**

# Suggestions

- This community to endorse data citation

- Use real world use cases that long tail scientists relate to, and recognized increase in efficiency through semantics is likely to promote data sharing (aspects of cost benefit)