

Data Infrastructure

Carelyn Campbell, Ben Blaiszik, Laura Bartolo

November 1, 2016

Data Landscape

Collaboration Tools
(e.g. Google Drive,
DropBox, Sharepoint,
Github, MatIN)

Data Sharing
Communities
(e.g. Dryad, FigShare,
NanoHub, Kaggle, NDS)

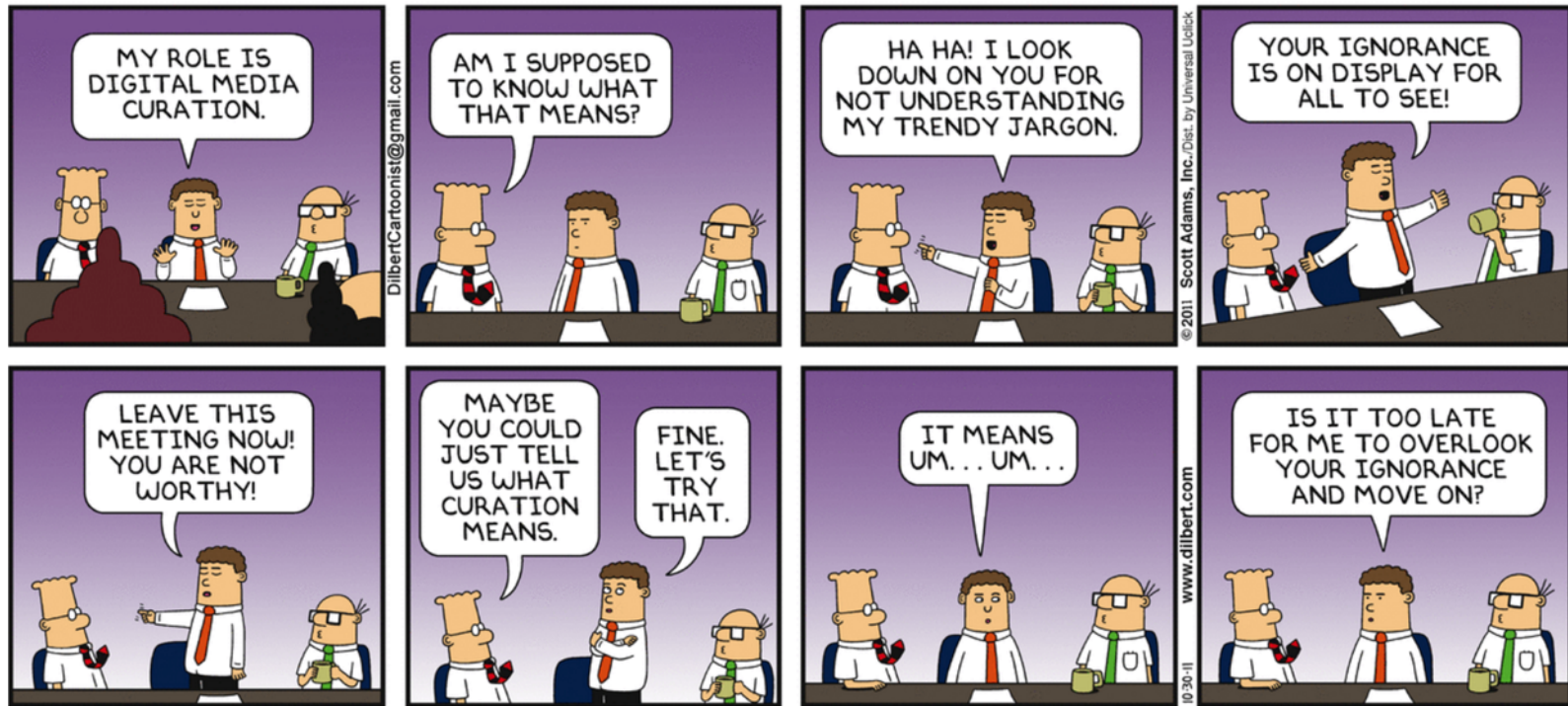
Data Repositories
(e.g. Aflow, MaterialsProject,
OQMD, NIMS MaterialNavi,
NoMaD, Materials Universe)

Data
Curation

Software

Data Analysis
Tools

What is Data Curation?

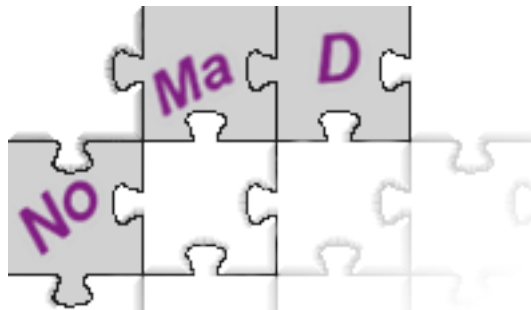


Scott Adams, October 30, 2011

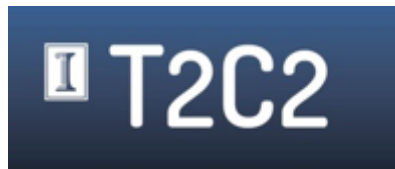
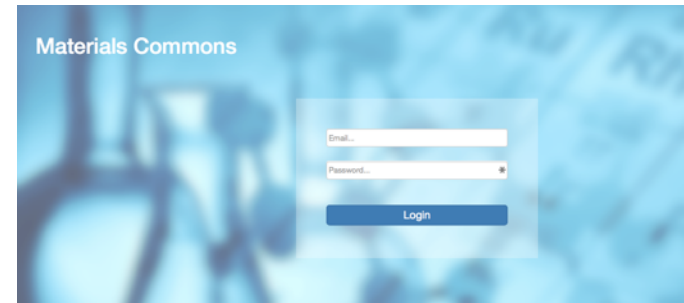
Data curation is the active and ongoing management of data through its lifecycle of interest and usefulness to scholarship, science, and education.

http://ischool.illinois.edu/academics/degrees/specializations/data_curation

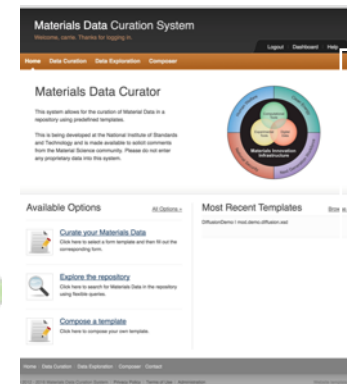
Materials Data Curation Tools



ICE INTEGRATED COLLABORATIVE ENVIRONMENT



Citrine



Data Model Definition

Defines the structure of **metadata** and **data**

Measurement Data Model

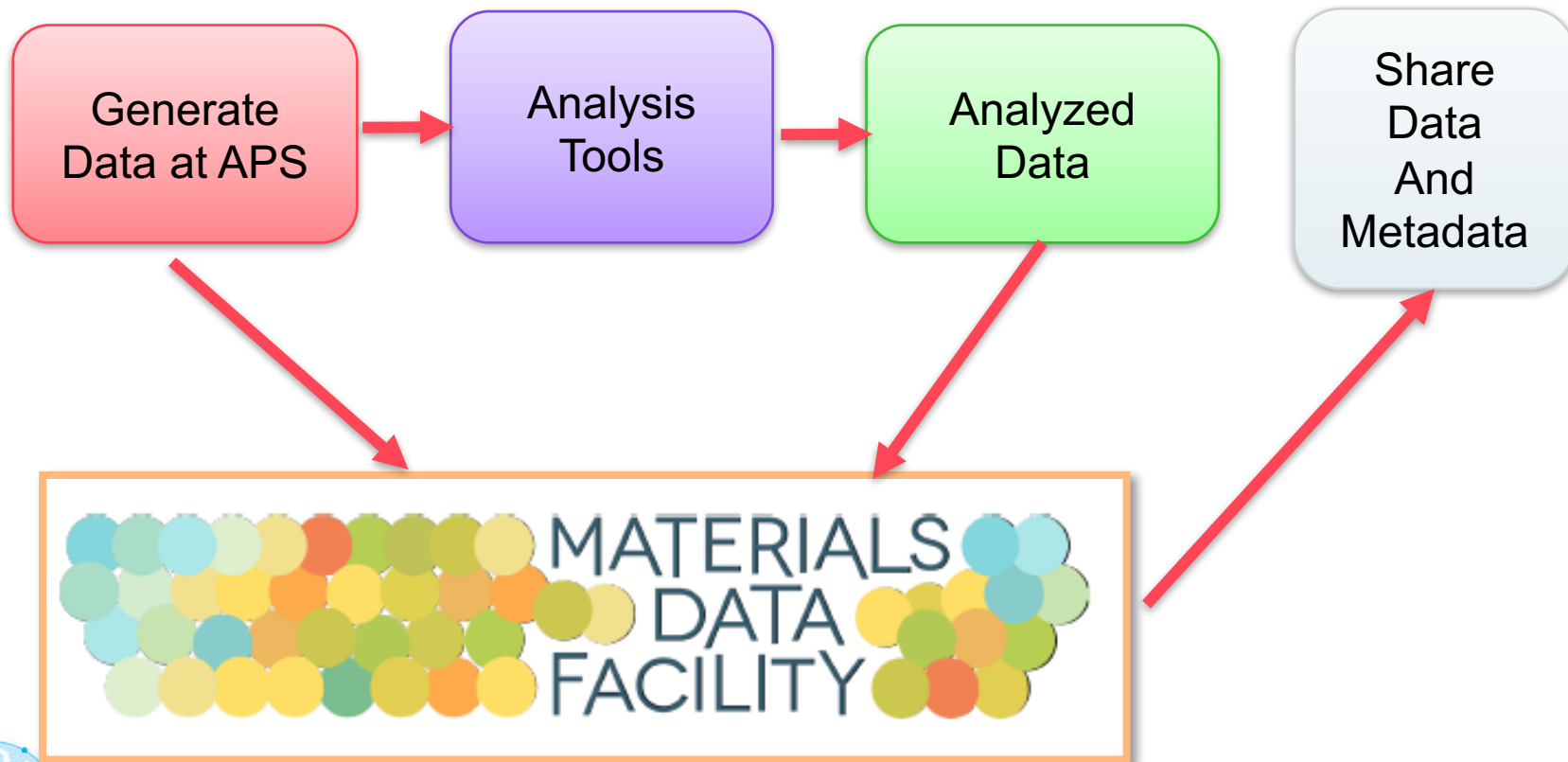
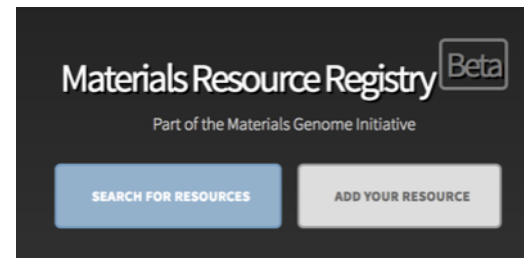
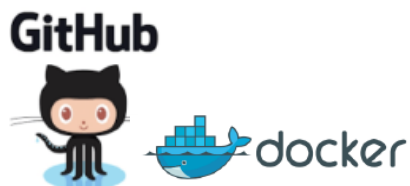
Metadata e.g.

- Sample owner
- Date of measurement $K\alpha_1$
- Sample stage position
- Apparatus temperature

Data e.g.

- As XML
- Raw data (text, ASCII, binary)
- Imported table
- *Link to image or raw data

Example APS Data (Large Data Example)



Workflows

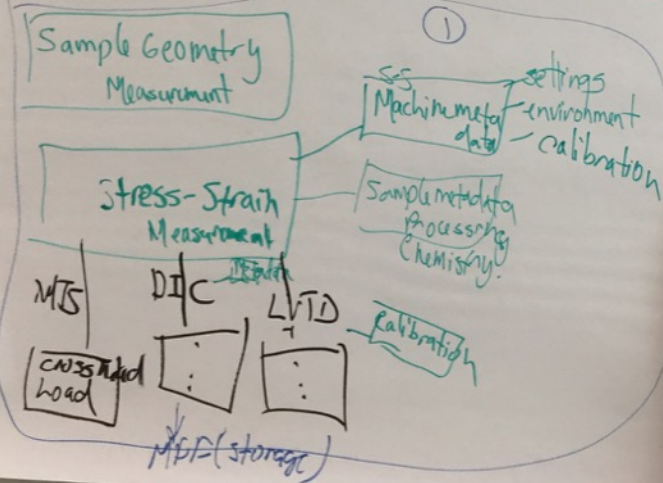
- Large Data sets: Single Point Source (e.g. APS)
- Experimental data (small to medium size), multiple source generation
- Computational Data
- Infrastructure Selection Tool

Stress-Strain Measurement

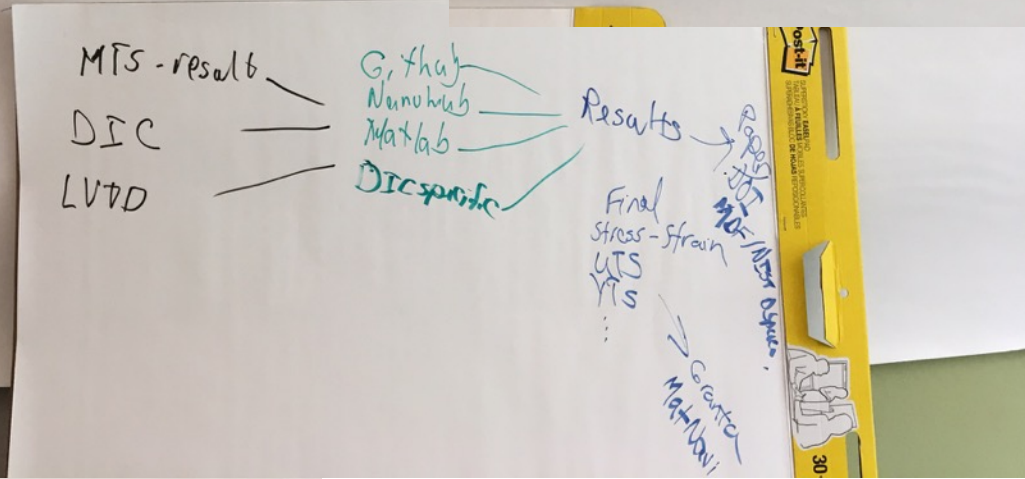
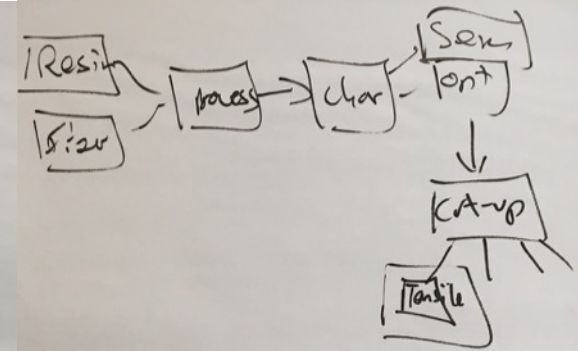
Experimental Data Workflow

- * Small \rightarrow Medium
- * Multiple modalities

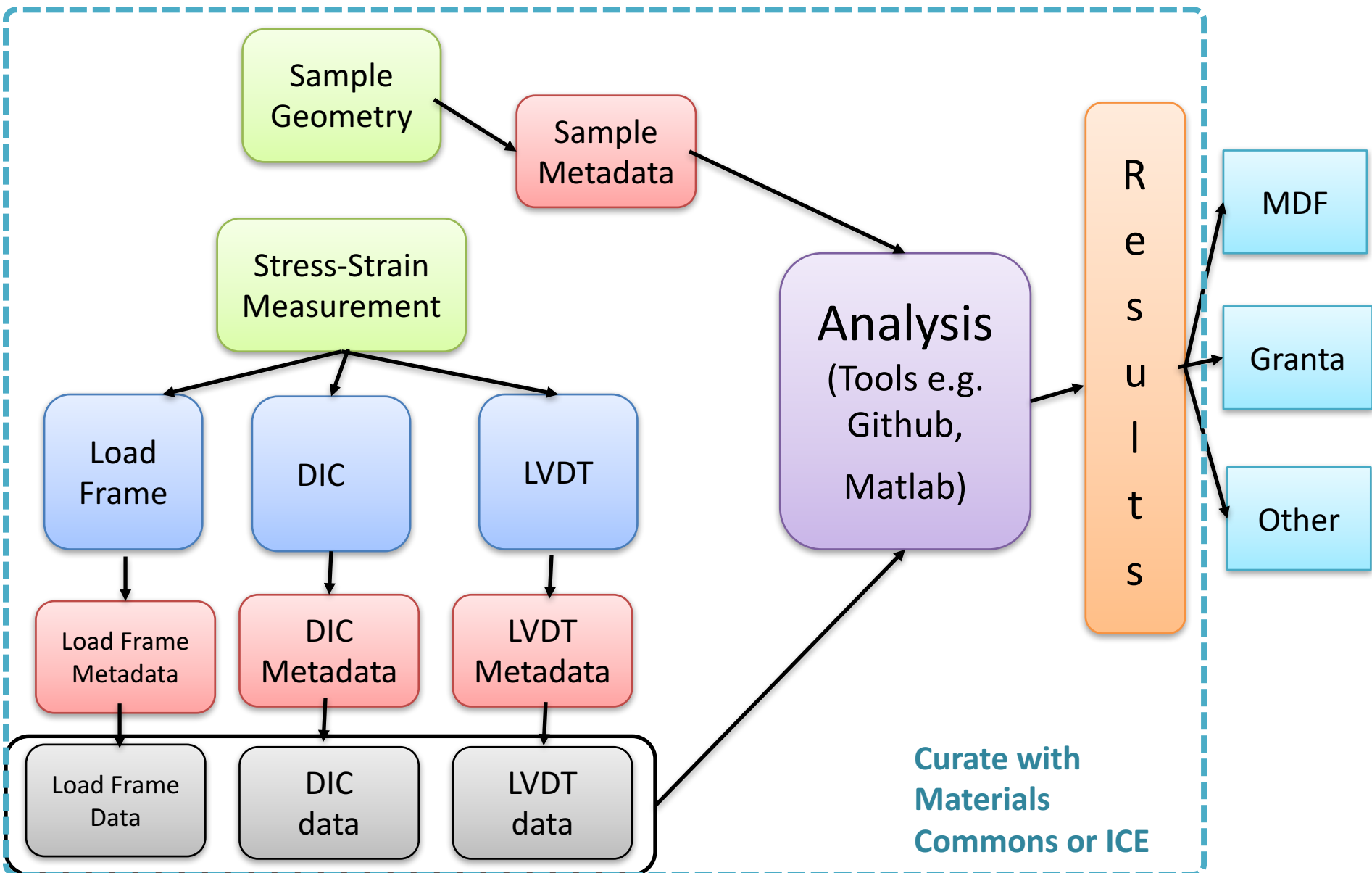
Example: Stress-Strain Curve



- ① Measurement Data
 - * ICF
 - * Material Comms
 - MDS
 - 4CEED (Microscopy)
- ② Measurement Analysis



Experimental Workflow: Stress-strain Measurement

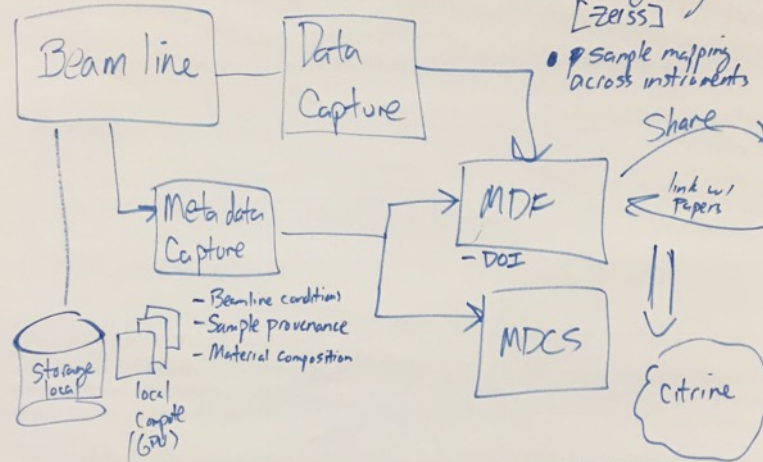


Big Data Workflow

Big Data Workflow

* APS, SNS, ... experimental or computational
~~Leadership~~ Computing Facilities

Schemas (per technique)



- IoT, capture everything [Zeiss]
- sample mapping across instruments

- How Long?
NSF-5yr
- Which data to keep?
 - ML to validate and value data
- What does it take to re-use?

- Experts
 - * full w/ all metadata
- Non-experts (technique)
 - * segmented data e.g.
- End user
 - * Aggregated
 - * Tank Data search

Computational Data Workflow

- Lots of different techniques
 - PhaseField modeling: no standards.
 - Community standards needed
-
- Codes changes quickly
 - Social change needed.
 - FEM - more benchmark. -- more standardize

How do I select a
Materials Data Infrastructure Tool?

Example: Workflow Tool Selection

A Taxonomy of Workflow Management Systems for Grid Computing

Jia Yu and Rajkumar Buyya*

Grid Computing and Distributed Systems (GRIDS) Laboratory, Department of Computer Science and Software Engineering, The University of Melbourne, Melbourne, Australia
E-mail: raj@cs.mu.oz.au

Received 28 May 2005; accepted in revised form 6 December 2005

Key words: Grid computing, resource management, scheduling, taxonomy, workflow management

Abstract

With the complex resources. Therefore, computing building development onomy n workflow

1. Intro

Grids [5] structure cations t heteroge such as h ics, geop utilizing sets. In c ments, d devices, e need to b tion wor [92].

Work procedur tween pa

* Corres

Table 2. Workflow design taxonomy mapping.

Project name	Structure	Model	Composition systems	QoS constraints
DAGMan	DAG	Abstract	User-directed • Language-based	User specified rank expression for desired resources
Pegasus	DAG	Abstract	User-directed • Language-based Automatic	N/A
Triana	Non-DAG	Abstract	User-directed • Graph-based	N/A
ICENI	Non-DAG	Abstract	User-directed • Language-based • Graph-based	Metrics specified by users
Parana	DAG	Abstract	User-directed • Language-based • Graph-based	N/A
GridAnt	Non-DAG	Concrete	User-directed • Language-based	N/A
GrADS	DAG	Abstract	User-directed • Language-based	Estimated application execution time

Table 2. Workflow design taxonomy mapping.

Project name	Structure	Model	Composition systems
DAGMan	DAG	Abstract	User-directed • Language-based
Pegasus	DAG	Abstract	User-directed • Language-based Automatic
Triana	Non-DAG	Abstract	User-directed • Graph-based
ICENI	Non-DAG	Abstract	User-directed • Language-based • Graph-based

Example: Hardware Store Website

- + Brand
- + Price
- + Counter Depth (Yes/No)
- + Refrigerator Width (In.)
- + Height to Top of Door Hinge
- + Depth (Excluding Handles)
- + Color/Finish Family
- + Capacity (cu. ft.) - Refrigerators
- + Review Rating
- + General Features
- + Refrigeration Dispenser Features
- + Ice/Water Dispenser
- + Appliance Series
- + ENERGY STAR CERTIFIED
- + Eco Options

Refrigerator Width (In.)

- 28 - 29 (151)
- 29 - 29.9 (209)
- 30 - 31.9 (64)
- 32 - 32.9 (217)
- 33 - 34.9 (23)
- 35 - 36 (456)
- 36 or Greater (138)

Height to Top of Door Hinge

- 66.25 in (21)
- 66.31 in (1)
- 66.44 in (4)
- 66.5 (11)
- 66.5 in (9)
- 66.62 in (2)
- 66.625 (11)

Credit: homedepot.com

Any mention of commercial products is for information only; it does not imply recommendation or endorsement by NIST.

Example: Used Car Website

- Price +
- Mileage +
- Years +
- Used / New +
- Make +
- Type +
- Features +
- Size +
- Exterior Color +
- Interior Color +
- MPG City +
- MPG Highway +
- Cylinders +
- Transmission +
- Packages +
- Domestic / Import +
- Store +

MPG City —

- Over 20 MPG city (4507)
- Over 25 MPG city (1875)
- Over 30 MPG city (459)
- Over 35 MPG city (256)
- Over 40 MPG city (191)

Features —

- Adjustable Suspension (98)
- Air Conditioning (8091)
- Alloy Wheels (6688)
- Auto Cruise Control (217)
- Automated Parking (46)

Credit: carmax.com

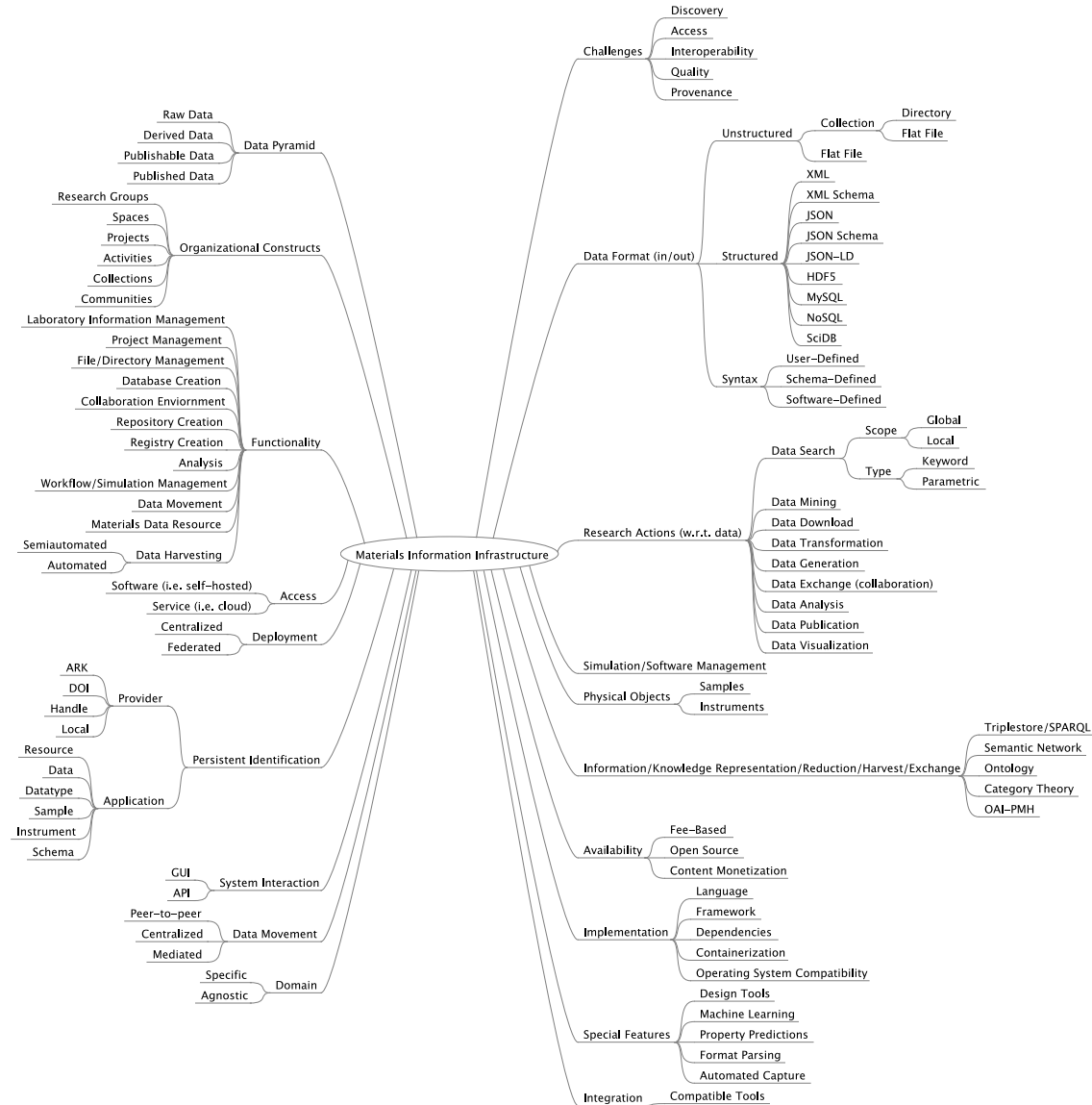
Any mention of commercial products is for information only; it does not imply recommendation or endorsement by NIST.

Registry: Materials Data Infrastructure Tools

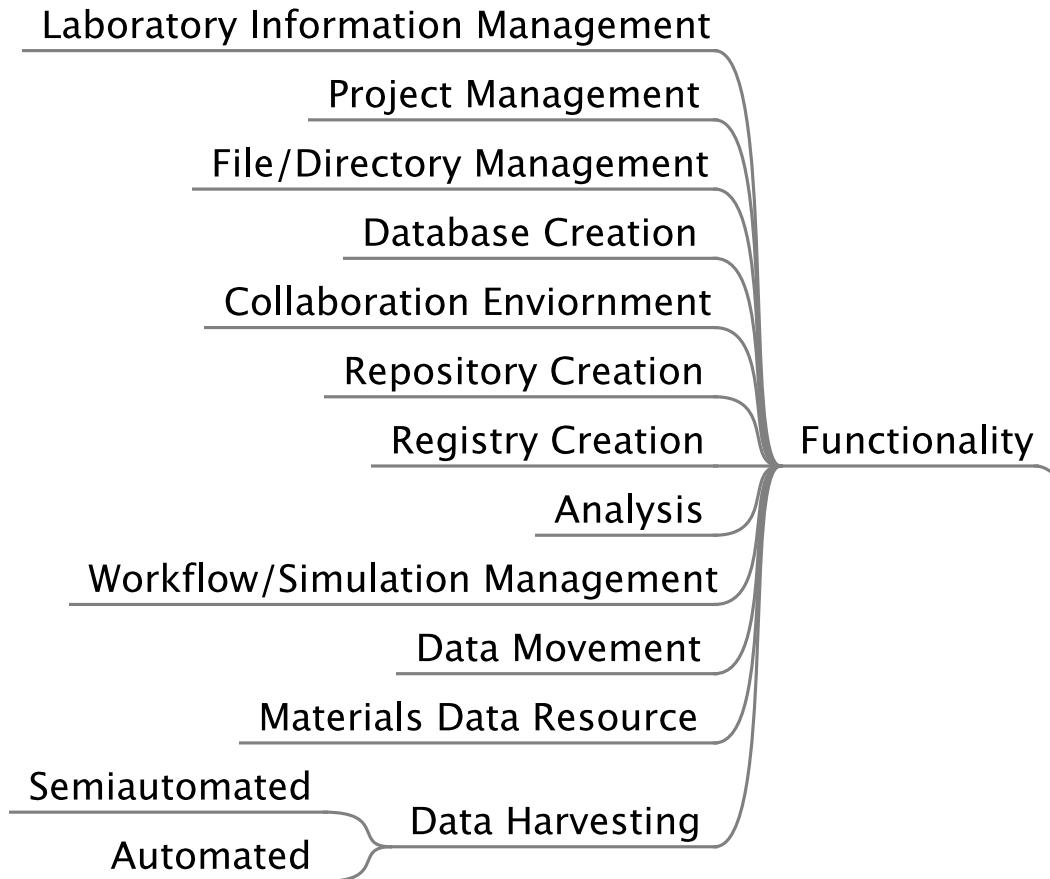
Material Types	<input type="checkbox"/> Metal <input type="checkbox"/> Semiconductor <input type="checkbox"/> Ceramic <input type="checkbox"/> Polymer <input type="checkbox"/> Biomaterial	<input type="checkbox"/> Organic <input type="checkbox"/> Inorganic <input type="checkbox"/> Oxide <input type="checkbox"/> Composite <input type="checkbox"/> Nanomaterials	<input type="checkbox"/> Superconductor <input type="checkbox"/> Non-Specific <input type="checkbox"/> Other	<input type="checkbox"/> ? (recommended)
Morphology/Structures	<input type="checkbox"/> Crystalline <input type="checkbox"/> Amorphous <input type="checkbox"/> Fluid <input type="checkbox"/> Quasi-periodic <input type="checkbox"/> Bulk <input type="checkbox"/> 2-Dimensional	<input type="checkbox"/> 1-Dimensional <input type="checkbox"/> Film <input type="checkbox"/> Nanotube <input type="checkbox"/> Fiber <input type="checkbox"/> Composite <input type="checkbox"/> Interfacial	<input type="checkbox"/> Interphase <input type="checkbox"/> Line Defect <input type="checkbox"/> Point Defect <input type="checkbox"/> Non-Specific <input type="checkbox"/> Other	<input type="checkbox"/> ? (recommended)
Material Property Classes	<input type="checkbox"/> Optical <input type="checkbox"/> Mechanical <input type="checkbox"/> Thermodynamic	<input type="checkbox"/> Structural <input type="checkbox"/> Simulation	<input type="checkbox"/> Other	<input type="checkbox"/> ? (recommended)
Experimental Data Acquisition Methods	<input type="checkbox"/> Electrochemical <input type="checkbox"/> Mechanical Testing	<input type="checkbox"/> Scanning Electron Microscopy <input type="checkbox"/> Spectroscopy <input type="checkbox"/> Optical Microscopy <input type="checkbox"/> Impact Testing	<input type="checkbox"/> Indentation <input type="checkbox"/> Dilatometry <input type="checkbox"/> Other	<input type="checkbox"/> ? (recommended)
Computational Data Acquisition Methods	<input type="checkbox"/> Density Functional Theory <input type="checkbox"/> Molecular Dynamics Simulation <input type="checkbox"/> Numerical Simulations <input type="checkbox"/> Multiscale <input type="checkbox"/> Finite Element Analysis <input type="checkbox"/> Computational Thermodynamics	<input type="checkbox"/> Statistical Mechanics <input type="checkbox"/> Dislocation Dynamics <input type="checkbox"/> Phase Field <input type="checkbox"/> Crystal Plasticity <input type="checkbox"/> Other	<input type="checkbox"/> ? (recommended)	
Sample Processing Methods	<input type="checkbox"/> Casting <input type="checkbox"/> Annealing <input type="checkbox"/> Vapor Deposition <input type="checkbox"/> Milling	<input type="checkbox"/> Extrusion <input type="checkbox"/> Pressing <input type="checkbox"/> Exfoliation <input type="checkbox"/> Melt Blending	<input type="checkbox"/> Polymerization <input type="checkbox"/> Curing <input type="checkbox"/> Evaporation <input type="checkbox"/> Other	<input type="checkbox"/> ? (recommended)

Need "checkboxes" for Materials Data Infrastructure Tools

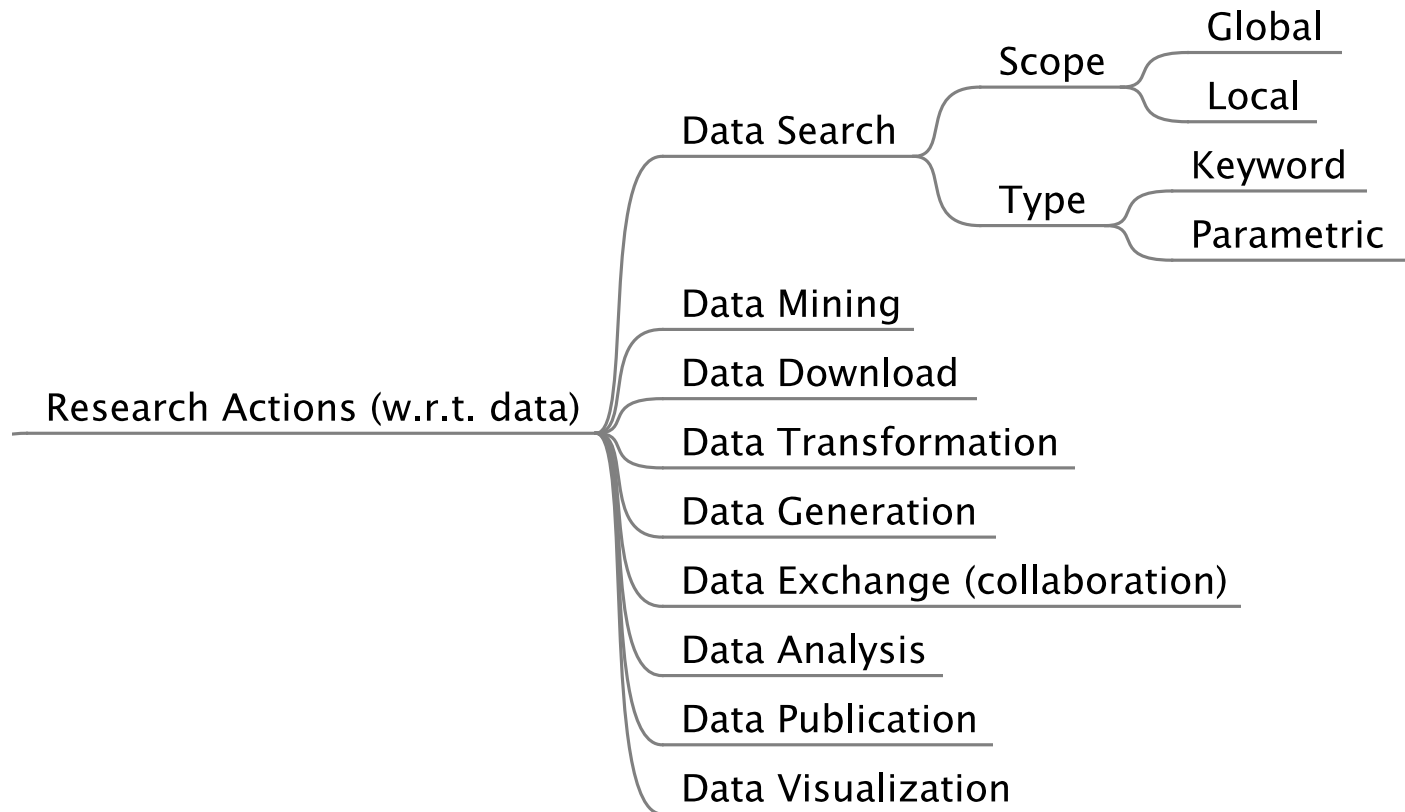
First Rough Draft Taxonomy



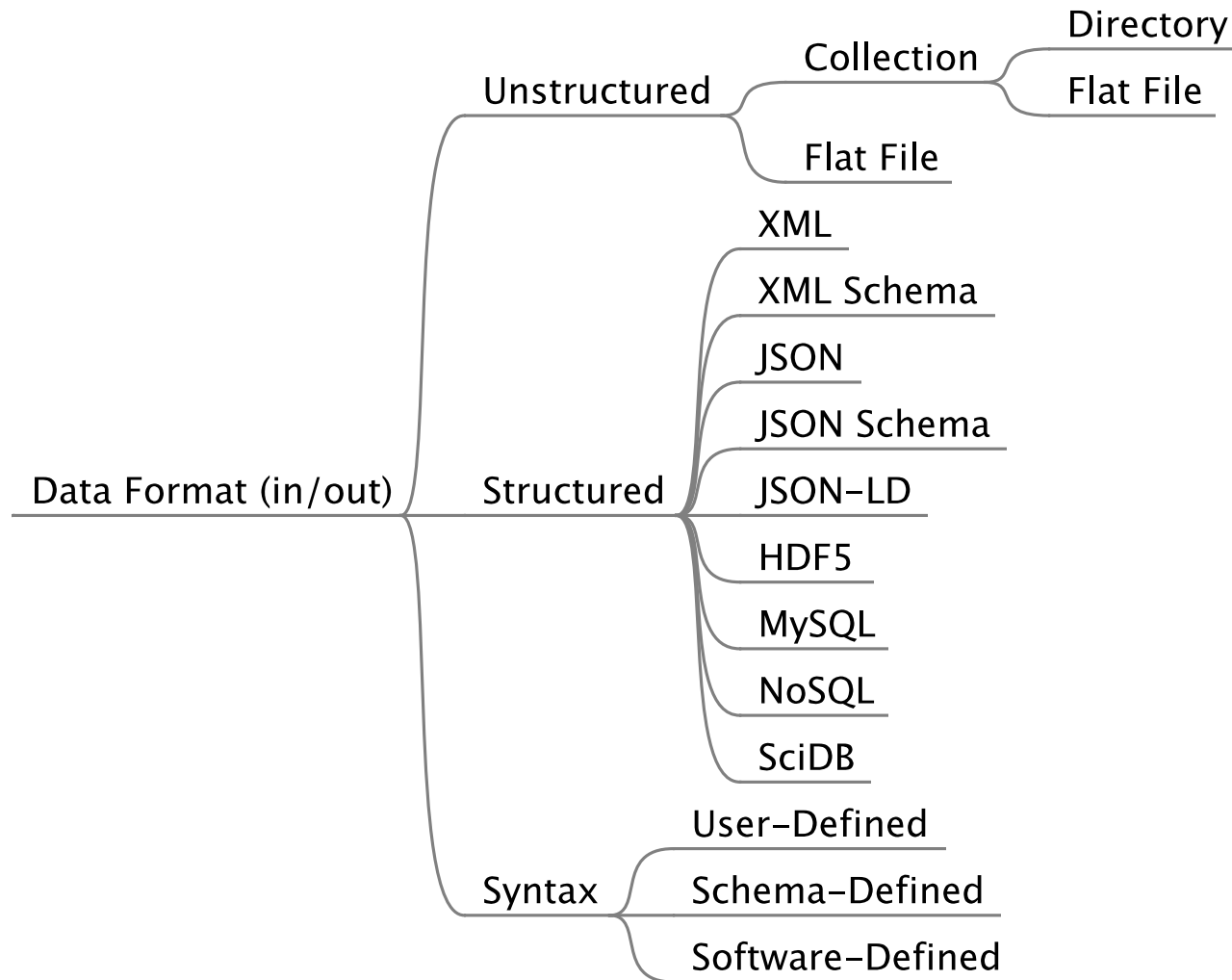
First Rough Draft Taxonomy



First Rough Draft Taxonomy



First Rough Draft Taxonomy



Notes from Summit Wrap-up Session

- Integrate tools into undergraduate education
 - Tools need to more user friendly
- Embed data experts into experimental groups
 - Alternate: floating data experts available for experimental groups.
 - Need to define skills needed for these data experts
- Encourage more conference exchanges at Data Analytics and Materials communities
- Define data curation guidelines/code
 - Benefit to users
- Data Challenge (Student)
 - Prize for data set
 - Best paper/DOI/PID
- Develop implementation path
- Improve peer recognition
- Develop data cite profile

Interest in following up with small working groups on specific issues.

Data Cit. Profile

Tool Integ. into undergrad edu
- More user-friendly

Embedding data into exp.
or floating data experts. groups

Data Analytics/Materials conf. exchange

Guidelines for data curation/
(Benefit to users)

Data Challenge (student-band) ^{code}
- prize
- best paper/DOF-PIID

Peer recognition

Role/Function Implementation Path