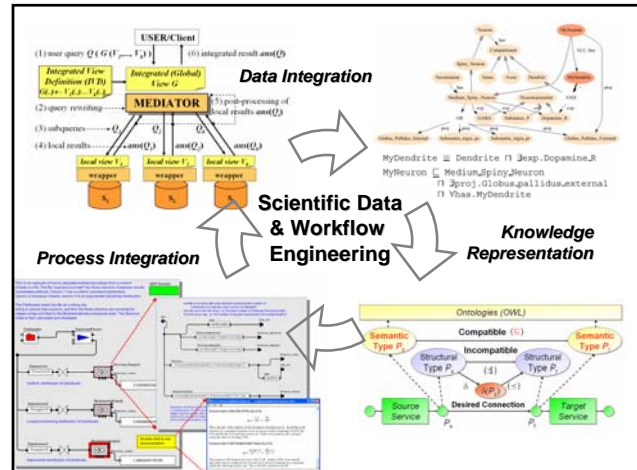


ECS289F Winter'05 Scientific Data Management

Bertram Ludäscher

Associate Professor
Dept. of Computer Science & Genome Center
University of California, Davis
ludaesch@ucdavis.edu



Logistics

- **Class:**
 - MWF, 4:10-5pm, 244 Olson
 - <http://www.sdsc.edu/~ludaesch/ECS289F-W05.html> (or Google with "Bertram ECS"; then navigate)
 - Notes will be posted there (& hand-outs available)
- **Office hours:**
 - General:
 - Wednesday & Friday, 11am-noon, room 3051, Kemper Hall
 - Extra:
 - email ludaesch@ucdavis.edu, subject ECS289F
- **Grading:**
 - Project 50%, presentation 30%, homework 20%



B. Ludäscher, ECS289F-W05, Topics in Scientific Data Management

Projects

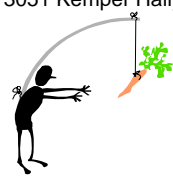
- **Implementation Projects (Hands-on)** with scientific data management and workflow tools:
 - [KEPLER](#) scientific workflow system
 - [SDSC Storage Resource Broker](#)
 - e.g. run a number of compute-intensive jobs (say from cheminformatics) on a cluster computer via KEPLER, or
 - access a real-time seismic sensor network client application using KEPLER and SRB
- **Research Projects (Theory):**
 - Readings in **data integration**: understand how a **database mediator** works i.e., *query rewriting*
 - Readings in **knowledge representation**: understand how *ontology* languages are used to capture and reason with domain knowledge
- **Combined Projects (Theory & Hands-on)**
 - Implementation of algorithms for query rewriting



B. Ludäscher, ECS289F-W05, Topics in Scientific Data Management

Research Assistant (RA)-ships

- Openings at the **CS department** and new **Genome Center** available; option for **summer internship** at the **San Diego Supercomputer Center**!
- **Mix of theory and practice** ideal:
 - Databases (& Information Systems) background helps a lot
 - No fear of ontologies or workflows (that's why we have this class ;-)
 - Problem-solving (thinking!) and Programming skills
- Come to my office hours (WF, 11am-noon, 3051 Kemper Hall) for details

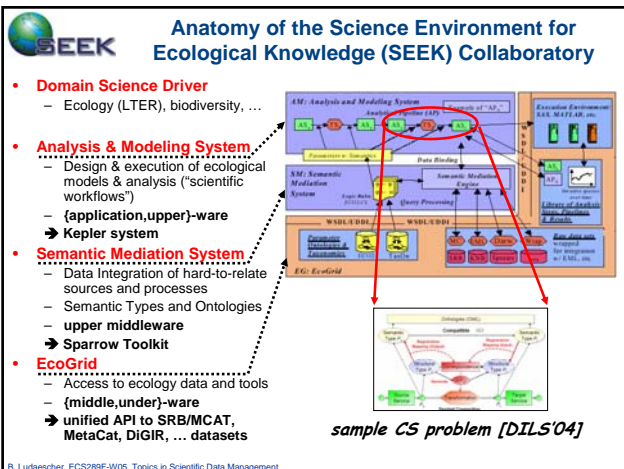


B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

Today's Outline

- Semantics & Scientific Data Integration
- Semantics & Scientific Workflow Management
- Conclusions

B. Ludascher, ECS289F-W05, Topics in Scientific Data Management



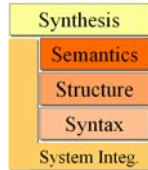
B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

Common Collaboratories / Distributed Science / Cyberinfrastructure Pieces

- Seamless and uniform data access ("Data-Grid")
 - data & metadata registry
 - distributed and high performance computing platform ("Compute-Grid")
 - service registry
 - User-friendly workbench / problem-solving environment
 - scientific workflow system
 - A common problem:
 - **integrating** (or at least **linking**) data from multiple sites, investigators, communities, ..., scales, ..., species, ...
- Federated, integrated, mediated databases
- often use of semantic extensions (e.g. ontologies)

B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

Interoperability & Integration Challenges



- reconciling S^5 heterogeneities
- “gluing” together resources
- bridging information and knowledge gaps computationally

- **System aspects: “Grid” Middleware**
 - distributed data & computing, SOA
 - web services, WSDL/SOAP, WSRF, OGSA, ...
 - *sources = functions, files, data sets ...*
- **Syntax & Structure: (XML-Based) Data Mediators**
 - wrapping, restructuring
 - (XML) queries and views
 - *sources = (XML) databases*
- **Semantics: Model-Based/Semantic Mediators**
 - conceptual models and declarative views
 - Knowledge Representation: ontologies, description logics (RDF(S), OWL ...)
 - *sources = knowledge bases (DB+CMs+ICs)*
- **Synthesis: Scientific Workflow Design & Execution**
 - Composition of declarative and procedural components into larger workflows
 - *(re)sources = services, processes, actors, ...*
 - *Semantic extensions needed here as well!*

B. Ludäscher, ECS289F-W05, Topics in Scientific Data Management

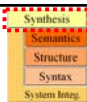
Information Integration Challenges: S^4 Heterogeneities



- **System aspects**
 - platforms, devices, data & service distribution, APIs, protocols, ...
 - ➔ **Grid middleware technologies**
 - + e.g. single sign-on, platform independence, transparent use of remote resources, ...
- **Syntax & Structure**
 - heterogeneous data formats (*one for each tool ...*)
 - heterogeneous data models (*RDBs, ORDBs, OODBs, XMLDBs, flat files, ...*)
 - heterogeneous schemas (*one for each DB ...*)
 - ➔ **Database mediation and warehousing technologies**
 - + XML-based data exchange, integrated views, transparent query rewriting, ...
- **Semantics**
 - descriptive metadata, different terminologies, implicit assumptions & hidden semantics (“context”) of experiments, simulations, observation, ...
 - ➔ **Knowledge representation & semantic mediation technologies**
 - + “smart” data discovery & integration
 - + e.g. ask about *X* (“mafic”); find data about *Y* (“diorite”); be happy anyways!

B. Ludäscher, ECS289F-W05, Topics in Scientific Data Management

Information Integration Challenges: S^5 Heterogeneities



- **Synthesis** of applications, analysis tools, data & query components, ... into “scientific workflows”
 - How to make use of these wonderful things & put them together to solve a scientist’s problem?
- ➔ **Scientific Problem Solving Environments (PSEs)**
 - ➔ **Portals, Workbench** (“scientist’s view”, end user)
 - + ontology-enhanced data registration, discovery, manipulation
 - + creation and registration of new data products from existing ones, ...
 - ➔ **Scientific Workflow System** (“engineer’s view”, tool maker)
 - + for designing, re-engineering, deploying analysis pipelines and scientific workflows; *a tool to make new tools ...*
 - + e.g., creation of new datasets from existing ones, dataset registration, ...

Not discussed here: the “6th S”: **Social challenges** ...

B. Ludäscher, ECS289F-W05, Topics in Scientific Data Management

Our Focus

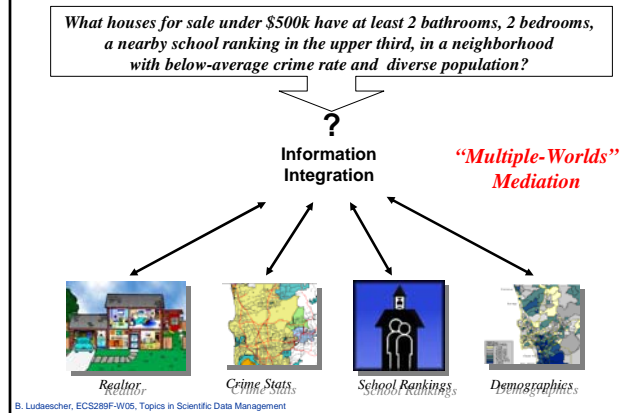
- **Scientific Data Integration:**
 - need DB/DI + KR (“semantic mediation”)
- **Automation of Scientific Data Analysis, Process & Application Integration**
 - need for scientific workflow systems
 - need for semantic extensions
- **But first:**
 - Some data & information integration problems

B. Ludäscher, ECS289F-W05, Topics in Scientific Data Management

An Online Shopper's Information Integration Problem



A Home Buyer's Information Integration Problem



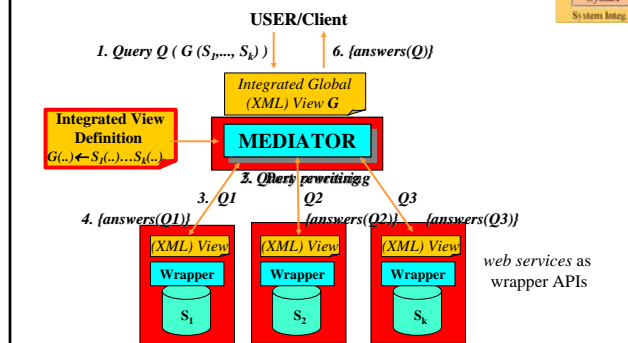
Information Integration from a Database Perspective

- Information Integration Problem
 - Given:** data sources S_1, \dots, S_k (databases, web sites, ...) and user questions Q_1, \dots, Q_n that can—in principle—be answered using the information in the S_i
 - Find:** the answers to Q_1, \dots, Q_n
- The Database Perspective: **source = "database"**
 - S_i has a **schema** (relational, XML, OO, ...)
 - S_i **can be queried**
 - define virtual (or materialized) **integrated (or global) view G** over local sources S_1, \dots, S_k using **database query languages** (SQL, XQuery, ...)
 - questions become queries** Q_i against $G(S_1, \dots, S_k)$



B. Ludaescher, ECS289F-W05, Topics in Scientific Data Management

Standard Mediator Architecture



B. Ludaescher, ECS289F-W05, Topics in Scientific Data Management

Query Planning in Data Integration

- **Given:**
 - Declarative user query Q: **answer(...)** \leftarrow ...G ...
 - ... & { G \leftarrow ... S ... } global-as-view (GAV)
 - ... & { S \leftarrow ... G ... } local-as-view (LAV)
 - ... & { ic(...) \leftarrow ... S ... G... } integrity constraints (ICs)
- **Find:**
 - equivalent (or minimal containing, maximal contained) query plan Q': **answer(...)** \leftarrow ... S ...
 - query rewriting (logical/calculus, algebraic, physical levels)
- **Results:**
 - A variety of results/algorithms; depending on classes of queries, views, and ICs: **P**, **NP**, ... , **undecidable**
 - hot research area in core CS (database community)



B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

A Neuroscientist's Information Integration Problem

Biomedical Informatics
Research Network
<http://abirn.net>



What is the cerebellar distribution of rat proteins with more than 70% homology with human NCS-1? Any structure specificity?
How about other rodents?

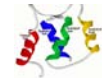
Information
Integration

"Complex
Multiple-Worlds"
Mediation

Inter-source links:
• unclear for the non-scientists
• hard for the scientist



protein localization
(NCMIR)



sequence info
(CaPROT)



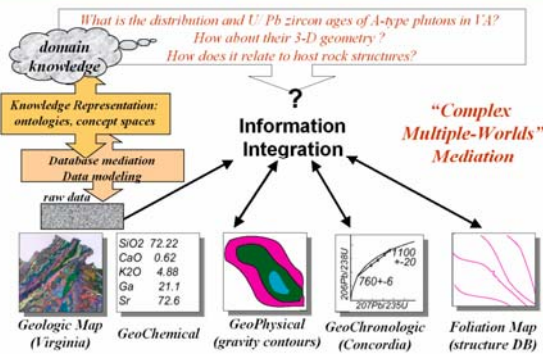
morphometry
(SYNAPSE)



neurotransmission
(SENSELAB)

B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

The Problem: Scientific Data Integration or: ... from Questions to Queries ...



CYBERINFRASTRUCTURE FOR THE GEOSCIENCES

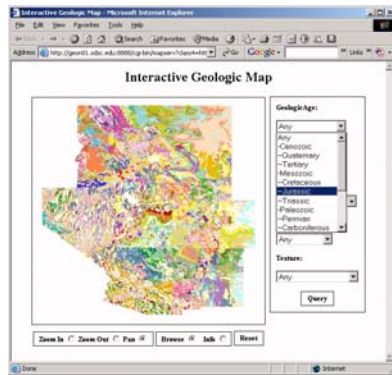
www.geongrid.org



Scientific Data Integration using Semantic Extensions

B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

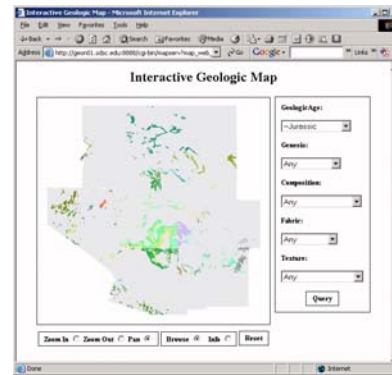
Querying by Geologic Age ...



B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

Synthesis
Semantics
Structure
Syntax
System Integ.

Querying by Geologic Age: Results

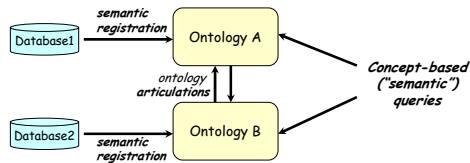


B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

Synthesis
Semantics
Structure
Syntax
System Integ.

Semantic Mediation (via "semantic registration" of schemas and ontology articulations)

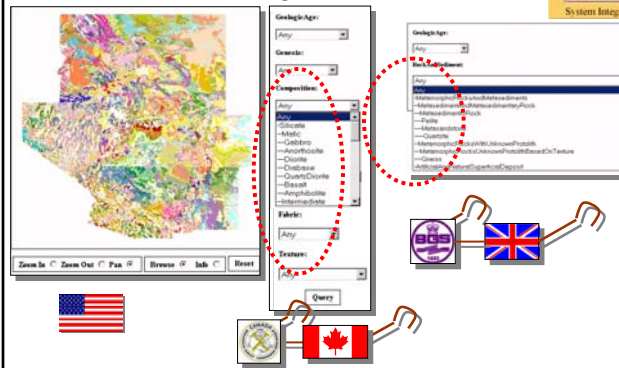
- Schema elements and/or data values are associated with concept expressions from the target ontology
→ conceptual queries "through" the ontology
- Articulation ontology
→ source registration to A, querying through B
- Semantic mediation: query rewriting w/ ontologies



B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

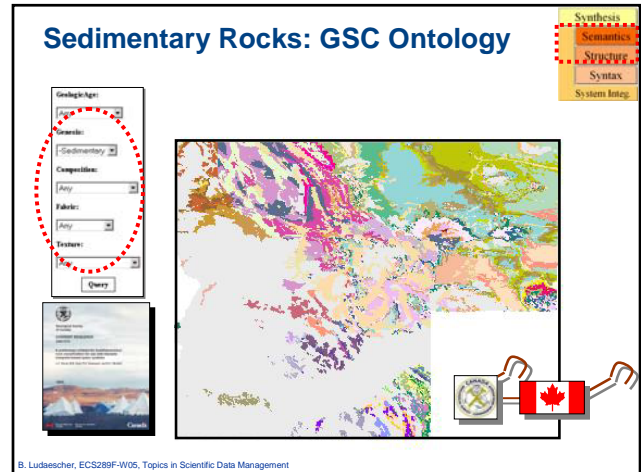
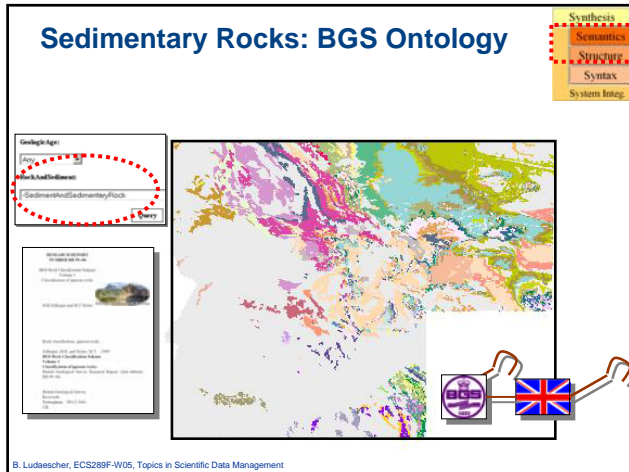
Synthesis
Semantics
Structure
Syntax
System Integ.

Different views on State Geological Maps



B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

Synthesis
Semantics
Structure
Syntax
System Integ.



Some Thoughts ...

- Translate this idea of multiple conceptual (ontology) views to **your domain!**
 - e.g. datasets \leftrightarrow biological pathways registration
- Your data is valuable (time & \$\$\$ spent in producing it)
 - data (re-)usability**
- Metadata helps to discover, localize, assess relevant data sets, given particular scientific questions & queries
- Does your system "understand" what to do with the metadata?
- Capturing more semantics** of a data set in a way that humans and systems can exploit it is an **investment in reusability**
 - "We are producing more and more data"
 - Today "we can store everything!"
 - But can we **use anything?** (i.e., is anyone looking at the data after the initial creation?)
- Design system, interfaces, data and metadata models with reusability in mind (think archives and "time capsules")
- This may even be pushed to the experiment/simulation/workflow design...

B. Ludaescher, ECS289F-W05, Topics in Scientific Data Management

Data Semantics and Ontologies should be useful for Humans and "The Machine"

Class MechanicalDeposition

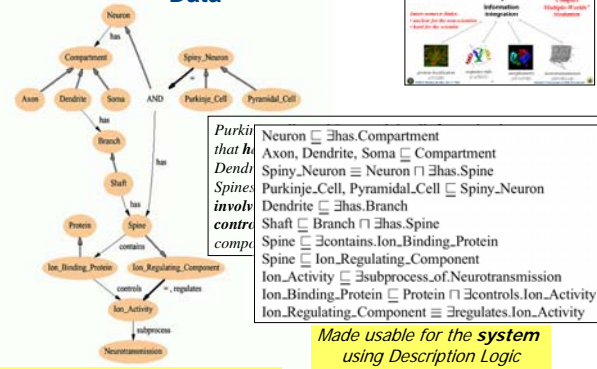
Mechanism of deposition: Mechanical deposition is the process of particles being transported by wind, water, or ice, and then settling to form a sediment. The particles are typically sand, silt, or clay, and the resulting sediment is called sand, silt, or clay. The process of deposition is often associated with the formation of dunes, sandbars, and other sedimentary features.

Direct Known Instances: [List of instances]

Properties: [List of properties]

B. Ludaescher, ECS289F-W05, Topics in Scientific Data Management

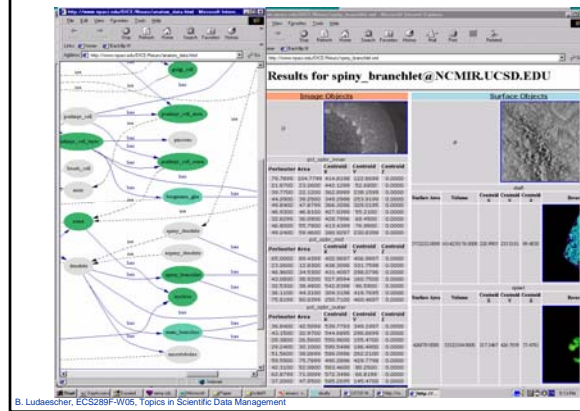
Example: Domain Knowledge to “glue” SYNAPSE & NCMIR Data



formalized as domain map/ontology

“Semantic Source Browsing”:

Domain Maps/Ontologies (left) and conceptually linked data (right)



B. Ludescher, ECS289F-W05, Topics in Scientific Data Management

A Semantic Mediation Result View

1. Interactive PROTEIN LOCALIZATION Query (UCSD/NCMIR data only)

Find the distribution of PROTEIN having ISOFORM in SPECIES and corresponding IMAGES

SPECIES = "chicken".
 protocol(PROTEIN, IZOPFORM, SPECIES, IMAGE_ID).

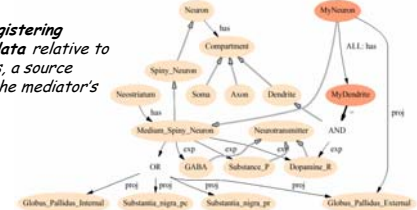
PROTEIN LOCALIZATION Results

PROTEIN	IZOPFORM	SPECIES	IMAGE_ID
Protein A	123456	chicken	Image 1
Protein B	234567	chicken	Image 2
Protein C	345678	chicken	Image 3
Protein D	456789	chicken	Image 4
Protein E	567890	chicken	Image 5

B. Ludescher, ECS289F-W05, Topics in Scientific Data Management

Source Contextualization through Ontology Refinement

In addition to **registering** ("hanging off") data relative to existing concepts, a source may also **refine** the mediator's domain map...



MyDendrite \sqsubseteq Dendrite \sqcap has.Dopamine_R
 MyNeuron \sqsubseteq Medium_Spiny_Neuron
 \sqcap has.Globus.pallidus.external
 \sqcap has.MyDendrite

⇒ sources can register new concepts at the mediator...
 ⇒ increase your data usability

B. Ludescher, ECS289F-W05, Topics in Scientific Data Management

Outline

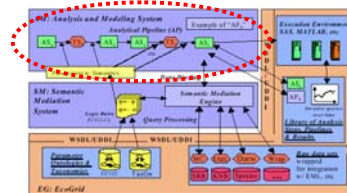
- Semantics & Scientific Data Integration
- Semantics & Scientific Workflow Management
- Conclusions

B. Ludäscher, ECS289F-W05, Topics in Scientific Data Management


- B. Ludaescher, ECS289F-W05, Topics in Scientific Data Management

[illegible]

- B. Ludaescher, ECS289F-W05, Topics in Scientific Data Management



Promoter Identification Workflow



Step 1
Microarray Analysis
microarray data

Step 2
Clusfavour Analysis
Gene ID

Step 3
GenBank sequence retrieval
cDNA sequence

Step 4
NCBI BLAST search
genomic sequence

Step 5
Transfac search

Step 6
Transcription factor binding
Promoter Identification

Step 7
Promoter Model generator
promoter data

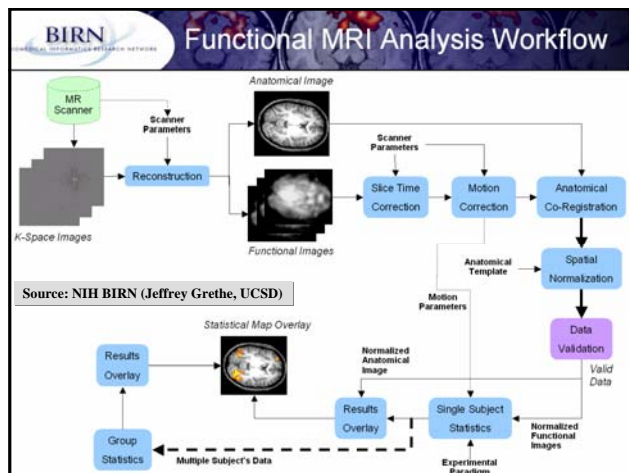
Step 8
NCBI BLAST search
Consensus sequence

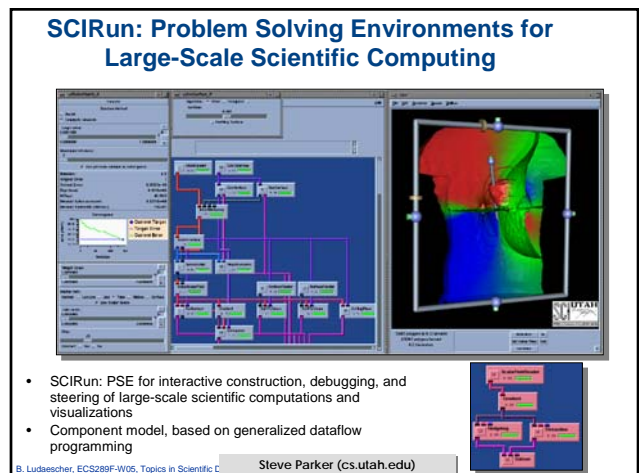
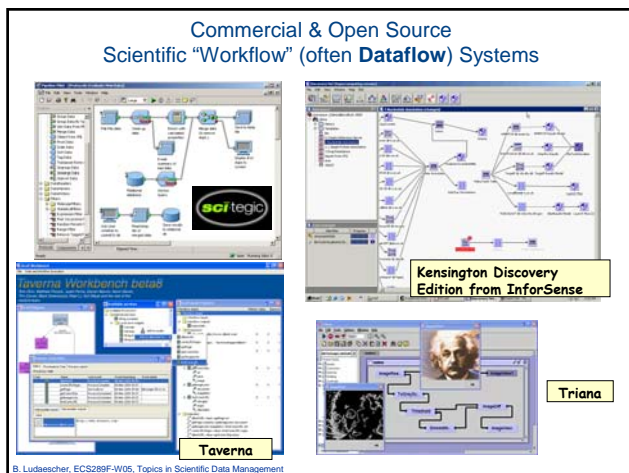
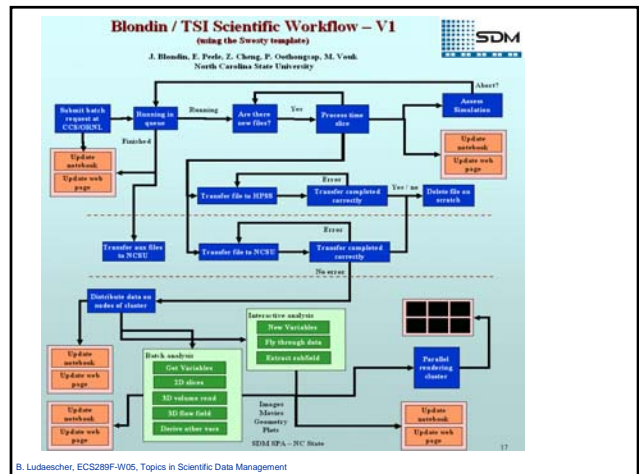
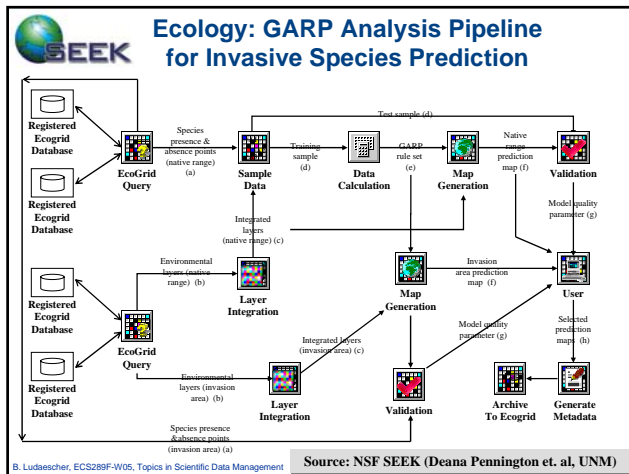
new candidate target genes

Source: Matt Coleman (LLNL)

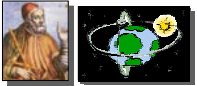
B. Ludaescher, ECS289F-W05, Topics in Scientific Data Management

- B. Ludaescher, ECS289F-W05, Topics in Scientific Data Management





Ptolemy II



Ptolemy II - Heterogeneous Modeling and Design in Java

The Ptolemy project studies modeling, simulation, and design of concurrent, real-time, embedded systems. The focus is on assembly of concurrent components. The key underlying principle in the project is the use of each model of computation in the interaction components.

read!

see!

try!

Source: Edward Lee et al. <http://ptolemy.eecs.berkeley.edu/ptolemyII/>

B. Ludaescher, ECS289F-W05, Topics in Scientific Data Management


Why Ptolemy II (and thus KEPLER)?

- Ptolemy II Objective:
 - “The focus is on **assembly of concurrent components**. The key underlying principle in the project is the use of **well-defined models of computation** that govern the interaction between components. A major problem area being addressed is the use of **heterogeneous mixtures of models of computation**.”
- Dataflow Process Networks w/ natural support for **abstraction, pipelining (streaming) actor-orientation, actor reuse**
- User-Orientation
 - Workflow design & exec console (Vergil GUI)
 - “Application/Glue-Ware”**
 - excellent modeling and design support
 - run-time support, monitoring, ...
 - not a middle-/underware** (we use someone else's, e.g. Globus, SRB, ...)
 - but middle-/underware is conveniently accessible through actors!
- PRAGMATICS
 - Ptolemy II is mature, continuously extended & improved, well-documented (500+pp)
 - open source system
 - Ptolemy II folks actively participate in KEPLER

B. Ludaescher, ECS289F-W05, Topics in Scientific Data Management

KEPLER/CSP: Contributors, Sponsors, Projects

(or loosely coupled Communicating Sequential Persons ;-)



www.kepler-project.org

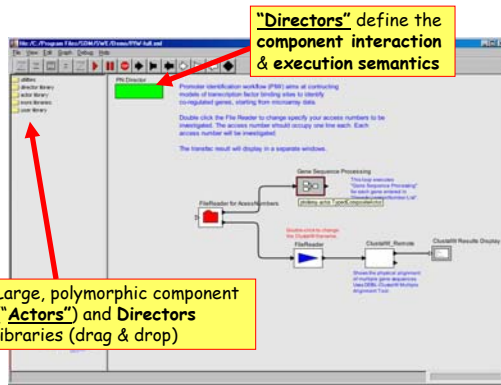
B. Ludaescher, ECS289F-W05, Topics in Scientific Data Management

KEPLER: An Open Collaboration

- Initiated by members from DOE SDM/SPA and NSF SEEK; now several other projects (GEON, Ptolemy II, EOL, Resurgence/NMI, ...)
- Open Source (BSD-style license)
- Intensive Communications:
 - Web-archived mailing lists
 - IRC (!)
 - Meetings, Hackathons
- Co-development:
 - via shared CVS repository
 - joining as a new co-developer (currently):
 - get a CVS account (read-only)
 - local development + contribution via existing KEPLER member
 - be voted “in” as a member/co-developer
- Software & social engineering
 - How to better accommodate new groups/communities?
 - How to better accommodate different usage/contribution models (core dev ... special purpose extender ... user)?

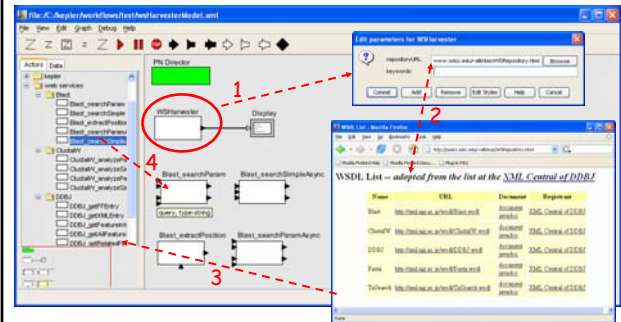
B. Ludaescher, ECS289F-W05, Topics in Scientific Data Management

Ptolemy II/KEPLER GUI (Vergil)



B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

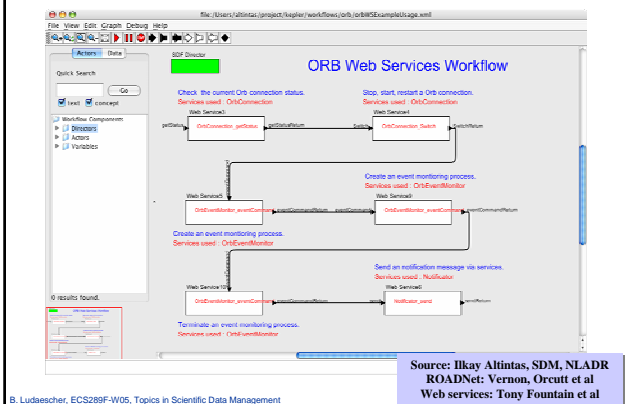
Web Services → Actors (WS Harvester)



→ "Minute-made" (MM) WS-based application integration
 • Similarly: MM workflow design & sharing w/o implemented components

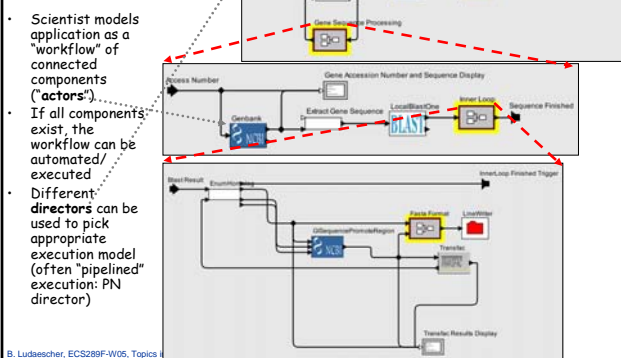
B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

Rapid Web Service-based Prototyping (Here: ROADNet Command & Control Services for LOOKING Kick-Off Mtg)



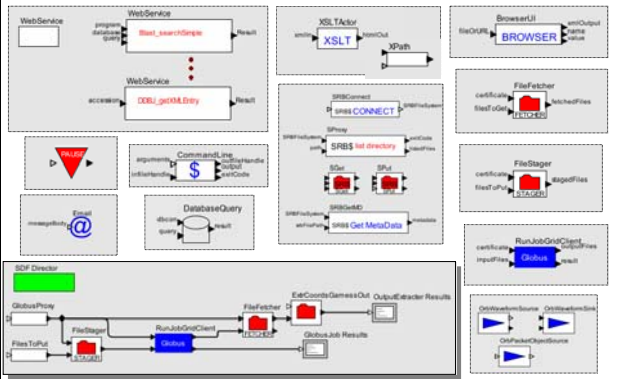
B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

An "early" example: Promoter Identification SSDBM, AD 2003



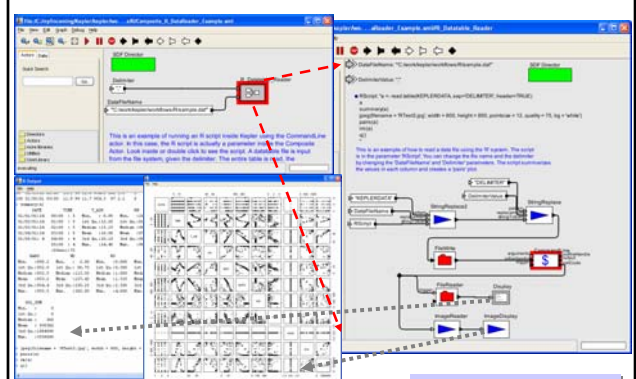
B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

Some Recent Actor Additions



B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

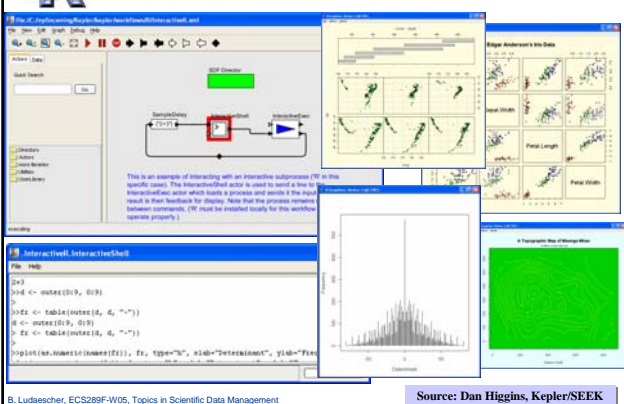
R in KEPLER (w/ editable script)



B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

Source: Dan Higgins, Kepler/SEEK

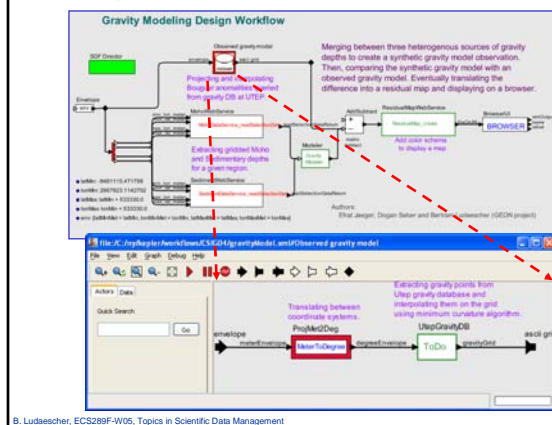
R in KEPLER (interactive session)



B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

Source: Dan Higgins, Kepler/SEEK

Blurring Design (ToDo) and Execution



B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

Some Scientific Workflow Challenges

- **Typical Features**
 - data-intensive and/or compute-intensive
 - plumbing-intensive (consecutive web services won't fit)
 - dataflow-oriented
 - distributed (remote data, remote processing)
 - user-interaction “in the middle”, ...
 - ... vs. (C-z; bg; fg)-ing (“detach” and reconnect)
 - advanced programming constructs (map(f), zip, takewhile, ...)
 - logging, provenance, “registering back” (intermediate) products...

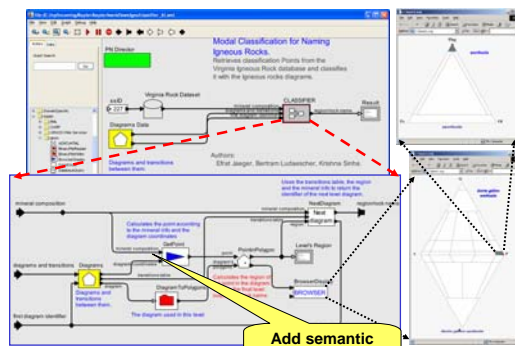
B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

Scientific Workflows & Semantics

- Registering data to ontologies: semantic types (in addition to structural data types)
- Smarter data set discovery & integration
- Now also:
 - Smarter workflow design
 - More “intelligent” (semantics-aware) component composition
 - Improved (re-)usability of data, services (actors), and workflows
 - Given semantic type of my input ports, what other data sets / actors produce such input

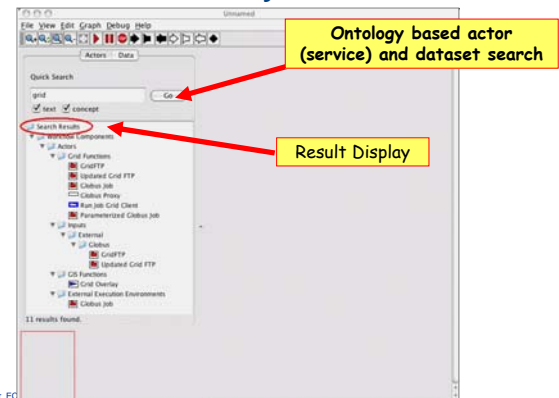
B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

Reengineering a Geoscientist's Mineral Classification Workflow



B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

Beginnings: Ontology-based Actor/Service Discovery



B. Ludascher, ECS289F-W05, Topics in Scientific Data Management

Semantics & Scientific Workflows

Data comes from **heterogeneous** sources

- Real-world observations
- Spatial-temporal contexts
- Collection/measurement protocols and procedures
- Many representations for the same information (count, area, density)
- Schematically heterogeneous

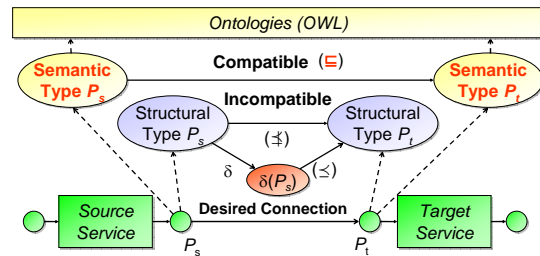
Data discovered and “synthesized” **manually**

Hard to reuse/repurpose existing analytical steps
(another form of heterogeneity)

B. Ludaescher, ECS289F-W05, Topics in Scientific Data Management

A KR+DI+Scientific Workflow Problem

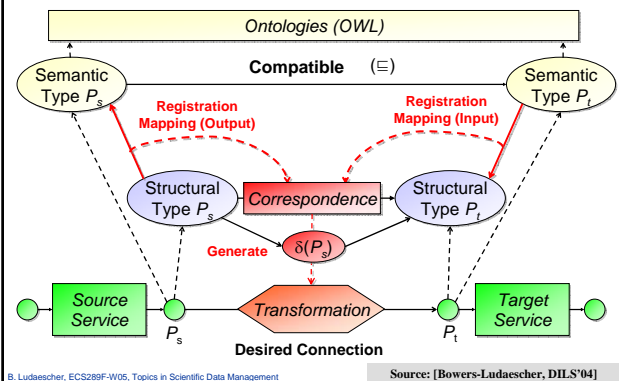
- **Services can be semantically compatible, but structurally incompatible**



B. Ludaescher, ECS289F-W05, Topics in Scientific Data Management

Source: [Bowers-Ludaescher, DILS'04]

Ontology-Informed Data Transformation (“Structure-Shim”)



B. Ludaescher, ECS289F-W05, Topics in Scientific Data Management

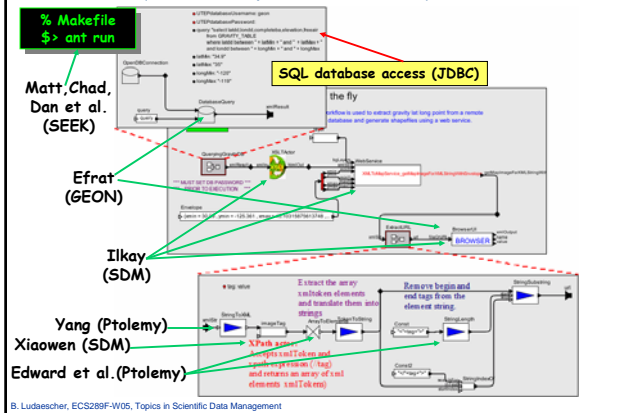
Source: [Bowers-Ludaescher, DILS'04]

Outline

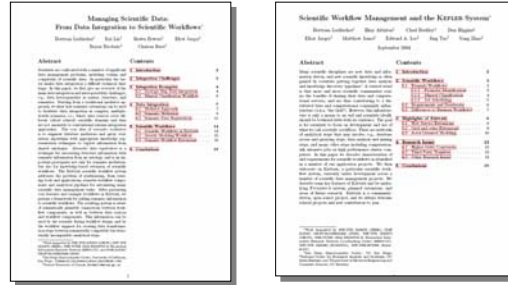
- Scientific Data Integration
- Scientific Workflow Management
- Musings & Conclusions

B. Ludaescher, ECS289F-W05, Topics in Scientific Data Management

GEON Dataset Generation & Registration



Further Reading



under review – available upon request from ludaesch@sdsc.edu

Related Publications

- **Semantic Data Registration and Integration**
 - [On Integrating Scientific Resources through Semantic Registration](#), S. Bowers, K. Lin, and B. Lüdäscher, 16th International Conference on Scientific and Statistical Database Management (*SSDBM04*), 21-23 June 2004, Santorini Island, Greece.
 - [A System for Semantic Integration of Geologic Maps via Ontologies](#), K. Lin and B. Lüdäscher. In *Semantic Web Technologies for Searching and Retrieving Scientific Data (SCISW)*, Sanibel Island, Florida, 2003.
 - [Towards a Generic Framework for Semantic Registration of Scientific Data](#), S. Bowers and B. Lüdäscher. In *Semantic Web Technologies for Searching and Retrieving Scientific Data (SCISW)*, Sanibel Island, Florida, 2003.
 - [The Role of XML in Mediated Data Integration Systems with Examples from Geological \(Map\) Data Interoperability](#), B. Brodicar, B. Lüdäscher, and K. Lin. In *Geological Society of America (GSA) Annual Meeting*, volume 35(6), November 2003.
 - [Semantic Mediation on Scientific Data Integration: A Case Study from the GEON Grid](#), K. Lin, B. Lüdäscher, B. Brodicar, D. Seber, C. Baru, and K. A. Sinha. In *Geological Society of America (GSA) Annual Meeting*, volume 35(6), November 2003.
- **Query Planning and Rewriting**
 - [Processing First-Order Queries under Limited Access Patterns](#), Alan Nash and B. Lüdäscher, Proc. 23rd ACM Symposium on Principles of Database Systems (*PODS04*) Paris, France, June 2004.
 - [Processing Unions of Conjunctive Queries with Negation under Limited Access Patterns](#), Alan Nash and B. Lüdäscher, 9th Intl. Conference on Extending Database Technology (*EDBT04*) in Athens, Greece, October 2004.
 - [Web Service Composition through Declarative Queries: The Case of Conjunctive Queries with Union and Negation](#), B. Lüdäscher and Alan Nash, Research abstract (poster), 20th Intl. Conference on Data Engineering (*ICDE04*) Boston, IEEE Computer Society, April 2004.

Related Publications

- **Scientific Workflows**
 - **Kepler: An Extensible System for Design and Execution of Scientific Workflows**, I. Altintas, C. Berkley, E. Jaeger, M. Jones, B. Ludäscher, S. Mock, *16th International Conference on Scientific and Statistical Database Management (SSDBM'04)*, 21-23 June 2004, Santorini Island, Greece.
 - **Kepler: Towards a Grid-Enabled System for Scientific Workflows**, Ilkay Altintas, Chad Berkley, Efraim Jaeger, Matthew Jones, Bertram Ludäscher, Steve Mock, *Workflow in Grid Systems (GGF10)*, Berlin, March 9th, 2004.
 - **An Open-Driven Framework for Data Transformation in Scientific Workflows**, S. Bowers and B. Ludäscher, *Intl. Workshop on Data Integration in the Life Sciences (DILS'04)*, March 25-26, 2004 Leipzig, Germany, LNCS 2994.
 - **A Web Service Composition and Deployment Framework for Scientific Workflows**, I. Altintas, E. Jaeger, K. Lin, B. Ludäscher, A. Memon, *In the 2nd Intl. Conference on Web Services (ICWS)*, San Diego, California, July 2004.