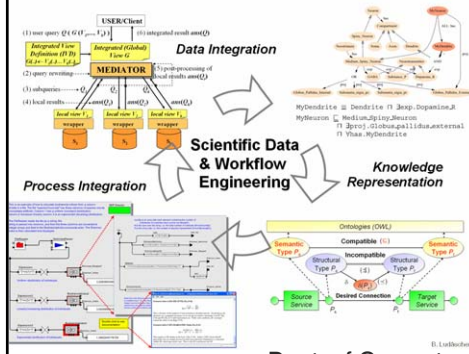


Scientific Data Integration: From the Big Picture to some Gory Details



Bertram Ludäscher
ludaesch@ucdavis.edu

Associate Professor
Dept. of Computer Science & Genome Center
University of California, Davis

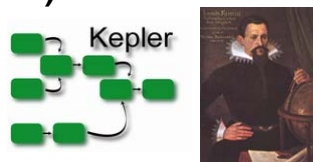
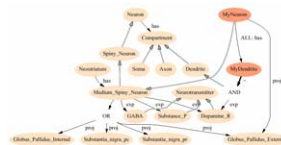
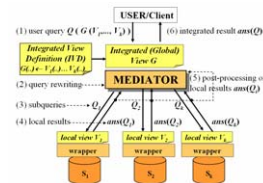


Fellow
San Diego Supercomputer Center
University of California, San Diego



Outline

- Data Integration & Mediation (DI)
- Challenges with Scientific Data
- Knowledge-based Extensions & Ontologies (DI+KR)
- Scientific Workflows (SWF+DI+KR)
- Scientific Workflow Design

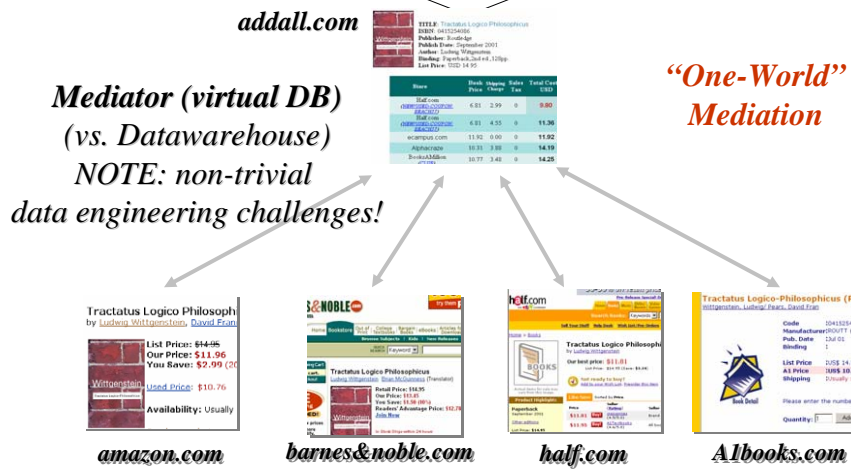


Bertram Ludäscher, UC DAVIS

An Online Shopper's Information Integration Problem



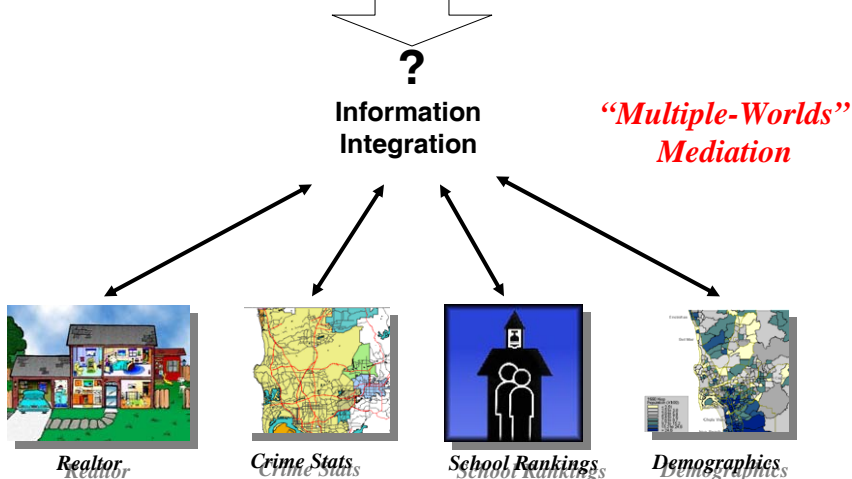
El Cheapo: "Where can I get the cheapest copy (including shipping cost) of Wittgenstein's Tractatus Logicus-Philosophicus within a week?"



A Home Buyer's Information Integration Problem



What houses for sale under \$500k have at least 2 bathrooms, 2 bedrooms, a nearby school ranking in the upper third, in a neighborhood with below-average crime rate and diverse population?



A Neuroscientist's Information Integration Problem

Biomedical Informatics
Research Network
<http://nbirn.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

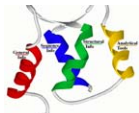
“Complex Multiple-Worlds” Mediation

Inter-source links:

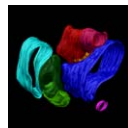
- unclear for the non-scientists
- hard for the scientist



protein localization
(NCMIR)



sequence info
(CaPROT)



morphometry
(SYNAPSE)



neurotransmission
(SENSELAB)

Bertram Ludäscher, UC DAVIS

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

What is the distribution and U/Pb zircon ages of A-type plutons in VA?
How about their 3-D geometry?
How does it relate to host rock structures?

?

Information Integration

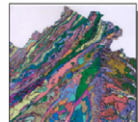
“Complex Multiple-Worlds” Mediation

domain knowledge

Knowledge Representation:
ontologies, concept spaces

Database mediation
Data modeling

raw data



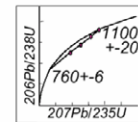
Geologic Map
(Virginia)

SiO ₂	72.22
CaO	0.62
K ₂ O	4.88
Ga	21.1
Sr	72.6

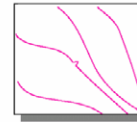
GeoChemical



GeoPhysical
(gravity contours)



GeoChronologic
(Concordia)



Foliation Map
(structure DB)

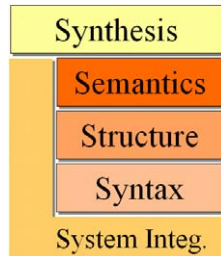


CYBERINFRASTRUCTURE FOR THE GEOSCIENCES

www.geongrid.org



Interoperability & Integration Challenges

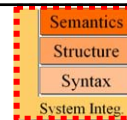


- reconciling S^5 heterogeneities
- “gluing” together resources
- bridging information and knowledge gaps computationally

- **System aspects: “Grid” Middleware**
 - distributed data & computing, SOA
 - resource discovery, authentication, authorization
 - web services, WSDL/SOAP, WSRF, OGSA, ...
 - *(re-)sources = services, files, data sets, nodes ...*
- **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, ...*

Bertram Ludäscher, UC DAVIS

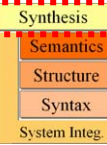
Information Integration Challenges: S^4 Heterogeneities



- **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 technologies**
 - + XML-based data exchange, integrated views, transparent query rewriting, ...
- **Semantics**
 - descriptive metadata, different terminologies, “hidden” semantics (context), implicit assumptions, ...
 - ➔ **Knowledge representation & semantic mediation technologies**
 - + “smart” data discovery & integration
 - + e.g. ask about **X** (*‘mafic’*); find data about **Y** (*‘diorite’*); be happy anyways!

Bertram Ludäscher, UC DAVIS

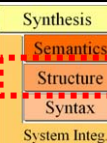
Information Integration Challenges: *S⁵ Heterogeneities*



- **Synthesis** of applications, analysis tools, data & query components, ... into “**scientific workflows**”
 - How to put together components to solve a scientist’s problem?
- ➔ **Scientific Problem Solving Environments (PSEs)**
 - ➔ **Portals, Workbench (“scientist’s view”)**
 - + ontology-enhanced data registration, discovery, manipulation
 - + creation and registration of new data products from existing ones, ...
 - ➔ **Scientific Workflow System (“engineer’s view”)**
 - + 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, ...

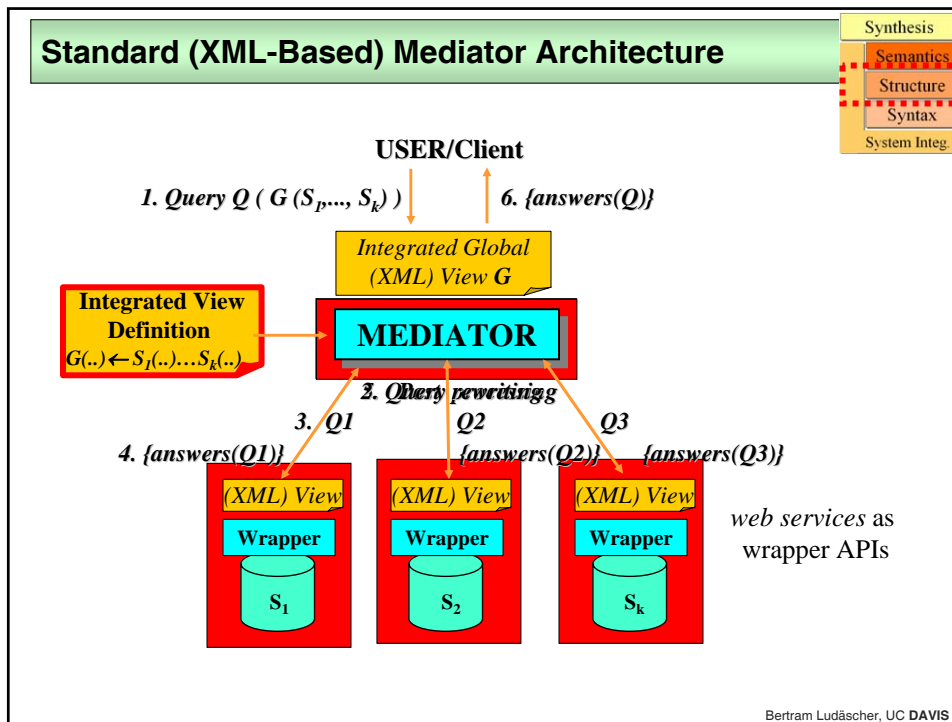
Bertram Ludäscher, UC DAVIS

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 $G(S_1, \dots, S_k)$

Bertram Ludäscher, UC DAVIS



Query Planning in Data Integration

Synthesis
Semantics
Structure
Syntax
System Integ.

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

Bertram Ludäscher, UC DAVIS

Data Integration: Limited Access Patterns+Views+ICs



We study the problem of rewriting a query Q in terms of a given set of views \mathcal{V} with limited access patterns \mathcal{P} , under a set Σ of integrity constraints. More precisely, we are interested in determining whether there exists a query plan Q' , expressed in terms of the views \mathcal{V} only, that is executable (i.e., observes \mathcal{P}) and equivalent to Q for all databases satisfying Σ . We say that Q is *feasible* if such Q' exists. For infeasible Q we seek the minimal containing and maximal contained executable queries, which provide the “best possible” executable query plans for approximating the answer to Q from above and below.

ICDT'05

Rewriting Queries Using Views with Access Patterns Under Integrity Constraints*

Alin Deutsch¹, Bertram Ludäscher², and Alan Nash³

University of California, San Diego
deutsch@cs.ucsd.edu, ludasch@sdsc.edu, anash@math.ucsd.edu

Abstract. We study the problem of rewriting queries using views in the presence of access patterns, integrity constraints, disjunction, and negation. We provide asymptotically optimal algorithms for finding minimal containing and maximal contained rewritings and for deciding whether an exact rewriting exists. We show that rewriting queries using views in this case reduces (a) to rewriting queries with access patterns and constraints without views and also (b) to rewriting queries using views under constraints without access patterns. We show how to solve (a) directly and how to reduce (b) to rewriting queries under constraints only (semantic optimization). These reductions provide two separate routes to a unified solution for all three problems, based on an extension of the relational chase theory to queries and constraints with disjunction and negation. We also handle equality and arithmetic comparison.

the set of listed publishers $L(p)$, repository $R^*(a, t)$, ACM anthology $A^*(a, t, c)$, and DBLP conference article $D^*(a, t, c)$. The relation symbols are annotated with access patterns, indicating which arguments must be given as inputs (marked \uparrow) and which ones can be retrieved as outputs (marked \circ) when accessing the relation. For example $C^{\circ(a, t)}$ means that an author a has to be given as input before one can retrieve the titles t of a's conference publications from $C(a, t)$.

Let Q be the query which asks for pairs of authors and titles of conference publications, journal publications, and magazines which are not PC-magazines:

Query Q

- $Q(a, t) \leftarrow C(a, t)$ (1)
- $Q(a, t) \leftarrow J(a, t)$ (2)
- $Q(a, t) \leftarrow M(a, t) \wedge \neg P(a, t, p), L(p)$ (3)

Q cannot be executed since no underlined literal is answerable e.g., the access patterns require a to be bound before involving $C(a, t)$ but no such binding is available. Worse yet, Q is not even feasible, i.e., there is no executable query Q' equivalent to Q . However, if the following set Σ of integrity constraints is given, an answerable Q' can be found that is equivalent under Σ :

ICs Σ

- View $C(a, t) \rightarrow \exists c D(a, t, c)$ (4)
- View $J(a, t) \rightarrow \exists p R(a, t) \wedge \neg P(a, t, p) \wedge L(p) \vee \exists c A(a, t, c), D(a, t, c)$ (5)
- View $M(a, t) \rightarrow \exists p \neg P(a, t, p), L(p)$ (6)

Constraint (4) states that every conference publication is a DBLP conference publication; (5) states that every journal publication is available from a repository, comes from a listed publisher and is not a PC magazine, or is available from the ACM anthology and from DBLP; and (6) states that magazine articles are not PC-magazine articles. We are only interested in databases which satisfy these constraints Σ . On these databases, Q is equivalent to Q^2 , obtained by “chasing” Q with Σ :

Chase result Q^2

- $Q^2(a, t) \leftarrow C(a, t), D(a, t, c)$
- $Q^2(a, t) \leftarrow J(a, t), R(a, t), \neg P(a, t, p), L(p)$
- $Q^2(a, t) \leftarrow J(a, t), A(a, t, c), D(a, t, c)$
- $Q^2(a, t) \leftarrow M(a, t), \neg P(a, t, p), L(p)$

Again, unanswerable literals are underlined. The answerable part $\text{ans}(Q^2)$ is obtained (roughly) by removing unanswerable parts:

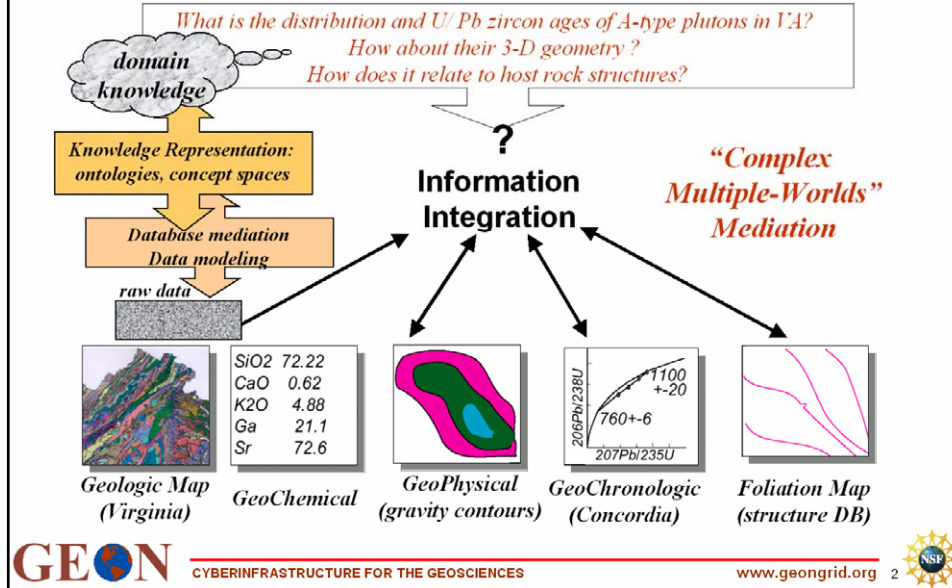
Answerable part $\text{ans}(Q^2)$

- $\text{ans}(Q^2)(a, t) \leftarrow C(a, t), D(a, t, c)$ (7)
- $\text{ans}(Q^2)(a, t) \leftarrow J(a, t), R(a, t)$ (8)
- $\text{ans}(Q^2)(a, t) \leftarrow J(a, t), D(a, t, c)$ (9)
- $\text{ans}(Q^2)(a, t) \leftarrow M(a, t)$ (10)

Scientific Data Integration using Semantic Extensions



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

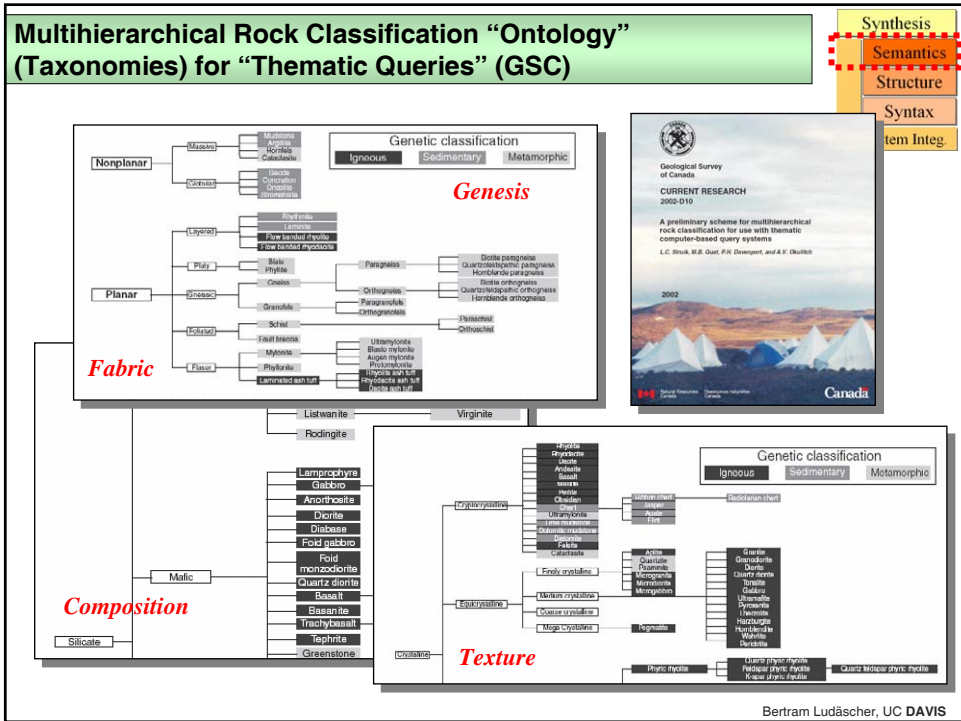
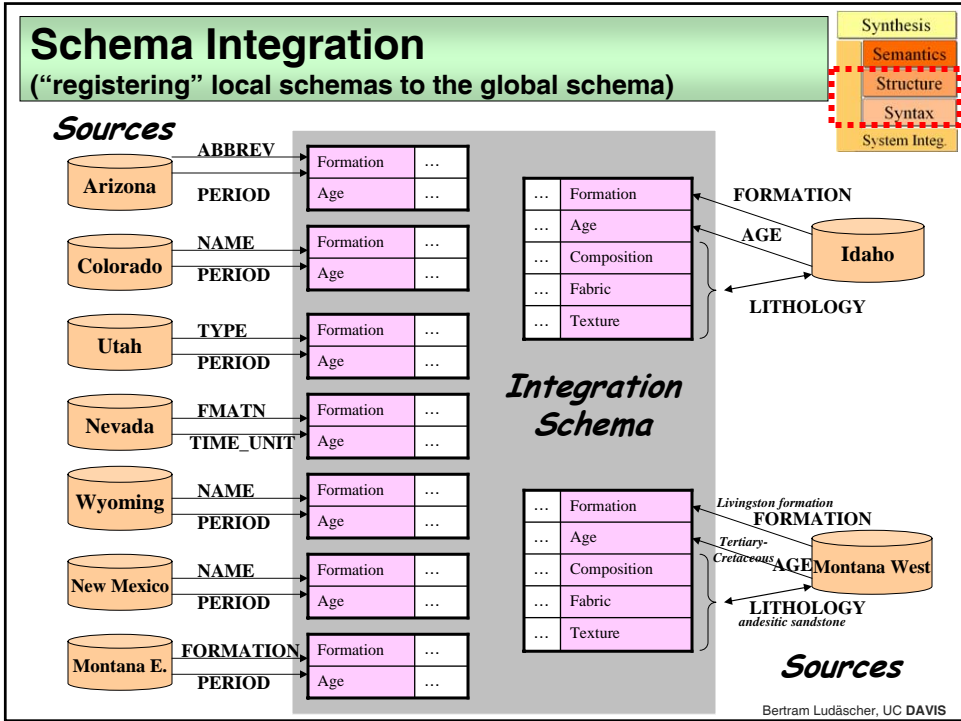


- **Given:**

- Geologic maps from different state geological surveys (shapefiles w/ different data schemas)
- Different ontologies:
 - Geologic age ontology (e.g. USGS)
 - Rock classification ontologies:
 - Multiple hierarchies (chemical, fabric, texture, genesis) from Geological Survey of Canada (GSC)
 - Single hierarchy from British Geological Survey (BGS)

- **Problem:**

- Support **uniform queries** across all maps
- ... using **different** ontologies
- Support **registration** w/ ontology A, **querying** w/ ontology B



Ontology-Enabled Application Example: Geologic Map Integration

- Synthesis
- Semantics
- Structure
- Syntax
- System Integ.

The screenshot shows a web browser window displaying a geologic map of Nevada. On the left, a 'domain knowledge' cloud contains a hierarchical list of geological terms such as 'Carboniferous', 'Permian', 'Triassic', and 'Jurassic'. A large orange arrow labeled 'Knowledge representation Geologic Age ONTOLOGY' points from this list to the map. A callout box over the map says '+/- a few hundred million years'. On the right, a 'GeologicAge:' control panel features a dropdown menu with 'Jurassic' selected, and a 'Query' button. The browser's address bar shows a URL from 'http://geon01.sdsc.edu:8888/...'. The name 'Bertram Ludäscher, UC DAVIS' is visible in the bottom right corner.

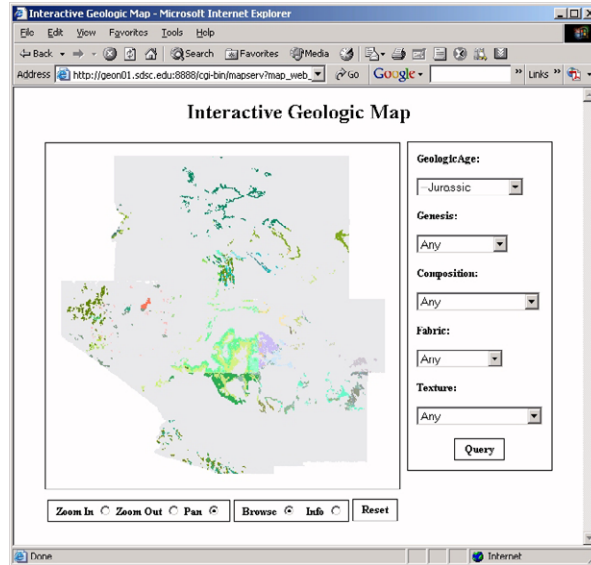
Querying by Geologic Age ...

- Synthesis
- Semantics
- Structure
- Syntax
- System Integ.

This screenshot shows the 'Interactive Geologic Map' application running in a Microsoft Internet Explorer browser. The main window displays a colorful geologic map of Nevada. To the right of the map is a control panel with a 'GeologicAge:' dropdown menu where 'Jurassic' is selected. Below the dropdown is a 'Texture:' dropdown menu set to 'Any' and a 'Query' button. At the bottom of the map area, there are buttons for 'Zoom In', 'Zoom Out', 'Pan', 'Browse', 'Info', and 'Reset'. The browser's address bar shows the URL 'http://geon01.sdsc.edu:8888/cg-bin/mapserv?class4=ht...'. The name 'Bertram Ludäscher, UC DAVIS' is visible in the bottom right corner.

Querying by Geologic Age: Results

Synthesis
Semantics
Structure
Syntax
System Integ.

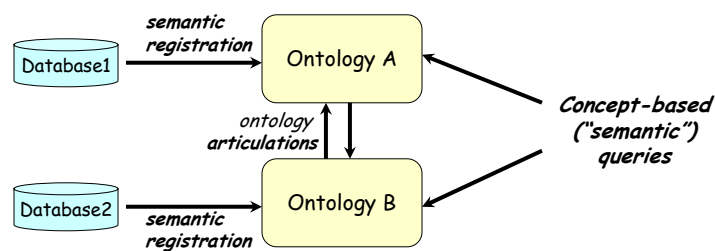


Bertram Ludäscher, UC DAVIS

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

Synthesis
Semantics
Structure
Syntax
System Integ.

- **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**



Bertram Ludäscher, UC DAVIS

Different views on State Geological Maps

Synthesis
 Semantics
 Structure
 Syntax
 System Integ.

Geologic Age: Any

Genesis: Any

Composition: Any

- Silicite
- Mafic
- Gabbro
- Anorthosite
- Diorite
- Diabase
- Quartz Diorite
- Basalt
- Amphibolite
- Intermediate

Fabric: Any

Texture: Any

Query

Geologic Age: Any

RockAndSediment: Any

- MetamorphicRocksAndMetasediments
- MetasedimentAndMetasedimentaryRock
- MetasedimentaryRock
- Pella
- Metasandstone
- Quartzite
- MetamorphicRocksWithUnknownProtolith
- MetamorphicRocksUnknownProtolithBasedOnTexture
- Gneiss
- ArtificialAndNaturalSuperficialDeposit

Zoom In Zoom Out Pan Browse Info Reset

USGS GSC

Bertram Ludäscher, UC DAVIS

Sedimentary Rocks: BGS Ontology

Synthesis
 Semantics
 Structure
 Syntax
 System Integ.

Geologic Age: Any

RockAndSediment: SedimentAndSedimentaryRock

Query

RESEARCH REPORT
 NUMBER RR 99-06
 BGS Rock Classification Scheme
 Volume 1
 Classification of igneous rocks
 M.E. Gillispie and M.T. Styles

Rock classification, igneous rocks
 Gillispie M.E. and Styles M.T. 1999
 BGS Rock Classification Scheme
 Volume 1
 Classification of igneous rocks.
 British Geological Survey Research Report, (last revised)
 RR 99-06.

British Geological Survey
 Keyworth
 Nottingham NG12 5GG
 UK

BGS UK

Bertram Ludäscher, UC DAVIS

Sedimentary Rocks: GSC Ontology

- Synthesis
- Semantics
- Structure
- Syntax
- System Integ.

GeologicAge: Any

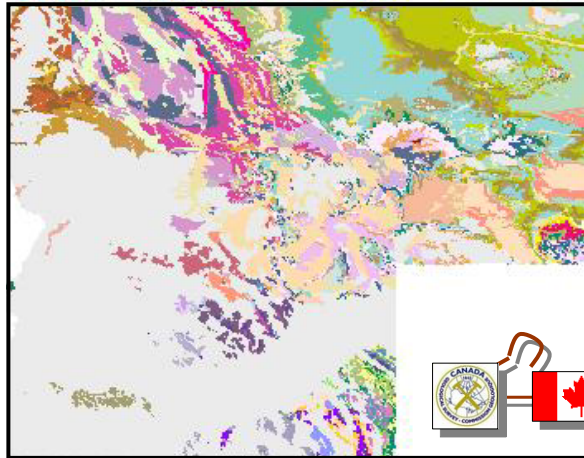
Genesis: Sedimentary

Composition: Any

Fabric: Any

Texture: Any

Query



Bertram Ludäscher, UC DAVIS

Implementation in OWL: Not only “for the machine” ...



Class Intrusive

Intrusive igneous rocks are formed from magma that cools and solidifies deep beneath the Earth's surface. The solidifying effect of the surrounding rock allows the magma to solidify very slowly.

Class Sedimentary

Sedimentary rocks are formed from loose particles or pieces of rock that are buried and compacted over time. They form from the gradual accumulation of sediments on the Earth's surface. Sedimentary rocks often have distinct layering or bedding. Many of the most important types of the most common sedimentary rocks are listed below.

Direct Known Subclasses:

- [Sedimentary](#)
- [Glacial](#)
- [Subaerial](#)
- [Subaqueous](#)

Class MechanicalDeposition

Mechanical deposition sedimentary rocks are made up of pieces of pre-existing rocks. Pieces of rock are loosened by weathering, then transported to some basin or depression where sediment is trapped. If the sediment is buried deeply, it becomes compacted and cemented, forming sedimentary rock. These sedimentary rocks may have particles ranging in size from microscopic clay to huge boulders. Their names are based on their clast or grain size. The smallest grains are called clay, then silt, then sand. Grains larger than 2 millimeters are called pebbles. Shale is a rock made mostly of clay, siltstone is made up of silt-sized grains, sandstone is made of sand-sized clasts, and conglomerate is made of pebbles surrounded by a matrix of sand or mud.

© US Geological Survey

Direct Known Supclasses:

- [Sedimentary](#)

Direct Known Subclasses:

- [Glacial](#)
- [Subaerial](#)
- [Subaqueous](#)

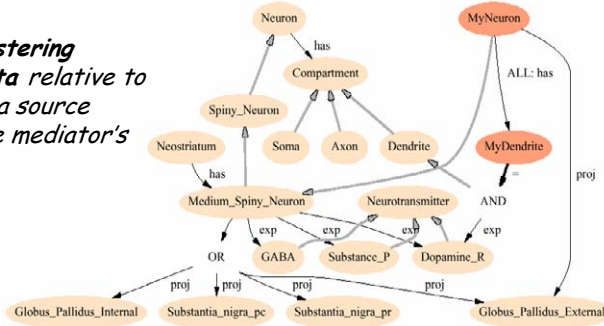
Properties:

Bertram Ludäscher, UC DAVIS

Source Data 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...*



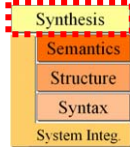
$MyDendrite \equiv Dendrite \sqcap \exists exp.Dopamine_R$
 $MyNeuron \sqsubseteq Medium_Spiny_Neuron$
 $\sqcap \exists proj.Globus_pallidus_external$
 $\sqcap \forall has.MyDendrite$

*⇒ sources can register **new concepts** at the mediator ...*



Scientific Workflows

Motivation: Scientific Workflows, Pre-Cyberinfrastructure



- **Data Federation & Grid “Plumbing”:**
 - access, move, replicate, query ... data (**Data-Grid**)
 - authenticate ... SRB Sget/Sput ... OPeNDAP, ... Antelope/ORBs
 - schedule, launch, monitor jobs (**Compute-Grid**)
 - Globus, Condor, Nimrod, APST, ...
- **Data Integration:**
 - Conceptual querying & integration, structure & semantics, e.g. mediation w/ SQL, XQuery + OWL (*Semantics-enabled Mediator*)
- **Data Analysis, Mining, Knowledge Discovery:**
 - manual/textbook (e.g. ternary diagrams), Excel, R, simulations, ...
- **Visualization:**
 - 3-D (volume), 4-D (spatio-temporal), n-D (conceptual views) ...



- **one-of-a-kind custom apps., detached (island) solutions**
- workflows are **hard to reproduce, maintain**
- **no/little** workflow design, automation, reuse, documentation
- need for an **integrated scientific workflow environment**

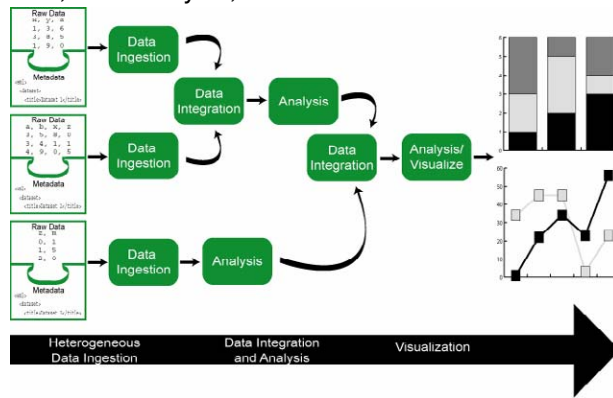
Bertram Ludäscher, UC DAVIS

What is a Scientific Workflow (SWF)?



- **Model the way scientists work with their data and tools**
 - Mentally coordinate data export, import, analysis via software systems
- **Scientific workflows emphasize data flow (≠ business workflows)**
- **Metadata** (incl. provenance info, semantic types etc.) is crucial for automated data ingestion, data analysis, ...

- **Goals:**
 - SWF automation,
 - SWF & component reuse,
 - SWF design & documentation
 - making scientists' data analysis and management easier!



Bertram Ludäscher, UC DAVIS

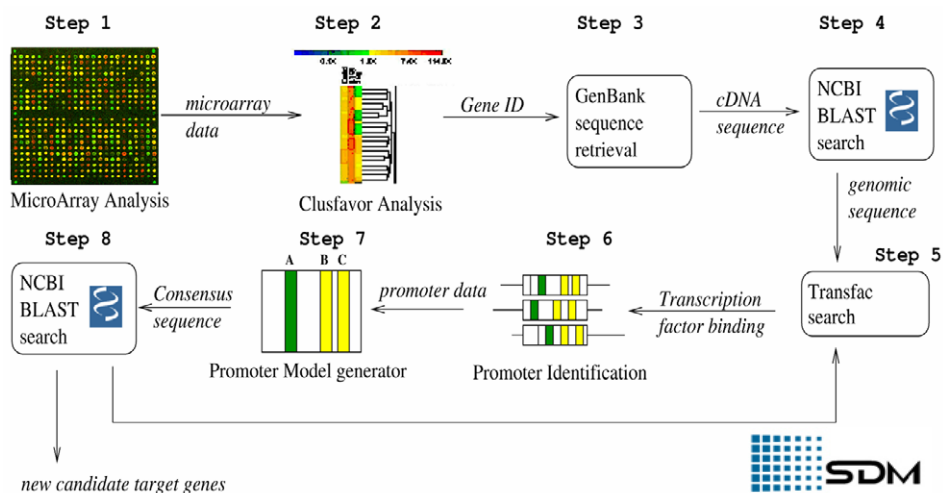
Some Scientific Workflow Features



- Typical requirements/characteristics:
 - data-intensive and/or compute-intensive
 - plumbing-intensive
 - dataflow-oriented
 - distribution (data, 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
 - ...
- ... easy to recognize a SWF when you see one!

Bertram Ludäscher, UC DAVIS

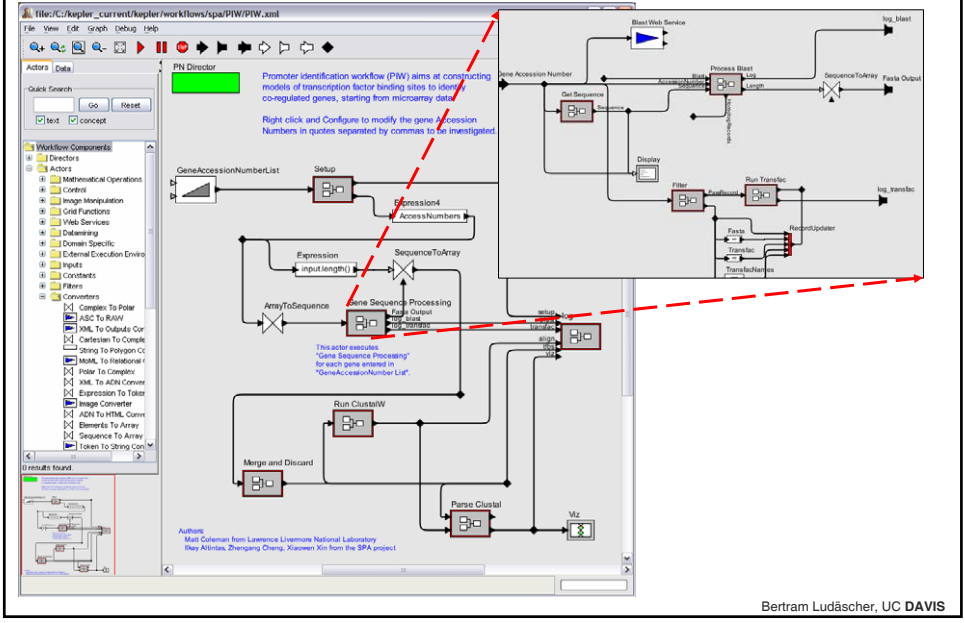
Promoter Identification Workflow (Napkin Drawing)



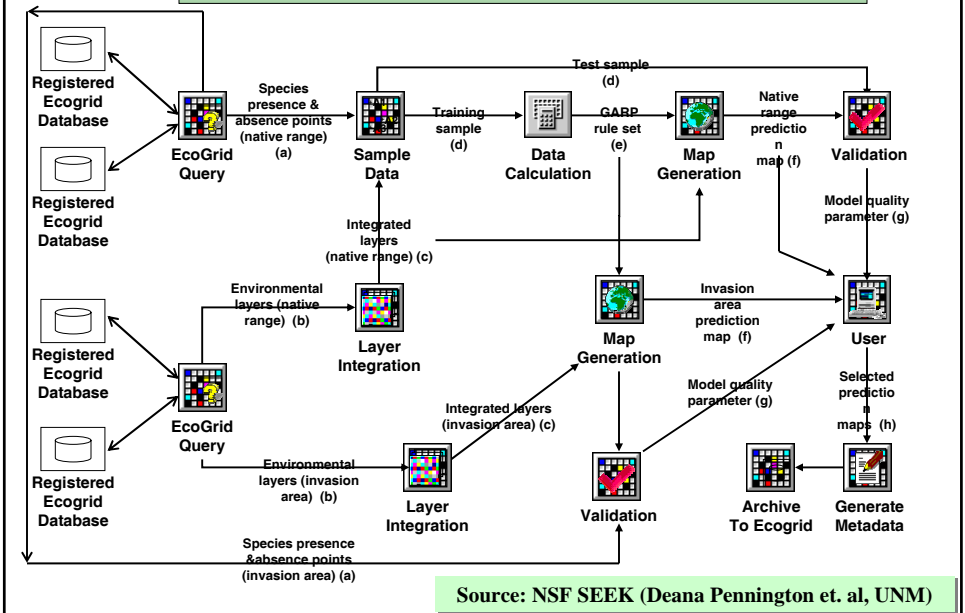
Source: Matt Coleman (LLNL)

Bertram Ludäscher, UC DAVIS

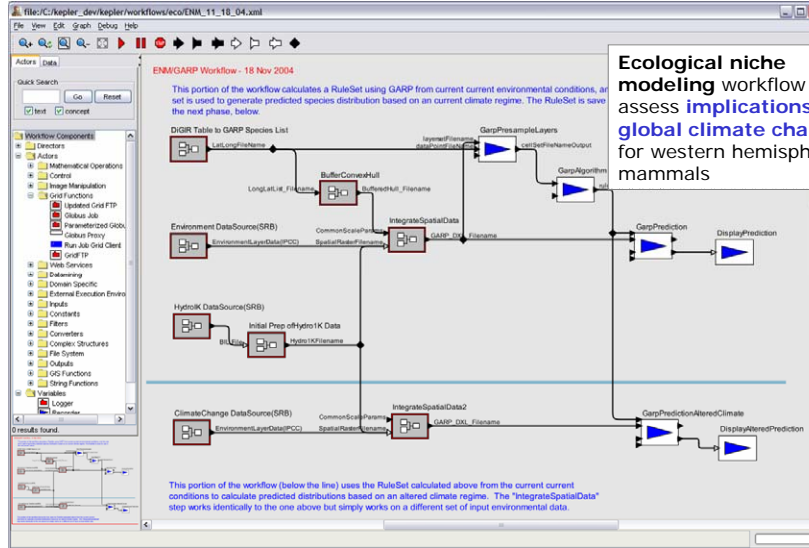
Promoter Identification Workflow in Kepler



Ecology: Invasive Species Prediction (Napkin Drawing)

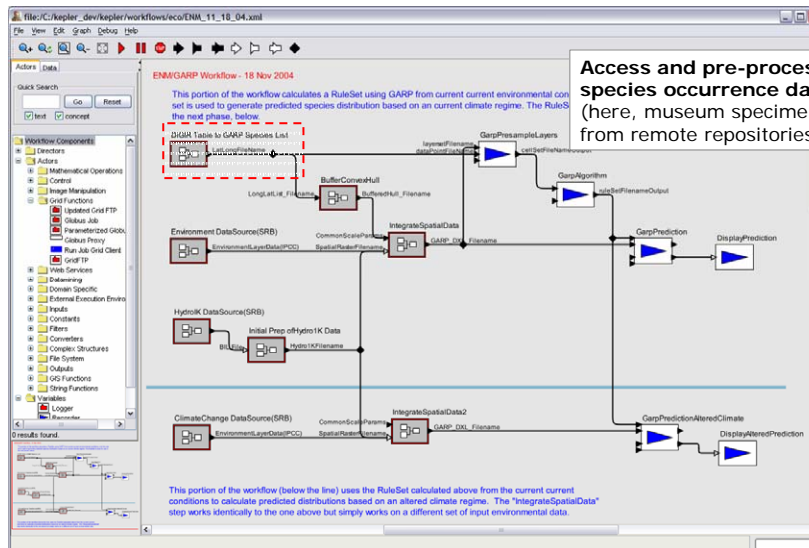


Ecological Niche Modeling in Kepler



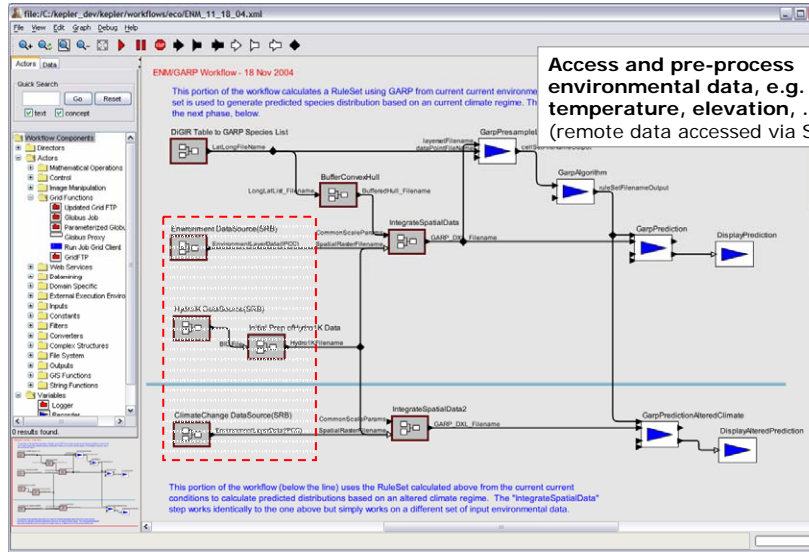
Bertram Ludäscher, UC DAVIS

Ecological Niche Modeling in Kepler



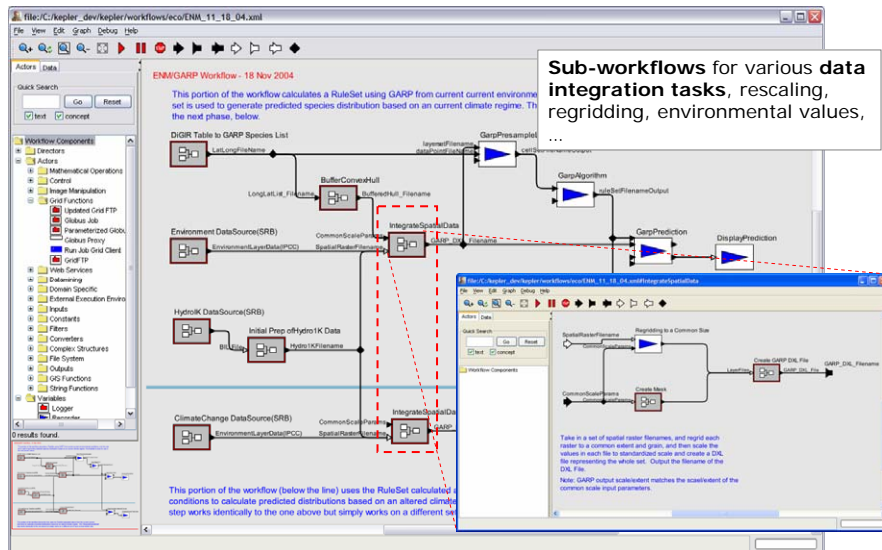
Bertram Ludäscher, UC DAVIS

Ecological Niche Modeling in Kepler



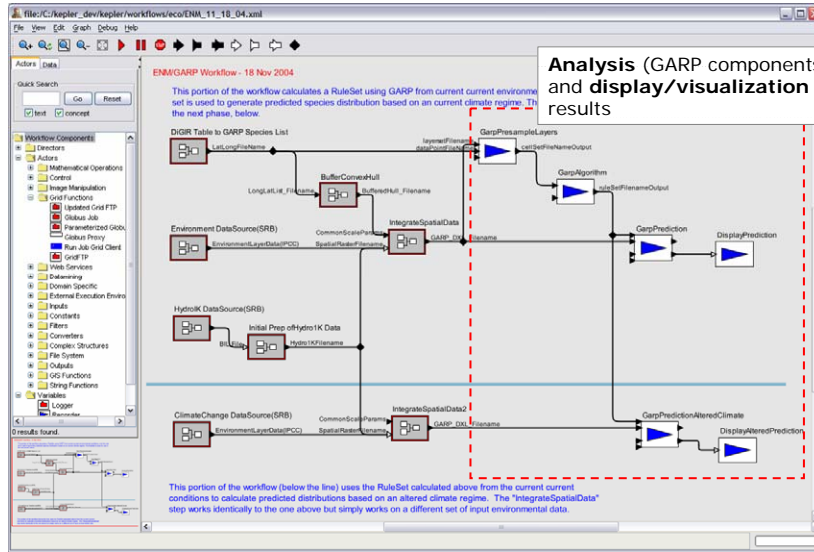
Bertram Ludäscher, UC DAVIS

Ecological Niche Modeling in Kepler



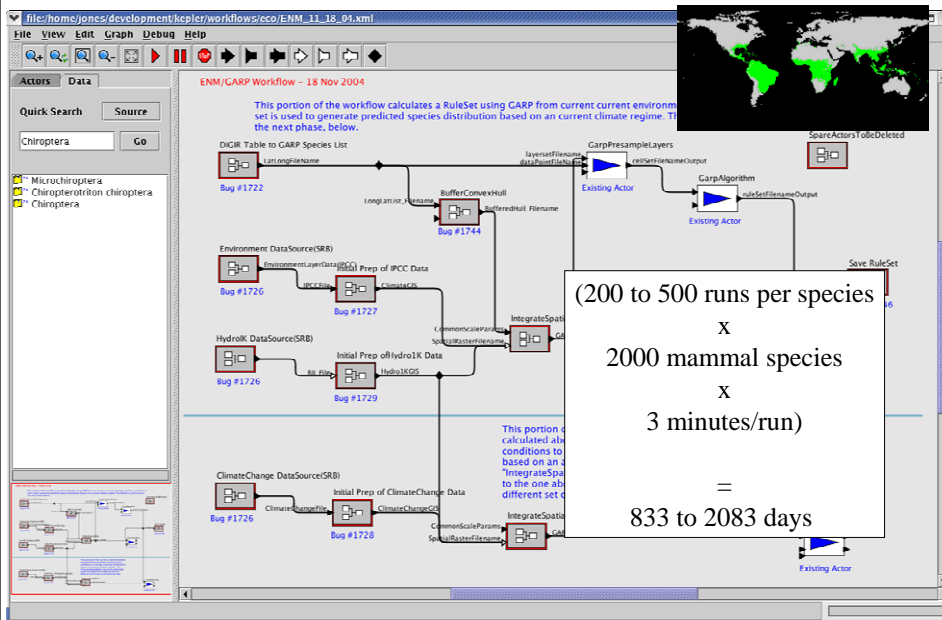
Bertram Ludäscher, UC DAVIS

Ecological Niche Modeling in Kepler



Bertram Ludäscher, UC DAVIS

Ecological Niche Modeling in Kepler



Query Builder

A simple example of using EML data. First, a search is done in the Data pane to locate an EML-described data set, which is dragged onto the workflow canvas. The EML data source is added to the workflow, and then it contacts the EcoGrid server to download the data and configure the ports. After being configured, it displays the ports from the EML data source, which are then mapped into an XY scatterplot.

Available Table Schemas:

Table Name	Fields	Data Type	Operator	Criteria
Datos Meteorologicos	BARO	DOUBLE	GT	GREATER THAN 9.53
Datos Meteorologicos	SEW	DOUBLE	LT	

Query Configuration:

SELECT: Where

AND

Datos Meteorologicos (BARO GREATER THAN 9.53)

Control: Add AND, Add OR, Add Condition, Remove

Table: Datos Meteorologicos, Field: BARO, Comparator: GREATER THAN, Value: 9.53

Bertram Ludäscher, UC DAVIS

EML Metadata Display in Kepler

A simple example of using EML data. First, a search is done in the Data pane to locate an EML-described data set, which is dragged onto the workflow canvas. The EML data source is added to the workflow, and then it contacts the EcoGrid server to download the data and configure the ports. After being configured, it displays the ports from the EML data source, which are then mapped into an XY scatterplot.

EML Metadata Display:

Data Set Description:

Identifier: lml-her-gra-1005
 Catalog System: lml
 Alternative Identifier: lml-her-gra-1005a1.1.1
 Title: Mollusc population abundance monitoring: Fall 2000 mid-marsh and creekbank infarund and epifaunal mollusc abundance collections from CCF marsh monitoring, site 1-10

Organization: Georgia Coastal Ecosystems LTER Project
 Dept. of Marine Sciences, University of Georgia, Athens, Georgia 30602-3626 USA
 Email Address: gcoast@uga.edu
 Web Address: http://gce-her.marine.uga.edu/ter/

Individual: Dr. Mervyn Abber
 Organization: University of Georgia
 Email Address: mabber@uga.edu
 Individual: Mr. Kenneth Hinds
 Organization: University of Georgia Marine Institute
 Email Address: khinds@uga.edu

This data set is the Fall 2000 estimate of infaunal and epifaunal mollusc abundance at the CCF-LTER marsh sites used for pond monitoring. Species abundance was determined by hand-collecting all the infaunal and epifaunal molluscs from five within quadrats in mid-marsh and creekbank zones ($n = 5$ quadrats per area) at all sites. The abundance was estimated to the 1/8th, rounded to the nearest 1/8th, counted and measured (size data is reported separately). The counts were converted to number per 250g.

Bertram Ludäscher, UC DAVIS

Our Starting Point: 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 well-defined models of computation in the interaction components.

Principal Investigator:

Edward A. Lee

Technical Staff:

Christopher Anagnostopoulos

Mary P. Stewart

Postdocs and Researchers:

Jorn Janneck

James Soder

Grad Students:

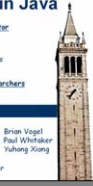
Elaine Chung

Chamberlain Fong

Lin Liu

Xiaojun Liu

Steve Rauscher



DATAFLOW PROCESS NETWORKS

Edward A. Lee
Thomas M. Parks

read!

Published in *Proceedings of the IEEE*, May, 1995.
© 1995, IEEE - All Rights Reserved

ABSTRACT

We review a model of computation used in industrial practice in signal processing software environments and experimentally in other contexts. We give this model the name "dataflow process networks," and study its formal properties as well as its utility as a basis for programming language design. Variants of this model are used in commercial visual programming systems such as SPW from the Alta Group of Cadence (formerly Comdisco Systems), COSSAP from Synopsys (formerly Cadis), the DSP Station from Mentor Graphics, and Hypersignal from Hyperception. They are also used in research software such as Khovos from the University of New Mexico and Ptolemy from the University of California at Berkeley, among many others.

Dataflow process networks are shown to be a special case of Kahn process networks, a model of computation where a number of concurrent processes communicate through unidirectional FIFO channels, where writes to the channel are non-blocking, and reads are blocking. In dataflow process networks, each process consists of repeated "firings" of a dataflow "actor." An actor defines a (often functional) quantum of computation. By dividing processes into actor firings, the considerable overhead of context switching incurred in most implementations of Kahn process networks is avoided.

We relate dataflow process networks to other dataflow models, including those used in dataflow machines, such as static dataflow and the tagged-token model. We also relate dataflow process networks to functional languages such as Haskell, and show that modern language concepts such as higher-order functions and polymorphism can be used effectively in dataflow process net-

see!

try!

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

Bertram Ludäscher, UC DAVIS

Why Ptolemy II ?



- **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
 - many research results
 - Ptolemy II participation in Kepler

Bertram Ludäscher, UC DAVIS

Kepler Today: Some Numbers



- **#Actors:**
 - Kepler: ~160 new + ~120 inherited (PTII)
 - soon there can be thousands (harvested from web services, R packages, etc.)
- **#Developers:**
 - ~ 24+, ~10 very active; more coming... (we think :-)
- **#CVS Repositories: ~2**
 - hopefully not increasing... :-{
- **# “Production-level” WFs:**
 - currently ~8, expected to increase quite a bit ...

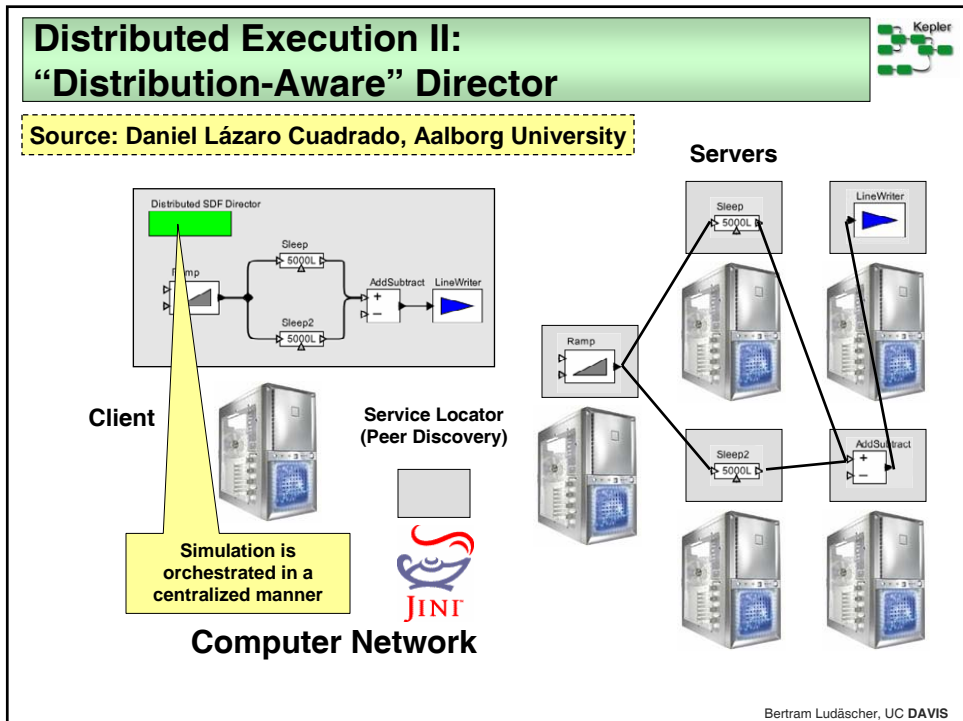
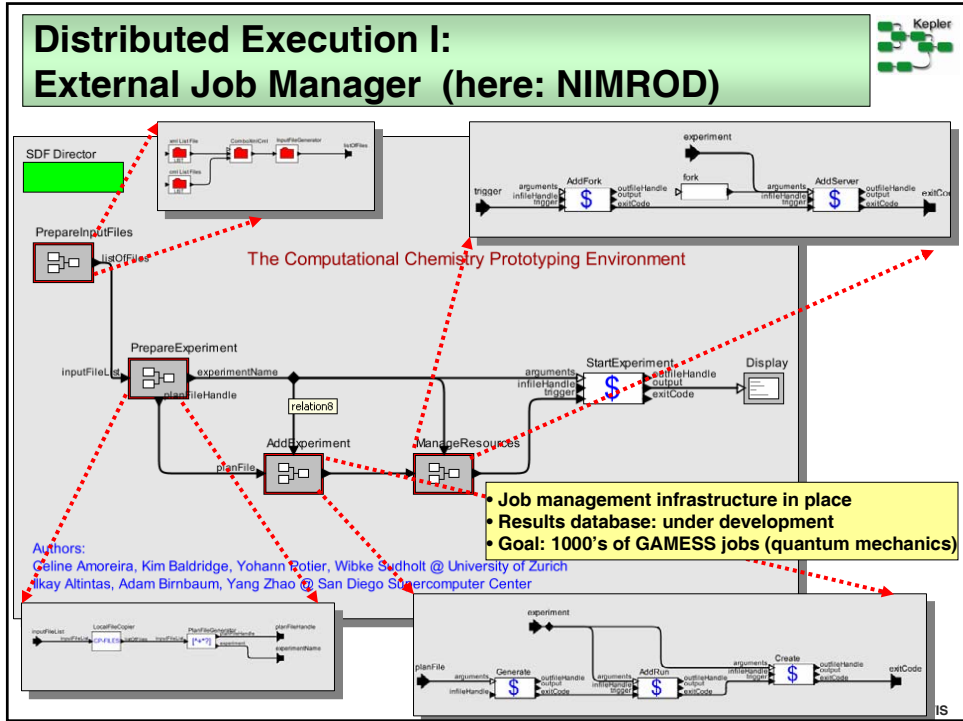
Bertram Ludäscher, UC DAVIS

A User’s Wish List



- Usability
- Closing the “lid” (cf. vnc)
- Dynamic plug-in of actors (cf. actor & data registries/repositories)
- Distributed WF execution
- Collection-based programming
- Grid awareness
- Semantics awareness
- WF Deployment (as a web site, as a web service, ...)
- “Power apps” (→ SCIRun)
- ...

Bertram Ludäscher, UC DAVIS



Separation of Concerns



- **A shining example:**
 - Ptolemy Directors – “factoring out” the concern of workflow “orchestration” (MoC)
 - common aspects of overall execution **not** left to the actors
- **Similarly:**
 - The “Black Box” (“flight recorder”)
 - a kind of “recording central” to avoid wiring 100’s of components to recording-actor(s)
 - The “Red Box” (error handling, fault tolerance)
 -
 - The “Yellow Box” (type checking)
 -
 - The “Blue Box” (shipping-and-handling)
 - central handling of data transport (by value, by reference, by scp, SRB, GridFTP, ...)

SDF/PN/DE/...

Recorder

On Error

Static Analysis

SHA @

Bertram Ludäscher, UC DAVIS

Separation of Concerns: Port Types



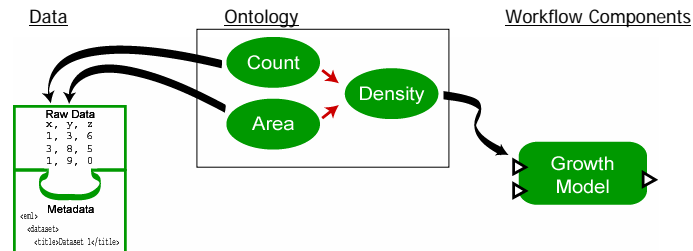
- **Token consumption (& production) “type”**
 - a director’s concern
- **Token “transport type”**
 - by value, reference (which one), protocol (SOAP, scp, GridFTP, scp, SRB, ...)
 - a SHA concern
- **Structural and semantic types**
 - SAT (static analysis & typing) concern
 - built after static unit type system...
 - static unit type system as a special case!?

Bertram Ludäscher, UC DAVIS

Need for Semantic Annotations of data & actors



- Label **data** with **semantic types** (concept expressions from an ontology)
- Label **inputs and outputs of analytical components** with **semantic types**

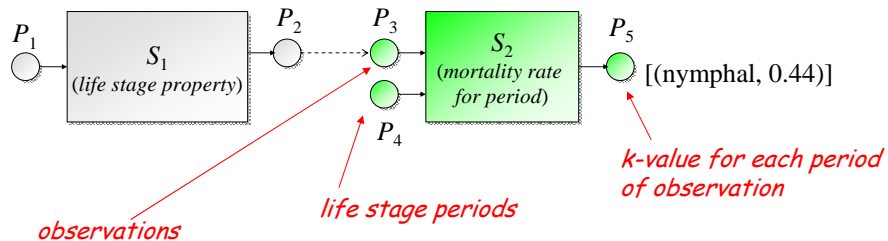


Example: Data has COUNT and AREA; workflow wants DENSITY

- → via ontology, system “knows” that data can still be used (because DENSITY := COUNT/AREA)
- Use **reasoning engines** to generate transformation steps
- Use **reasoning engine** to discover relevant components

Bertram Ludäscher, UC DAVIS

A Scientist’s “Semantic” View of Actors



Phase	Observed
Eggs	44,000
Instar I	3,513
Instar II	2,529
Instar III	1,922
Instar IV	1,461
Adults	1,300

Period	Phases
Nymphal	{Instar I, Instar II, Instar III, Instar IV}

Periods of development in terms of phases

Population samples for life stages of the common field grasshopper [Begon et al, 1996]

Source: [Bowers-Ludaescher, DILS'04]

Structural Type (XML DTD) Annotations



structType(P_2)

```

root population = (sample)*
elem sample     = (meas, lsp)
elem meas       = (cnt, acc)
elem cnt        = xsd:integer
elem acc        = xsd:double
elem lsp        = xsd:string
    
```

```

<population>
  <sample>
    <meas>
      <cnt>44,000</cnt>
      <acc>0.95</acc>
    </meas>
    <lsp>Eggs</lsp>
  </sample>
  ...
</population>
    
```

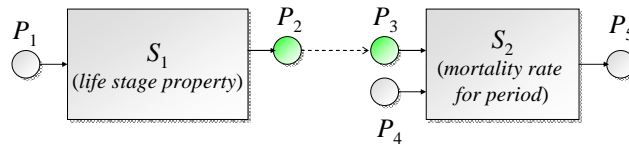
structType(P_3)

```

root cohortTable = (measurement)*
elem measurement = (phase, obs)
elem phase       = xsd:string
elem obs         = xsd:integer
    
```

```

<cohortTable>
  <measurement>
    <phase>Eggs</cnt>
    <obs>44,000</acc>
  </measurement>
  ...
</cohortTable>
    
```

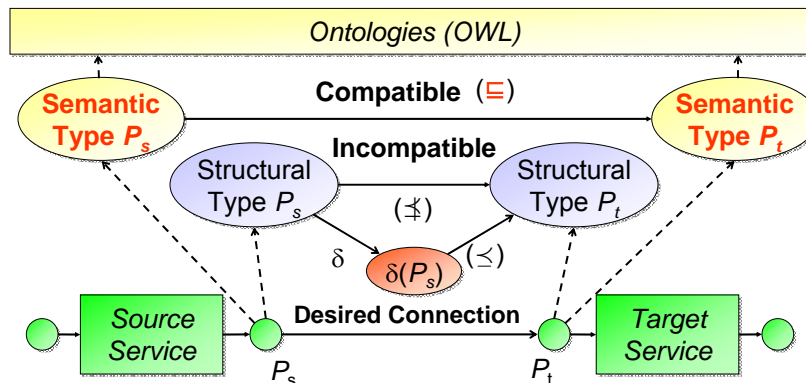


Source: [Bowers-Ludaescher, DILS'04]

A KR+DI+Scientific Workflow Problem

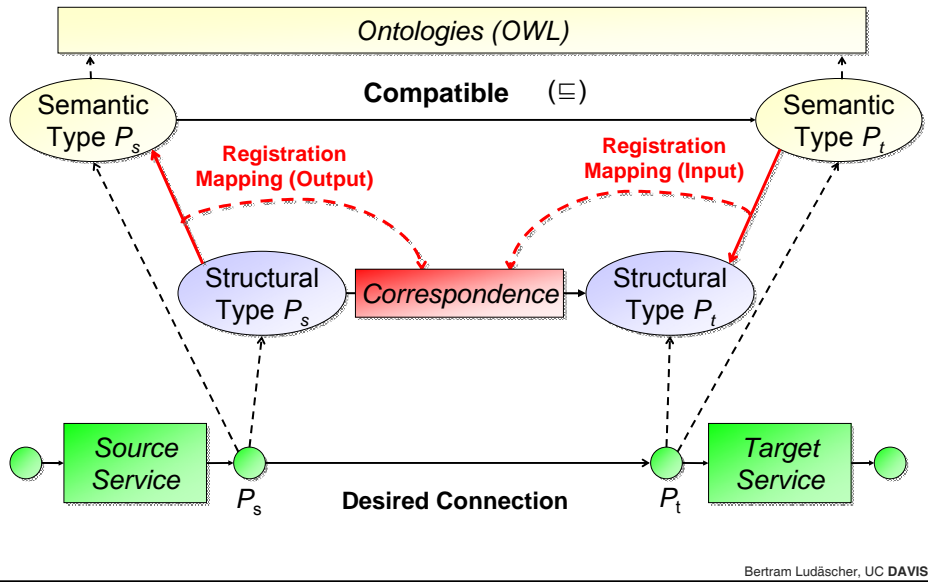


- Services can be **semantically compatible**, but **structurally incompatible**

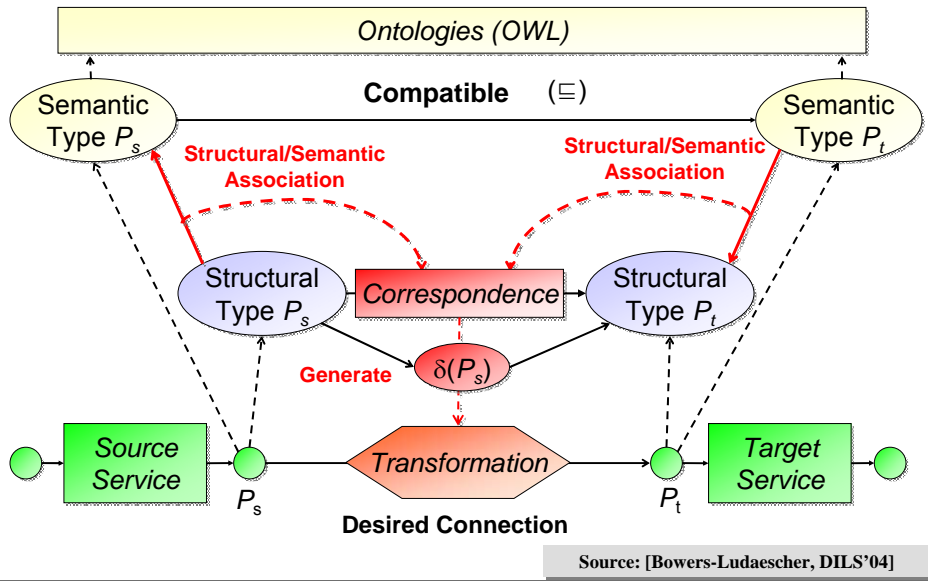


Source: [Bowers-Ludaescher, DILS'04]

The Ontology-Driven Framework



Ontology-Guided Data Transformation



Use of Semantics in SWF (DI+KR+SWF)



“Smart” Search

- Concept-based, e.g., “find all datasets containing biomass measurements”

Improved Linking, Merging, Integration

- Establishing links between data *through* semantic annotations & ontologies
- Combining heterogeneous sources based on annotations
- Concatenate, Union (merge), Join, etc.

Transforming

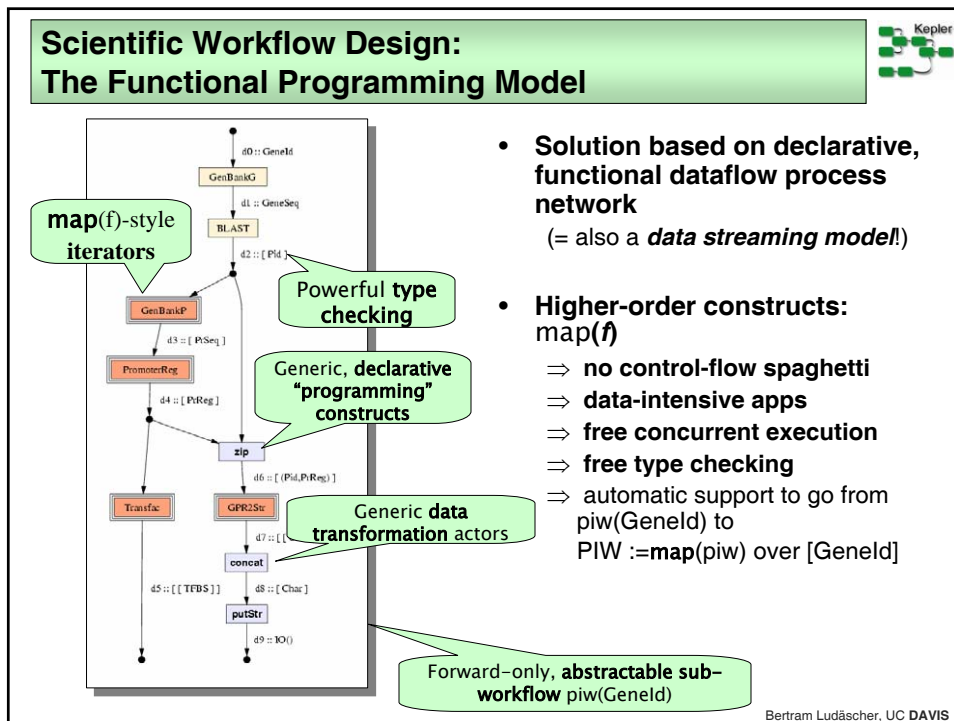
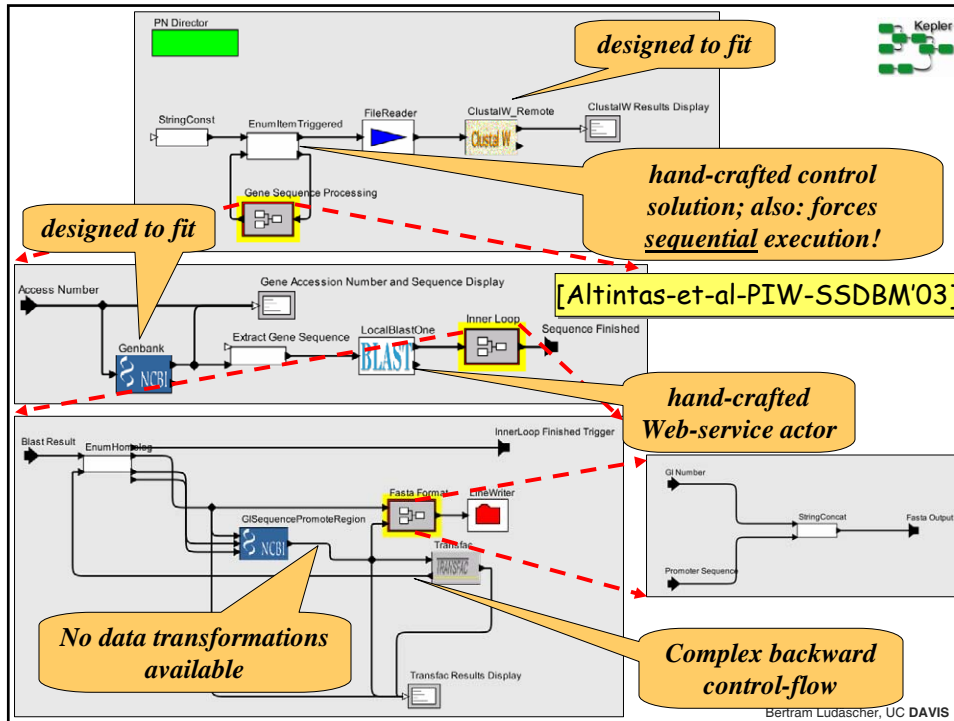
- Construct mappings from schema S1 to S2 based on annotations

Semantic Propagation

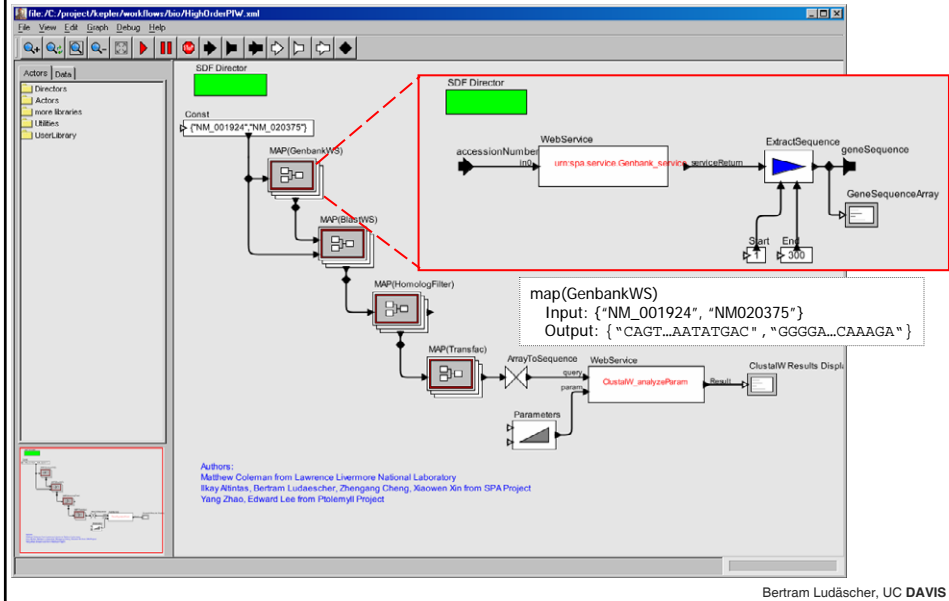
- “Pushing” semantic annotations through transformations/queries

Bertram Ludäscher, UC DAVIS

Scientific Workflow Design

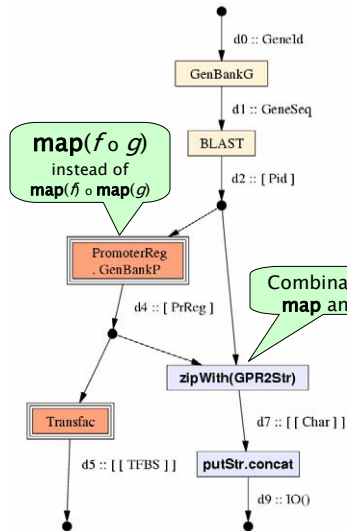


Scientific Workflow Design: The Functional Programming Model



Bertram Ludäscher, UC DAVIS

Research Problem: Optimization by Rewriting



- **Example: PIW as a declarative, referentially transparent functional process**
 \Rightarrow optimization via functional rewriting possible
 e.g. $\text{map}(f \circ g) = \text{map}(f) \circ \text{map}(g)$
- **Technical report & PIW specification in Haskell**



<http://kbis.sdsc.edu/SciDAC-SDM/scidac-tn-map-constructs.pdf>

Bertram Ludäscher, UC DAVIS

Scientific Workflow Design: Challenges



While many systems (including Kepler) support execution ... support for **SWF conceptual modeling and design** is lacking

- Formal models for scientific workflows
- Mechanisms for discovery, reuse, and adaptation of existing workflows and components
- End-to-end workflow development (methods and frameworks), especially for early stages

Bertram Ludäscher, UC DAVIS

Scientific Workflow Design: Contributions



While many systems (including Kepler) support execution ... support for SWF conceptual modeling and design is lacking

- Formal models for scientific workflows
- Mechanisms for discovery, reuse, and adaptation of existing workflows and components
- End-to-end workflow development (methods and frameworks), especially for early stages

Based on Actor-Oriented Modeling

plus a rich Type System ("hybrid" types)

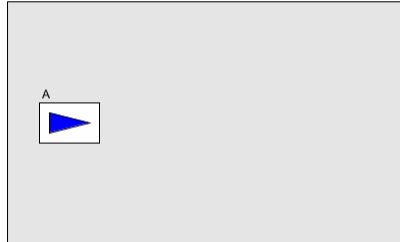
Modeling Primitives

(Adapters & Replacement; strategies)

Bowers-Ludaescher-ER-2005

UC DAVIS

Actor-Oriented Modeling



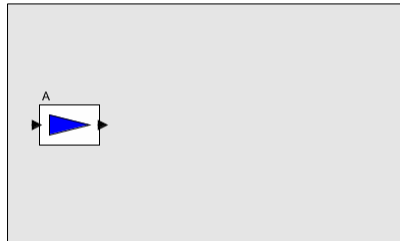
Actors

- single component or task
- well-defined interface (signature)
- dataflow view: given input data, produce output data

Bowers-Ludaescher-ER-2005

UC DAVIS

Actor-Oriented Modeling



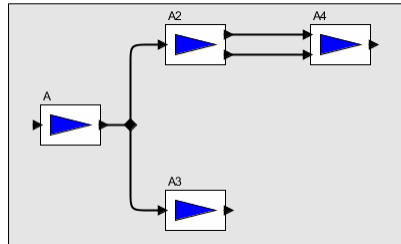
Ports

- each actor has a set of input and output ports
- denote the actor's signature
- produce/consume data (a.k.a. **tokens**)
- **parameters** are special “static” ports

Bowers-Ludaescher-ER-2005

UC DAVIS

Actor-Oriented Modeling



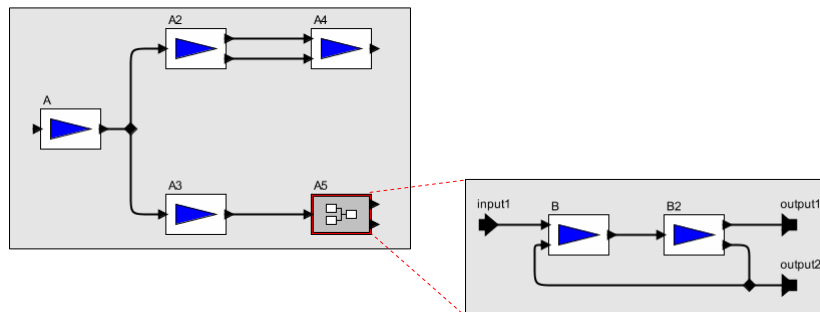
Dataflow Connections

- actor “communication” channels
- directed (hyper) edges
- connect output ports with input ports
- merge step + distribute step

Bowers-Ludaescher-ER-2005

UC DAVIS

Actor-Oriented Modeling



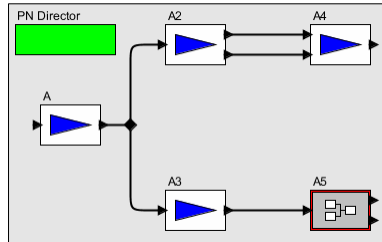
Sub-workflows / Composite Actors

- composite actors “wrap” sub-workflows
- like actors, have signatures (i/o ports of sub-workflow)
- hierarchical workflows (arbitrary nesting levels)

Bowers-Ludaescher-ER-2005

UC DAVIS

Actor-Oriented Modeling



Directors

- define the **Model of Computation (MoC)** of workflow graphs
- executes workflow
- sub-workflows may have different directors
- enables reusability

Bowers-Ludaescher-ER-2005

UC DAVIS

Models of Computation



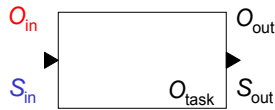
Directors **separate the concerns** of WF orchestration from Actor execution

- **Process Networks (PN)**
 - Actors execute as independent processes; blocking reads; non-blocking writes; FIFO buffers (queues) of unbounded size
- **Synchronous Dataflow (SDF)**
 - Fixed buffer size, since schedules are statically precomputed (based on token production/consumption rates/**firing**). SDF MoC is highly analyzable and used often in SWFs.
- Continuous Time (CT)
 - Connections represent the value of a continuous time signal at some point in time ... Often used to model physical processes.
- Discrete Event (DE)
 - Actors communicate through a queue of events in time. Used for instantaneous reactions in physical systems.
- ...

Bowers-Ludaescher-ER-2005

UC DAVIS

“Hybrid” Types



$O_{in} : \text{Measurement} \sqcap$
 $\quad \forall \text{ItemMeasured.SpeciesOccurrence}$
 $S_{in} : r(\text{site, day, spp, occ})$

Structural Types: Given a type language* \mathcal{S}

- Any port can be associated with a type $S \in \mathcal{S}$
- Kepler types include atomic (int, double, string) and complex types (record, list, etc.)

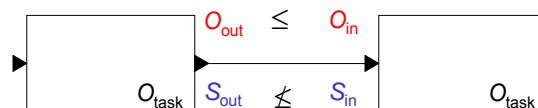
Semantic Types: Given an ontology language \mathcal{O}

- Any port can be associated a type $O \in \mathcal{O}$
- We consider description logic ontologies (e.g., OWL-DL)

* e.g., XML Schema, DTDs, relational, Kepler data types

Bowers-Ludaescher-ER-2005 UC DAVIS

Advantages of Hybrid Types



Semantically compatible
but structurally incompatible

Separates concerns of structural data typing and conceptual data typing

- Motivated by interoperating **legacy** and **independently** created components
- Allows one to be specified first (e.g., semantic)
- Structural **validation** independent from semantic validation
- Enables **discovery**
- ...

Bowers-Ludaescher-ER-2005 UC DAVIS

Semantic Annotations



Semantic and structural types can be “glued” together
... via logical constraints

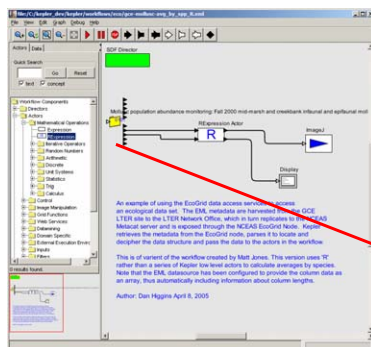
$\forall \text{site, day, spp, occ} \quad R(\text{site, day, spp, occ}) \rightarrow$
 $\exists y \text{ Measurement}(y), \text{ItemMeasured}(y, \text{occ}), \text{SpeciesOccurrence}(\text{occ})$

- Can provide **structural correspondences** between semantically compatible, structurally incompatible ports
- When something is known about the input/output dependencies of an actor, semantic types can be **propagated**

Bowers-Ludaescher-ER-2005

UC DAVIS

Semantic Annotations and Hybrid Types



Semantic Type Editor is used to assign one or more semantic types to the component or to the component's input and output ports.

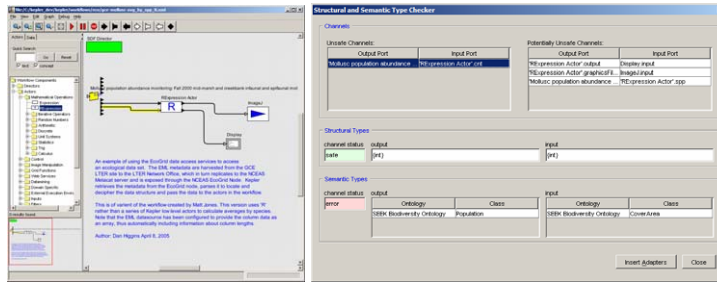


An **ontology browser** is provided in Kepler to navigate a classified OWL-DL ontology. Classes can be searched for and selected as a semantic type.

Bowers-Ludaescher-ER-2005

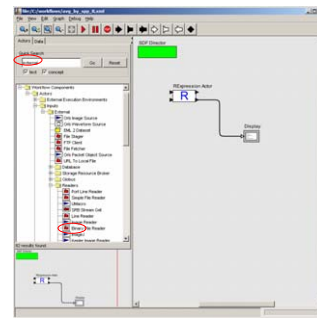
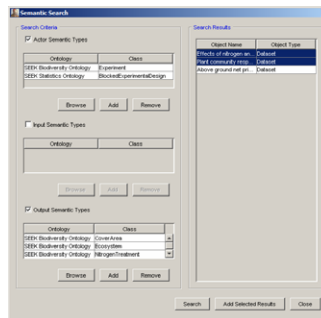
UC DAVIS

Semantic Annotations and Hybrid Types



Verifying structural and semantic compatibility of workflow connections in Kepler

Searching based on actor-level and input/output port **Semantic Types** in Kepler

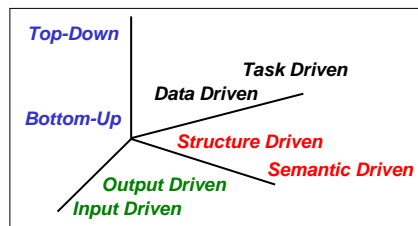
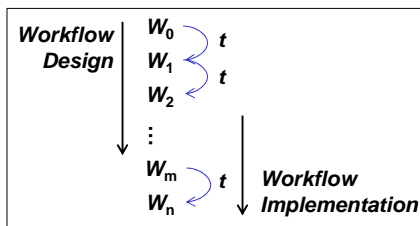


Workflow Design Primitives



End-to-End Workflow Design and Implementation

- Viewed as a series of primitive “transformations”
- Each takes a SWF and produces a new SWF
- Can be combined to form design “strategies”



E.g., re-engineering SWFs often is top-down, structure driven, whereas new SWFs are often a mix of semantic, input, and output driven

Basic Actor-Oriented Primitives



Basic Transformations	Starting Workflow	Resulting Workflow	Resulting Workflow
t_1 : Entity Introduction (actor or data connection)			
t_2 : Port Introduction			
t_3 : Datatype Refinement ($s' \preceq s, t' \preceq t$)			
t_4 : Hierarchical Abstraction			
t_5 : Hierarchical Refinement			
t_6 : Data Connection			
t_7 : Director Introduction			

Bowers-Ludaescher-ER-2005

UC DAVIS

Additional Primitives



Transformations	Starting Workflow	Resulting Workflow	Resulting Workflow
t_9 : Actor Semantic Type Refinement ($T' \leq T$)			
t_{10} : Port Semantic Type Refinement ($C' \leq C, D' \leq D$)			
t_{11} : Annotation Constraint Refinement ($\alpha \rightarrow \alpha'$)			
t_{12} : I/O Constraint Strengthening ($\psi \rightarrow \phi$)			
t_{13} : Data Connection Refinement			
t_{14} : Adapter Insertion			
t_{15} : Actor Replacement			
t_{16} : Workflow Combination (Map)			



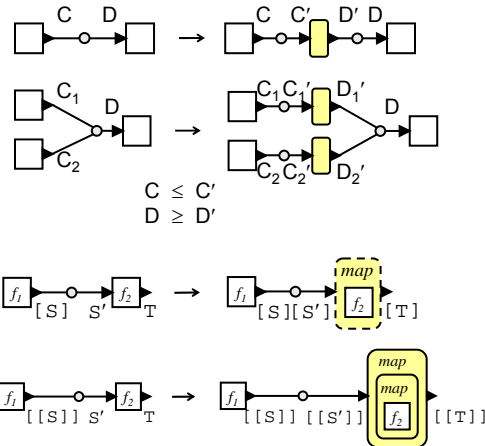
Bowers-Ludaescher-ER-2005

UC DAVIS

Adapters (Semantic and Structural Incompatibility)



Adapters:

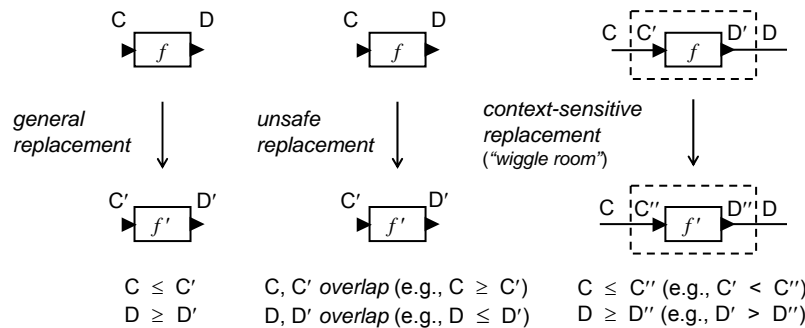


- Can be abstract (no implementation) or concrete
- Can bridge “semantic gaps” or fix structural mismatches
- Can be generated automatically (e.g., Taverna’s “list mismatch”)
- Can be reused (based on signatures)

Bowers-Ludaescher-ER-2005

UC DAVIS

The Replacement Primitive



- General replacement doesn’t consider **surrounding connections**
- Context-sensitive replacement gives more “wobble room” by “tuning” the actors semantic types based on connections

Bowers-Ludaescher-ER-2005

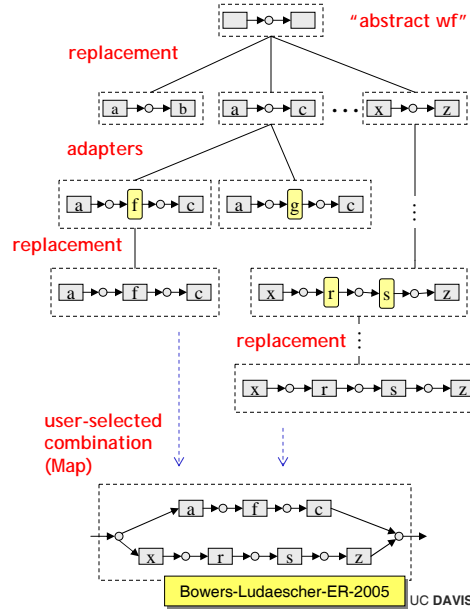
UC DAVIS

Applying adapter insertion and replacement



- **Adapter insertion and replacement can enable simple SWF elaboration**

- Given an initial set of connected **abstract** actors
- Repeatedly search for replacement concrete actors (atomic/composite)
- At each step, insert adapters when necessary
- Allow user to combine desired “variations”



The KEPLER Project

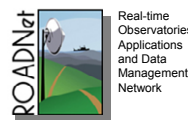
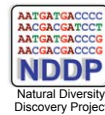


The Kepler Scientific Workflow System

- Built on Ptolemy II (UC Berkeley), developed by the Electrical Engineering community (circuit design and simulation)
- Open-source, Java
- Computation Models, Nested WFs, Loops
- Graphical Workflow Interface
- Workflow Execution
- Extensible Architecture
- Component Libraries
- Metadata, Discovery, Archival

The Kepler “Vision”

- End-to-end scientific workflow design and execution environment
- Data and compute intensive workflows
- Comprehensive component libraries for a wide range of scientific domains
- Enable collaboration, sharing across disciplines (“synergy”)



KEPLER kepler-project.org

NDDP www.nddp.org
 EOL eol.sdsc.edu
 GEON www.geongrid.org
 Ptolemy II ptolemy.eecs.berkeley.edu/ptolemyII
 ROADNet roadnet.ucsd.edu
 SEEK seek.ecoinformatics.org
 SciDAC www-casc.llnl.gov/sdm

Bertram Ludäscher, UC DAVIS

