

Internet Technology Group

# The Semantic Layered Research Platform Prototype and KIM

Sean Martin W3C Life Sciences Workshop October 29, 2004



#### Researchers often have a highly manual research process and no integrated information management systems

#### **Current Environment**

Vertical research methods and application silos
 Researchers think in terms of discrete tools needed to complete each individual research step

Piece-meal single-task oriented apparatus

Researchers think of and treat IT like any other lab equipment

- Highly skilled, highly manual processes
- Collaboration and community communication occur through word of mouth, conferences, published papers, and email
- Information often mislaid, lost, or dots not connected, context not trapped
- Or inordinate amount of time is spent manually keep track of information
- Increasing recognition that In Silico, or digital biology, is the way forward

_	-		_
-			1

#### Researchers often have a highly manual research process and no integrated information management systems

Entire endeavor is about creating information, but while it is created, it is not managed

- Time wasted manually moving data from one place to another
- Time spent checking and integrating data from different silos
- Information is lost
- In silico is treated like the physical world without seeing the power of data connections
- Collaboration is far more difficult
- Scientists are isolated from key data

This lack of an information integration backbone results in an inefficient research process



Semantic Layer Research Platform is designed to improve research process by providing an integrated information management system

#### Reduce barriers to HPC/GRID computing

- "Cost" (time/manual setup) of running/rerunning experiments is slashed Compute requirements of researchers increase dramatically as manual setup declines
- SLRP tracks, integrates and aids in all aspects of a research project, including data generated from experiment and literature searches with computational experiments, collaborations, and publication
- Highly repeatable experiments, better provenance
  - Clear Information trail
  - Exact algorithms, Exact data
  - Able to reference/send entire context to colleagues
- Better collaboration
  - Collaborative annotation
  - Seamless User Interfaces
  - Easier peer review of experiments (and their own "riffs")

_	
_	 

# The Semantic Layer Research Platform requires new technologies and new uses of existing technologies

- Cart Semantic data infrastructure
- Slingshot BPAL-like workflow capabilities for the Semantic Web
- DDR Write once, read-only object repository
- Jastor Code generation for RDF handling
- Telar Using RDF in Eclipse
- LSID System underpinning: All objects and their relationships uniquely named
- Standards-based RDF, LSID, CIM, WSDL, SOAP, etc.
- Open source

Adding up to the Semantic Layer Research Platform



### CART: Semantic Web Infrastructure for the Enterprise

Collections, Access Control, Revision Tracking

CART is a Client / Server RDF storage architecture built on the WebSphere Platform (DB2, EventBroker)

- Current situation: Simple RDF models published on web Client applications store metadata in RDF
- For the Enterprise we need: Multi-user
  - Distributed replicated/synchronized views
  - Distributed Events
  - Transaction support
  - Access Control

Reification Duplicate triple support Collections Scalability History / Auditable Data federation





## Slingshot

#### Enables dynamic semantic workflows

 Problem: Scientists run dynamic, "ad-hoc" workflows of tasks that execute in heterogeneous environments

Workflows are developed quickly and modified often

Workflows are carried out cooperatively in several types of participants including HPC/GRID applications, Web Service invocations and user-interactive client applications

#### Solution: OWL-S++ execution engine

- ++ == Groundings for GRID/HPC and workspace execution
- Current state of flow stored in CART (central RDF-store)
- Participants are notified when a new flow is submitted to the store and when the state of the flow changes (PUB-SUB)
- Each participant independently evaluates the flow to determine if it has to contribute anything at the moment
- When a participant commits a task, the state of the flow changes and all participants must reevaluate



#### **Distributed Data Repository**

Write once, read only file object store that provides secure access by LSID reference to data over AHPFS in grid clusters, WAN links, SOAP, and HTTP

- All aspects of DDR are described in RDF using OWL ontologies
- Repository based on Linux and ReiserFS with extended file system attributes
- Server stack based on Apache with mod\_perl providing access via SOAP and HTTP
- Remote administration of repository via XML/RDF in SOAP
- Granular access control via file system ACLs
- Java client API, including LSID protocol handler, provides online/offline writing, reading, and synchronizing of data and associated metadata



### Jastor

#### Code generation tool for creating custom libraries

- Developers don't actually want to deal with RDF
- Provides object-oriented view, including events, of RDF Conforms to a particular OWL ontology
- Extensible to any object oriented programming language Java classes
- OWL

Compliant with OWL-Lite

Supports an increasing number features of OWL-DL and OWL-Full



## Life Science Identifiers (LSID)

URN naming, discovery and access standard for data and metadata

- All system & data objects and concepts have a LSID URI
  - CART identifies documents/collections with LSIDs
  - LSIDs are used to name ontologys, formats & predicates
  - DDR provides access to files via LSID
- CVS <-> LSID Resolution service to track code modules
- Use LSIDs to track and access external references wherever possible

_	

#### Telar is a framework for building Semantic user interfaces

A framework designed to help developers build user interfaces capable of weaving data from a myriad of data sources into a single application for creating, exploring and organizing information

Java library for Eclipse

Enables mapping of dynamic RDF data to SWT widgets

Groups widgets into logical information clusters and manages their life cycle for UI forms presentation

Maps onto CART client-side API for synchronization support

Problems:

- Information about any subject cannot be contained within a single database or structure
- Users need a larger view of the information available on any given subject
- The number of applications users need to deal with all of today's data structures is overwhelming
- Existing applications are generally hard-coded to perform a limited set of functions on a few data structures

# Interdisciplinary Work

- Scientific Diversity: Multiple disciplines, methods and techniques
- Geographic Distance: Multiple sites (academia, government, corporate)
- Data Infrastructure:

Multiple distributed data bases (data sharing and integration, data safety and confidentiality)



Knowledge-Integrated Modeling (KIM)



# Knowledge-Integrated Modeling (KIM)

- Integrates peer-reviewed biomedical data into a 3D computational modeling, simulation and visualization platform (within a GRID environment).
- Enables large-scale biomedical collaborations.
- Guides experimental data acquisition and specifies data-base mining.
- Ultimately, is intended for both research *and* clinical use, i.e. for diagnostic purposes and treatment outcome predictions.















# **Semantic Layer**

- Rich network of objects and relationships
- Capture of information contexts in metadata
- Enabler of systems biology





## **Usability and Scenarios**

- Observed and interacted with customers Identified workbench personas Understood user activities and tasks
- Created scenarios based on observed working practices of scientists in their various personas Biologists, physicists, informaticians
- Created user interface designs



## Scenario: Creating simulation

👅 Life Sciences Workben	ch			
<u>File E</u> dit Se <u>a</u> rch <u>P</u> roject <u>R</u> un <u>Y</u>	<u>W</u> indow <u>H</u> elp		Tom 's Workbench	
🗢 🖘 🔹 Homer Re	search - Deve	lopment * Simulation * Analysis * Publish *		
Experiments 💌 🗙	New Simulation	6	▼ X	
New	Submit Stop		Runs: 1 Time: 1:00 Status:	
New Simulation 5 New Simulation 6 Tumor Simulation Image 2D Image 3D Analysis OpenDX Tiff2Movie	Title Template Version Owner Organization Creation date	New Simulation 6 Cell proliferation template 2 Tom Deisboeck Harvard-MIT Martinos Center for Biomedical Imaging February 4, 2004		
	Description	Scaling the simulation results according to Blankenberg et al., we see global exponential tumor growth dynamics until timepoint 'diagnosis. We further note local, asymmetric tumor growth		
	Clusters	Glioma		
	Belsted Info	Add Remove		
	Kelated IIIO			
	Alerts	<none></none>		
Collections Experiments Collections Experiments  Related Links  X New CaBlG CaBlG CaBlG CaBlG Cable Ca				
			<u>×</u>	
Alerts Run View	Overview Inputs 0	utputs   Workflow		























