

# The CombeChem Project

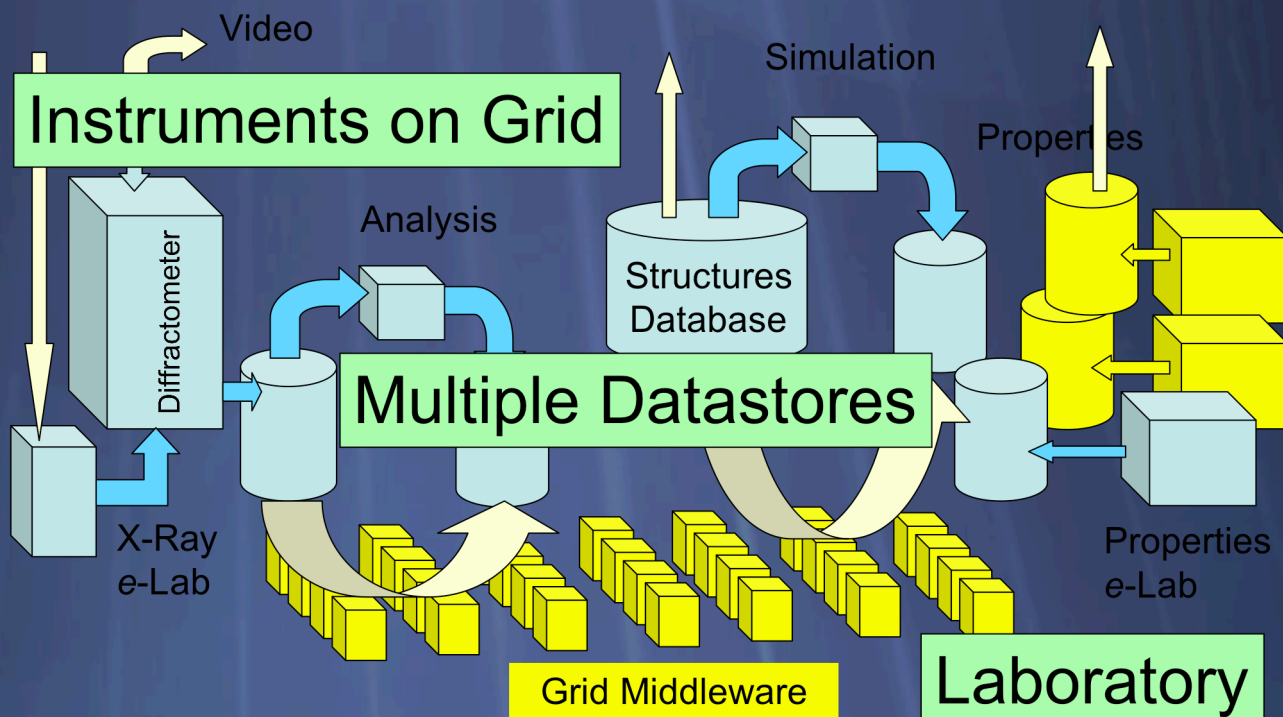
Semantic Support for the Chemistry Life cycle

Jeremy Frey Dave De Roure

Schools of Chemistry and  
Electronics and Computer Science  
University of Southampton

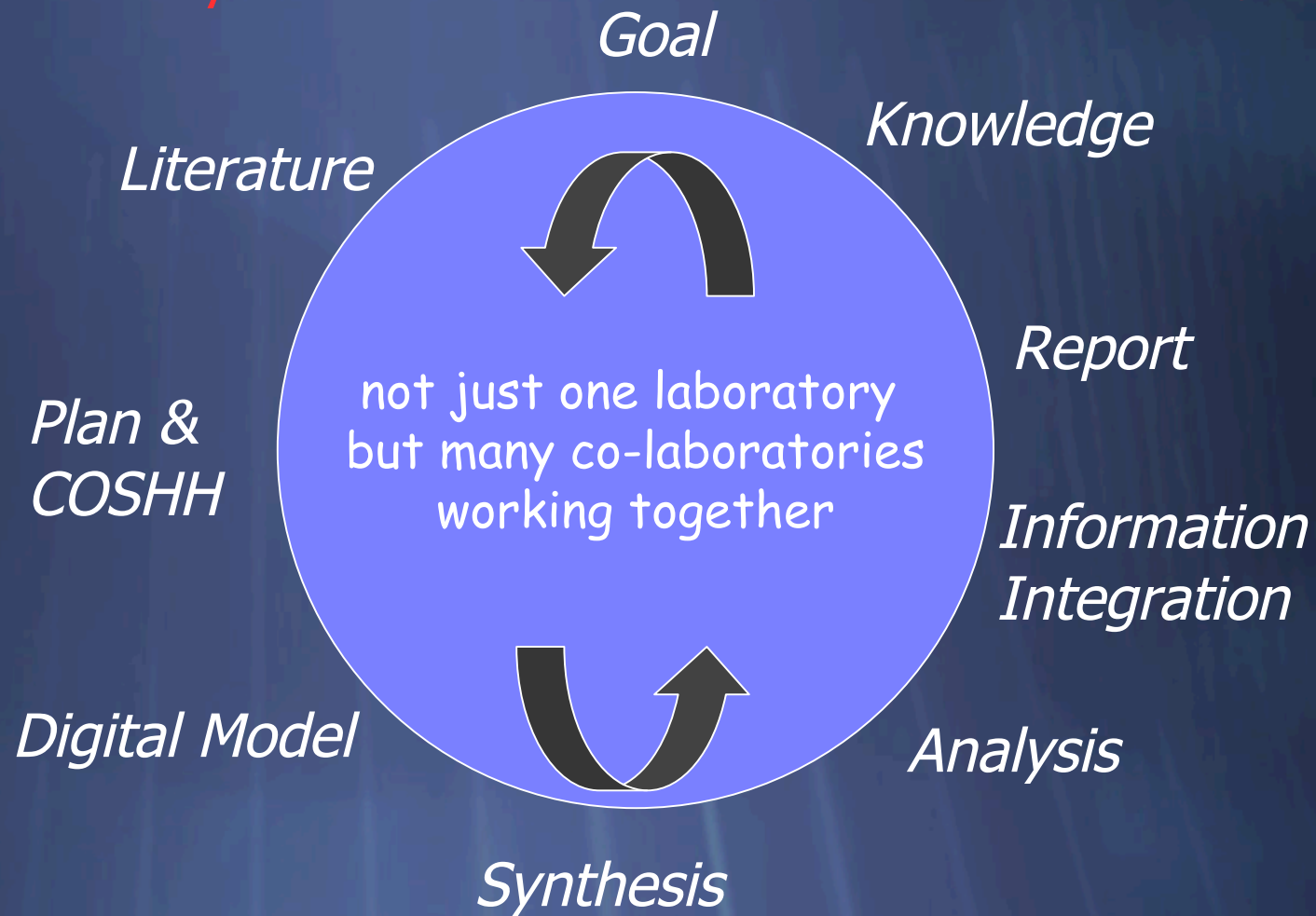
# The CombeChem Project

- ★ End to End linking of data and information: Laboratory to publication and back again
- ★ The exponential world of combinatorial synthesis and high throughput analysis meets the exponentially growing power of computing



Smart Laboratory

Smart HCI

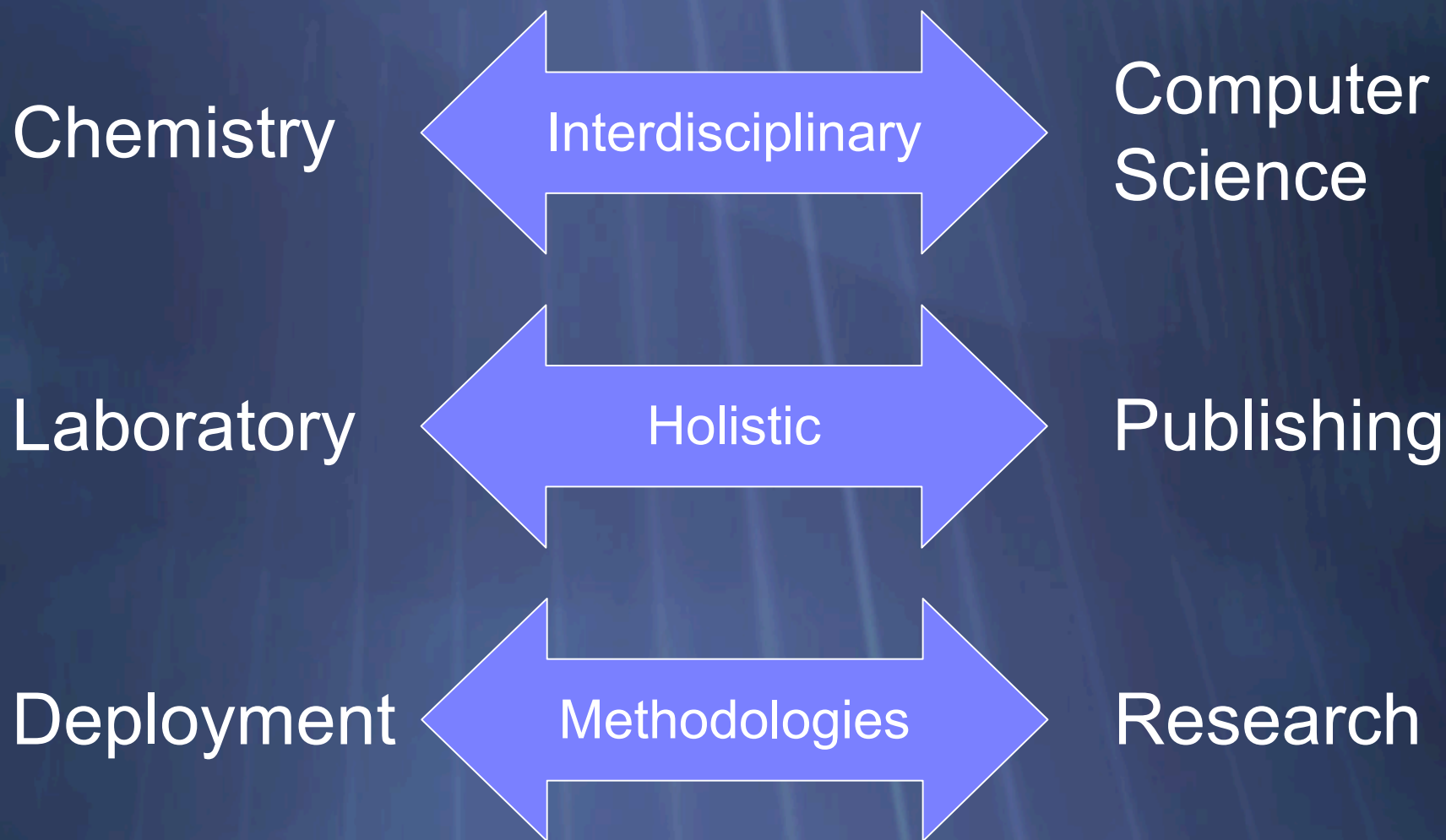


Smart Storage

Smart Dissemination

The concept of Publication @ Source

# The Stretches of CombeChem





# CombeChem Smart Tea

- ✦ Ethnography
- ✦ Electronic Lab Notebook
- ✦ Capture in RDF  
(Resource Description Framework)

the Smart Tea Project



"I can go anywhere and its, like, this is me and my data. Its all there, bang."

- Chris,  
a real chemist, on using Smart Tea  
instead of a paper lab book.

Smart Tea is about improving the information environment for chemists doing chemistry - within and beyond the lab. Smart Tea is about supporting chemists in the preparation, execution, analysis and dissemination of their experimental work.



necessary if a calculation or discussion is changed; the section to be deleted is simply removed by drawing a neat "x" through it.

In view of the fact that a notebook is a primary record, data are not copied into it from other sources (such as this manual or a lab partner's notebook, in a joint experiment) without clear acknowledgment of the source. Observations are never collected on note pads, filter paper, or other temporary paper for later transfer into a notebook. If you are caught using the "scrap of paper" technique, your improperly recorded data may be confiscated by your TA or instructor at any time. It is important to develop a standard approach to using a notebook routinely as the primary receptacle of observations.

Each week at the beginning of lab lecture, you will turn in your prelab problems from the manual for grading. Problems not turned in at the beginning of lab lecture will be

Observations are never collected on note pads, filter paper or other temporary paper for later transfer into a notebook

If you are caught using the "scrap of paper" technique, your improperly recorded data may be confiscated by your TA





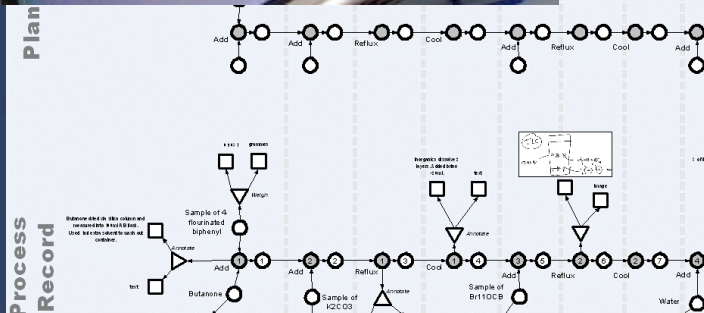
# PLANS

COSH Assessment Form			
SUBSTANCE NAME	PHYSICAL FORM	QUANTITY	NATURE OF HAZARD
Water	liquid	1000ml	None
Dextrose	Solid	<20g	possible contribution to sugar and stain
Caffeine	Solid (fine)	<1g	Respiratory & skin irritation, surface sensitizing.
Methyl	liquid	<100ml	No particular hazard

NATURE OF PROCESS  
liquid extraction of caffeine, followed by combination with dextrose to produce a sweet drink

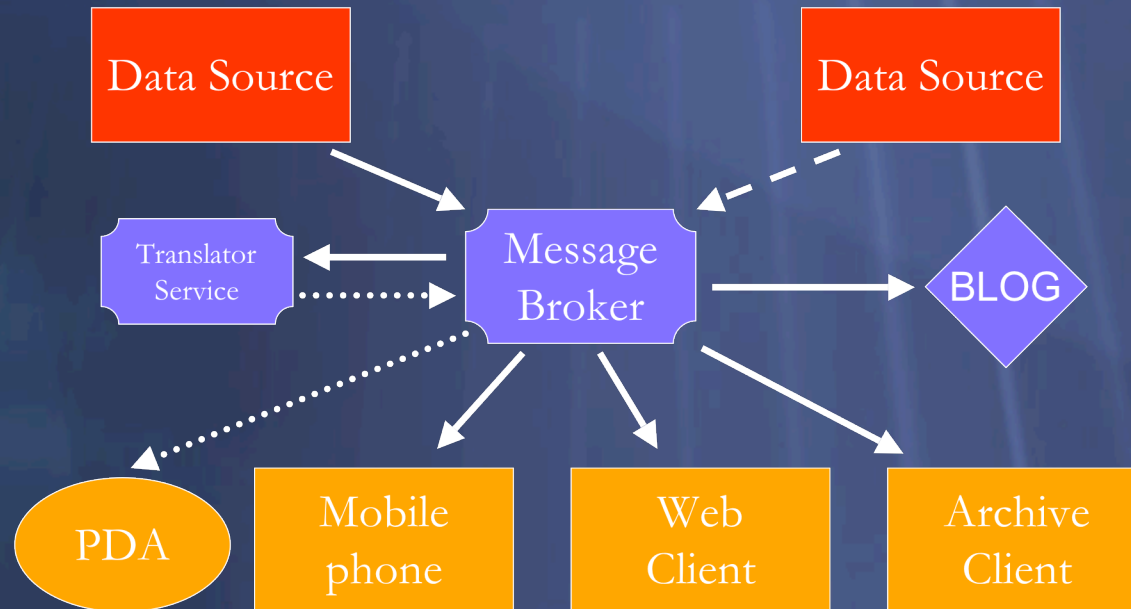
Is there a less hazardous substance? No  
If so, why not use it?

CONTROL MEASURES REQUIRED

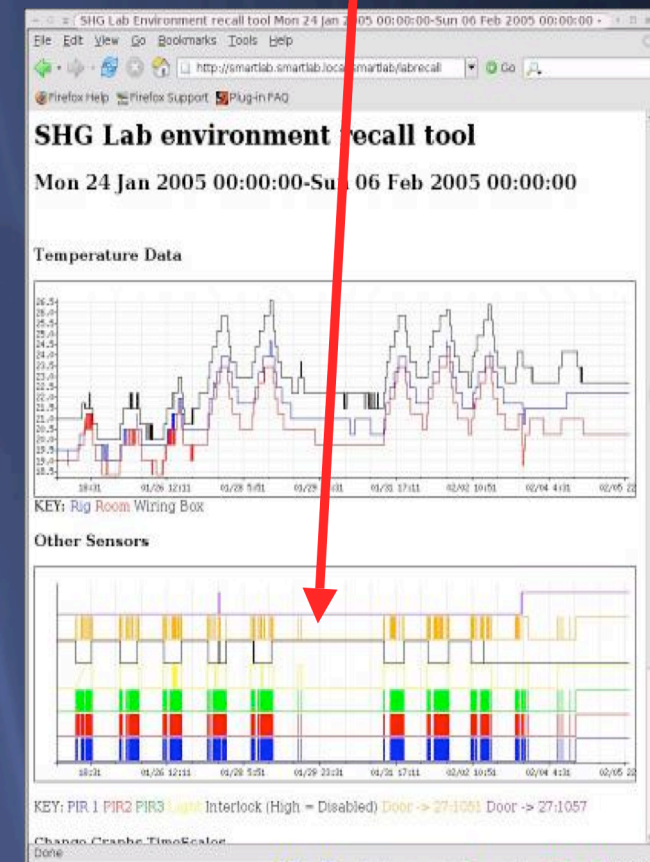


- ✦ Make use of Plans to inform the digital context - metadata in advance
- ✦ Have concern for the “End-to-End life cycle” of chemistry information from the start
- ✦ Understanding Usability and Human Computer Interaction is vital for adoption

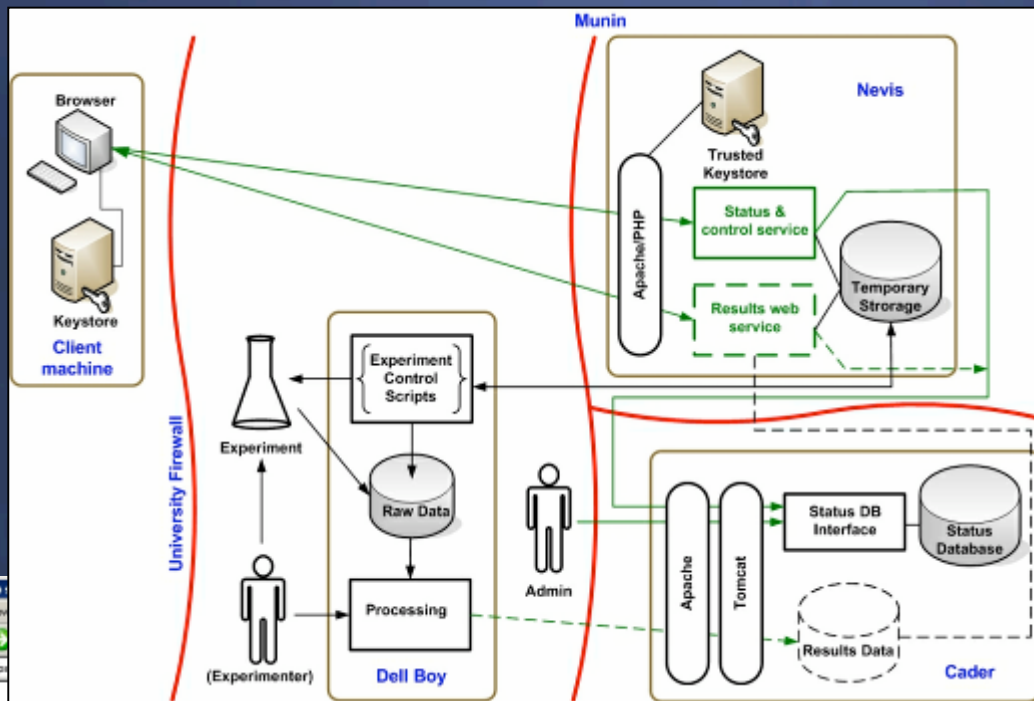
# Pub-Sub systems provide the flexible & extensible approach to distribution



Air Conditioning failed



# National Crystallography Service Grid Service Architecture



Users can follow and interact with experiment

The screenshot shows a web browser window with a workflow diagram at the top and a table of sample statuses below.

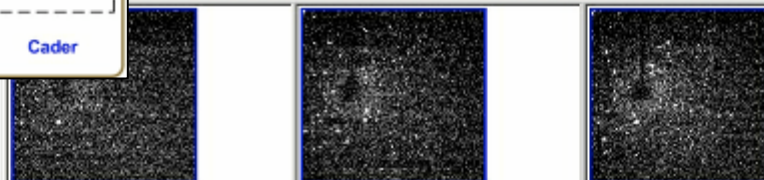
**Workflow Diagram:**

```

graph LR
    Experimenter["(Experimenter)"] --> Processing["Processing"]
    Processing --> Admin["Admin"]
    Admin --> Apache["Apache"]
    Admin --> Tomcat["Tomcat"]
    Apache --> Results["Results Data"]
    Tomcat --> Results
  
```

**Sample Status Table:**

NCS ID	Customer ID	Received	Collection	Status	Details
04MEL0098	2nd test	2004-02-12	001	Succeeded	<a href="#">HKL file</a> / <a href="#">Report</a>
04MEL0093	mel01	2004-02-06	001	Succeeded	<a href="#">HKL file</a> / <a href="#">Report</a>
04SRC0104	#13-123	2004-03-08	001	None	Due at 00:00:00 (est)
04SRC0103	#12-01	2004-03-08	001	Failed (Referred)	Diffraction too weak
			002	Failed (No Further Action)	Crystals too small
04SRC0105	HSP-HCI	2004-03-08	001	Added	





Chemical  
families  
polymorph  
similarities

Y X	NO <sub>2</sub>	CN	CF <sub>3</sub>	I	Br	Cl	F	H	Me	MeO
NO <sub>2</sub>										
CN										
CF <sub>3</sub>										
I										
Br										
Cl										
F										
H										
Me										
MeO										

A

B

B

C

D

E

F

G

H

I

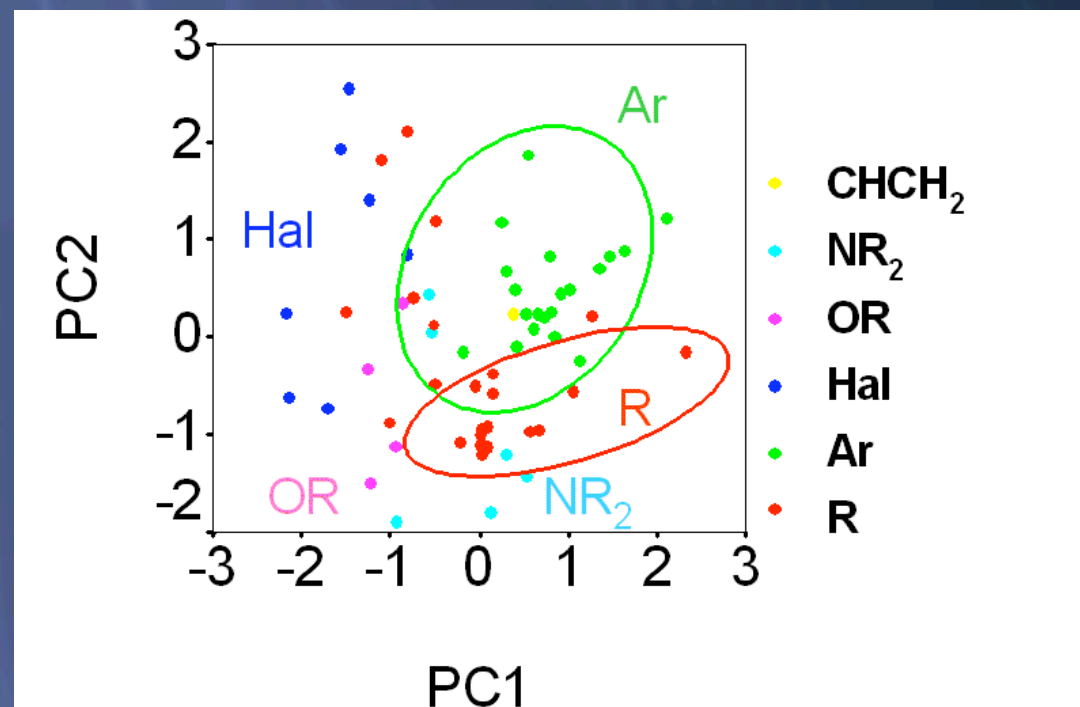
trans-  
lation  
(B1)

others  
(B2-B7)



# Statistics

- ★ Collect information about ligands and their (transition metal) complexes
- ★ Calculate descriptors with standard computational approach (DFT)
- ★ Robustness (computational, chemical, statistical)
- ★ Overlap with available experimental data



*Map of Ligand Space for  
Monodentate Phosphorus(III) Ligands*

Ligand Knowledge Base

# Dave's Chemistry Experiment

1. Take a building full of chemists
  2. Add RDF tools
  3. Stir occasionally
  4. See what's been made
- ★ A very big chunk of Semantic Web
  - ★ An ontology for units

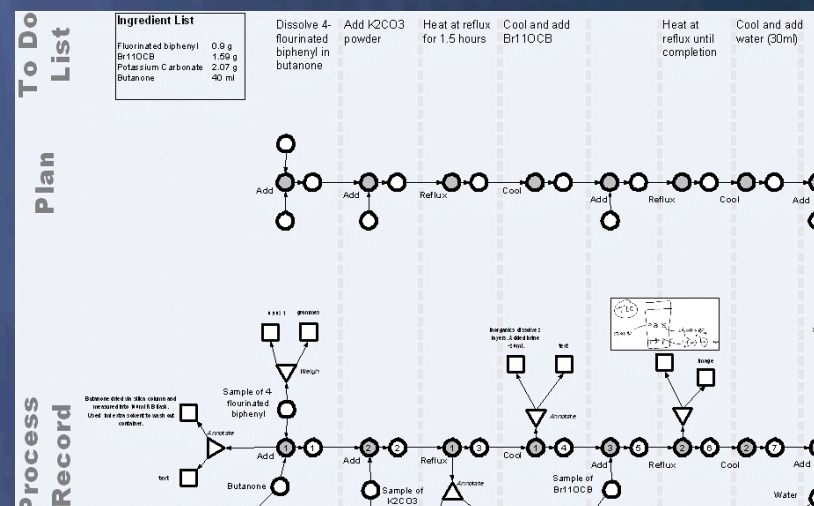
## Statistics on Green Triplestore

Thu Apr 21 11:37:17 2005

models	2573454
triples	84188993
inferred (FC)	24269915
ground facts	59919078
resources	9987377
literals	7974229
classes	88
properties	49

# Semantic DataGrid

- ★ CombeChem uses Semantic Web for
  - ★ Enhanced (annotated) DataGrid over multiple diverse stores
  - ★ Some Data Storage
  - ★ Storage of Provenance Information
  - ★ Annotated multimedia streams



# Triplestores

- ★ Started with the data hoarding approach of CSAKTive Space, using 3store from the AKT IRC
- ★ Scalability, lifecycle and the CombeChem sharing and publishing ethos led to the use of multiple triplestores to cache and query rather than store

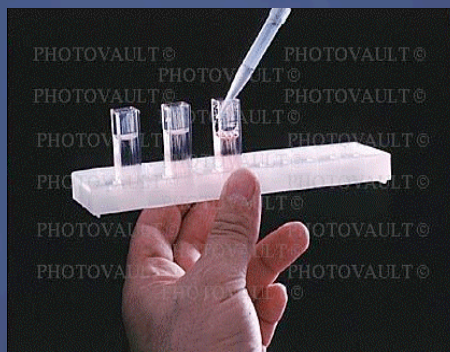




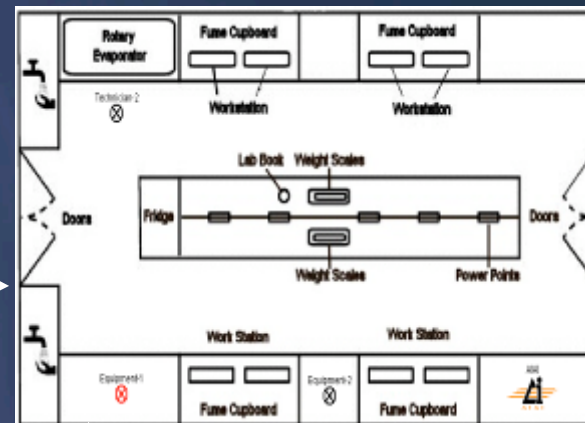
# Comb-e-Chem: Facility e-Science in Action



**Presence Awareness + Remote Participation**



**Resource + Floor Management**

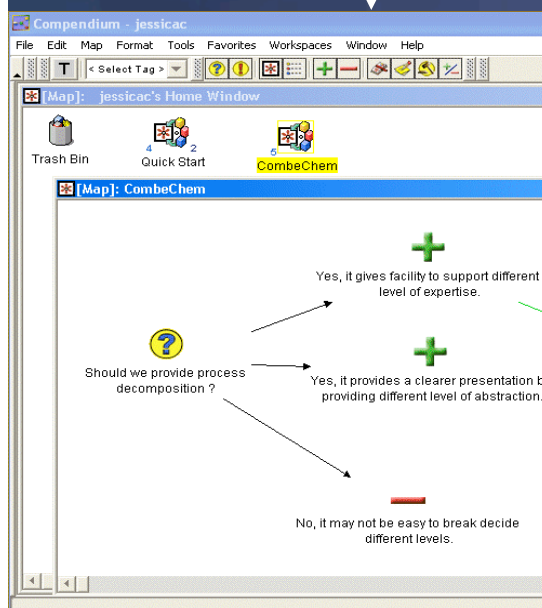


**Run-time tracking and control**

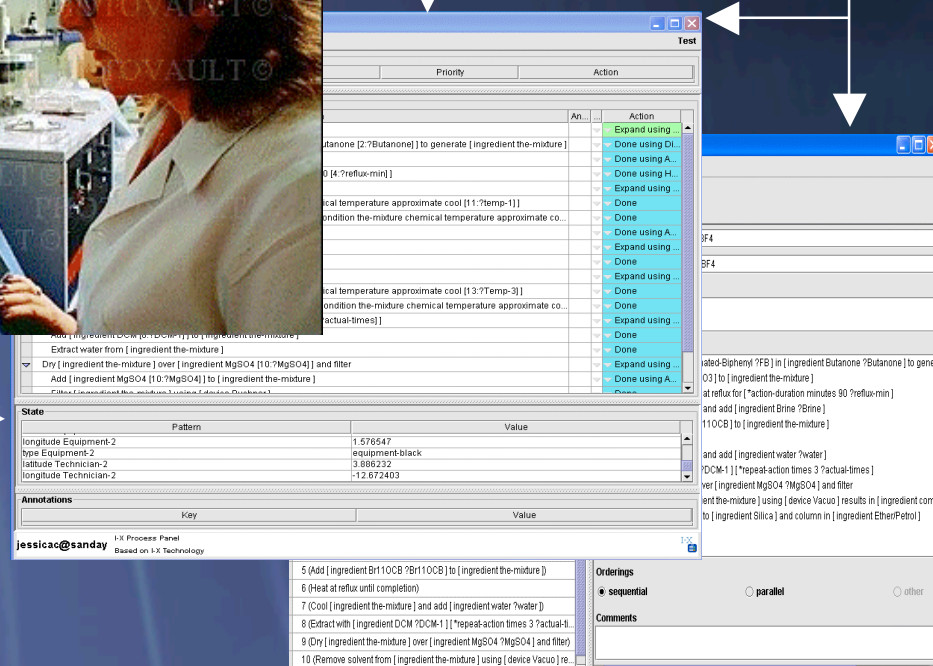


**Workflow Execution And Monitoring**

**Instrument Tracking And utilisation**

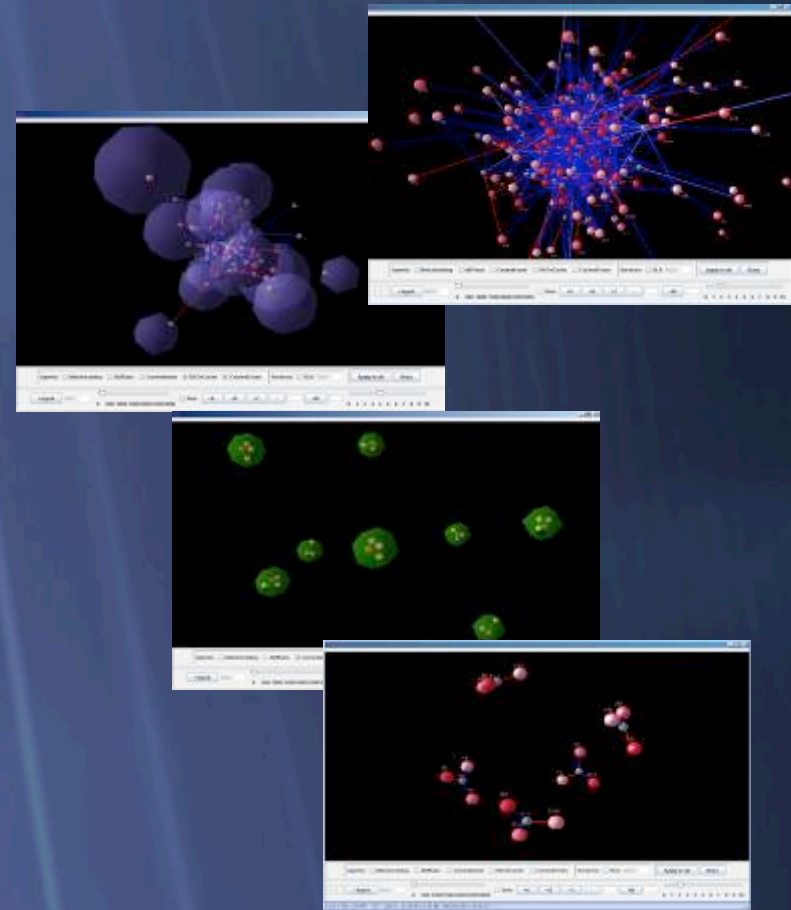


**Action Realisation + Rationale Feedback**



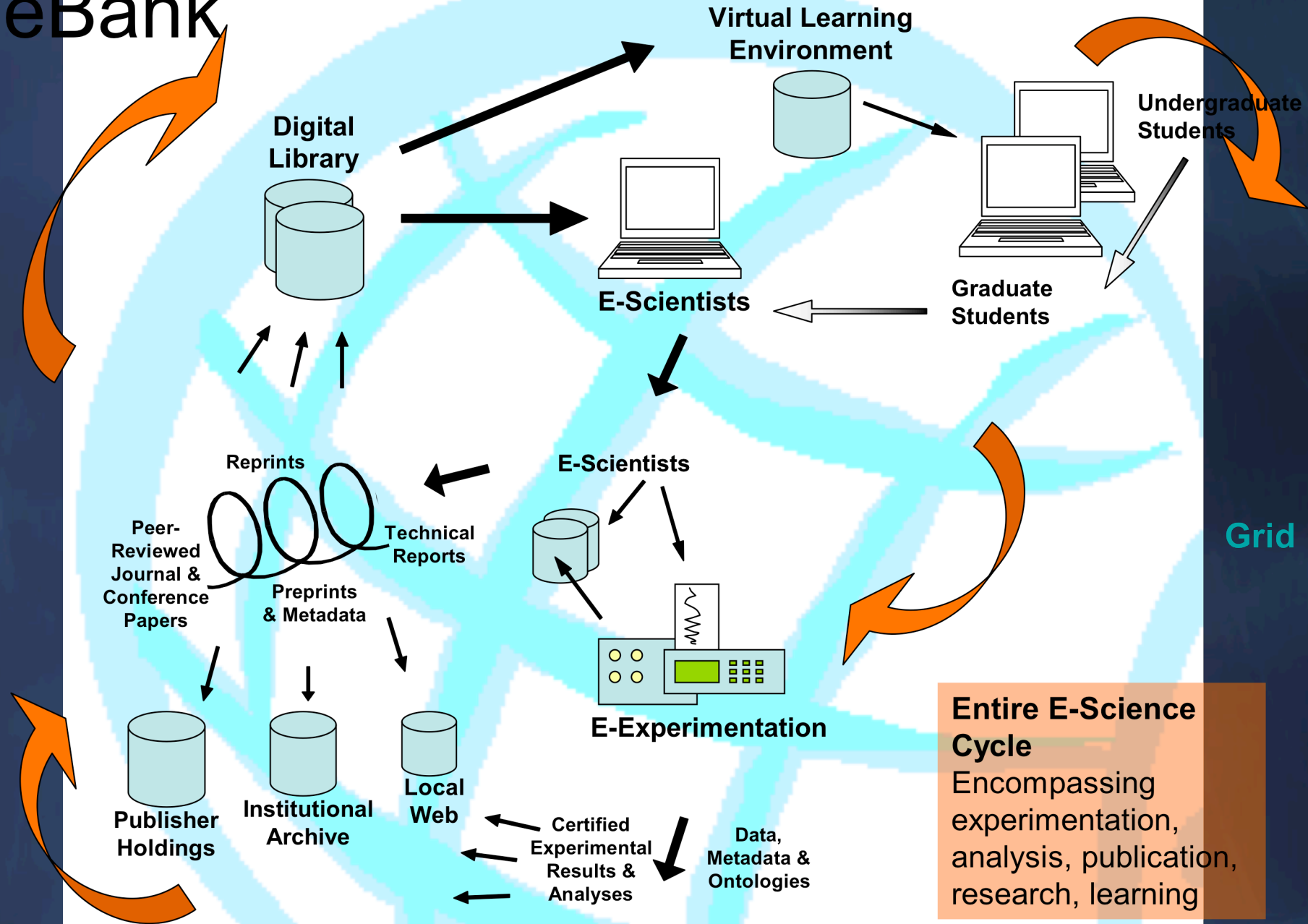
# Autonomic e-Science

- ✦ Built simulator of a future combechem in which 1000s of services are negotiating and self-organising
- ✦ Informed by combechem experience
- ✦ Article in IEEE Intelligent Systems on the Self-Organising Semantic Grid
- ✦ Raises questions about the future role of the scientist

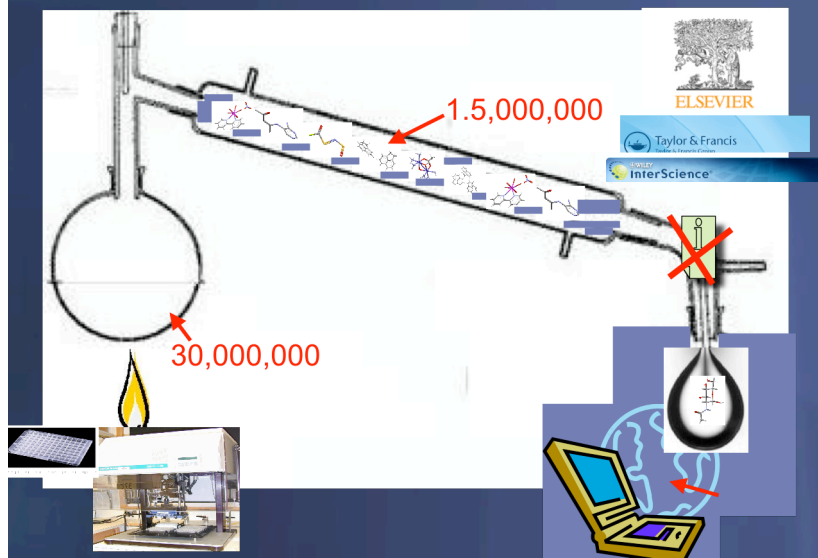




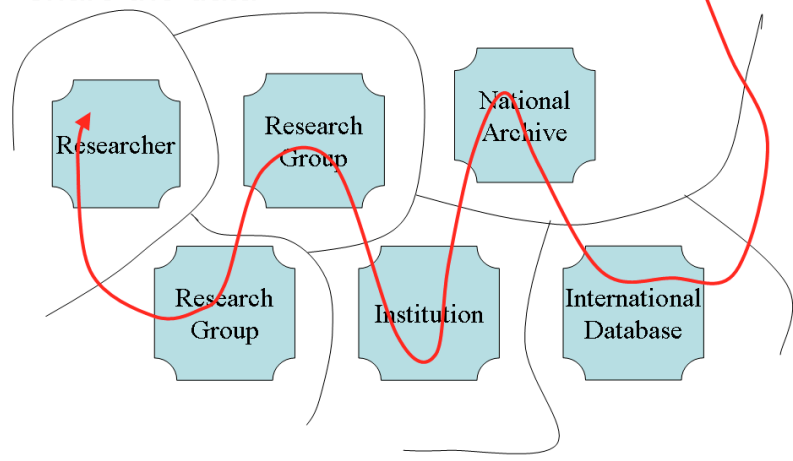
# eBank



# Access to the underlying data



Several groups making and analysing the library Administrative Domains transfer or share the data



**University of Southampton Crystal Structure Report Archive**

6,7,9,10,12,13,15,16-Octahydro-benzo-1,4,7,10,13-pentaoxacyclopentadecin

Simon J Coles, Michael B Hursthouse, Jeremy G Frey and Esther Rousay  
University of Southampton  
CuhbaOe

InChI:1C14H20O5/c1-2,4-14-13(3-1)18-11-9-16-7-5-15-6-8-17-10-12-19-14/h1-4H-12H2  
DOI: 10.594/crystals.chem.soton.ac.uk/145  
Compound Class: Organic  
Keywords: crown ethers  
Creation Date: 07 October 2004  
Deposited By: A.N. Admin  
Deposited On: 20 February 2006

**Available Files**

File Name	Size
02sot064.CIF	19k
02sot064.cml	8k
02sot064_checkcif.html	14k
02sot064.RES	9k
02sot064.PRQ	5k
02sot064.HTM	6k
02sot064.HKL	338k

**Refinement**

Parameter	Value
Cell angle beta	90.00
Cell angle gamma	90.00
Data collection temperature	120(2)

**Solution**

Parameter	Value
Solution figure of merit	0.0573
R Factor (Obs)	0.1185
Weighted R Factor (Obs)	0.1046
Weighted R Factor (All)	0.1243

**Processing**

Parameter	Value
02sot064.HTM	6k
02sot064.HKL	338k

**Other Files**

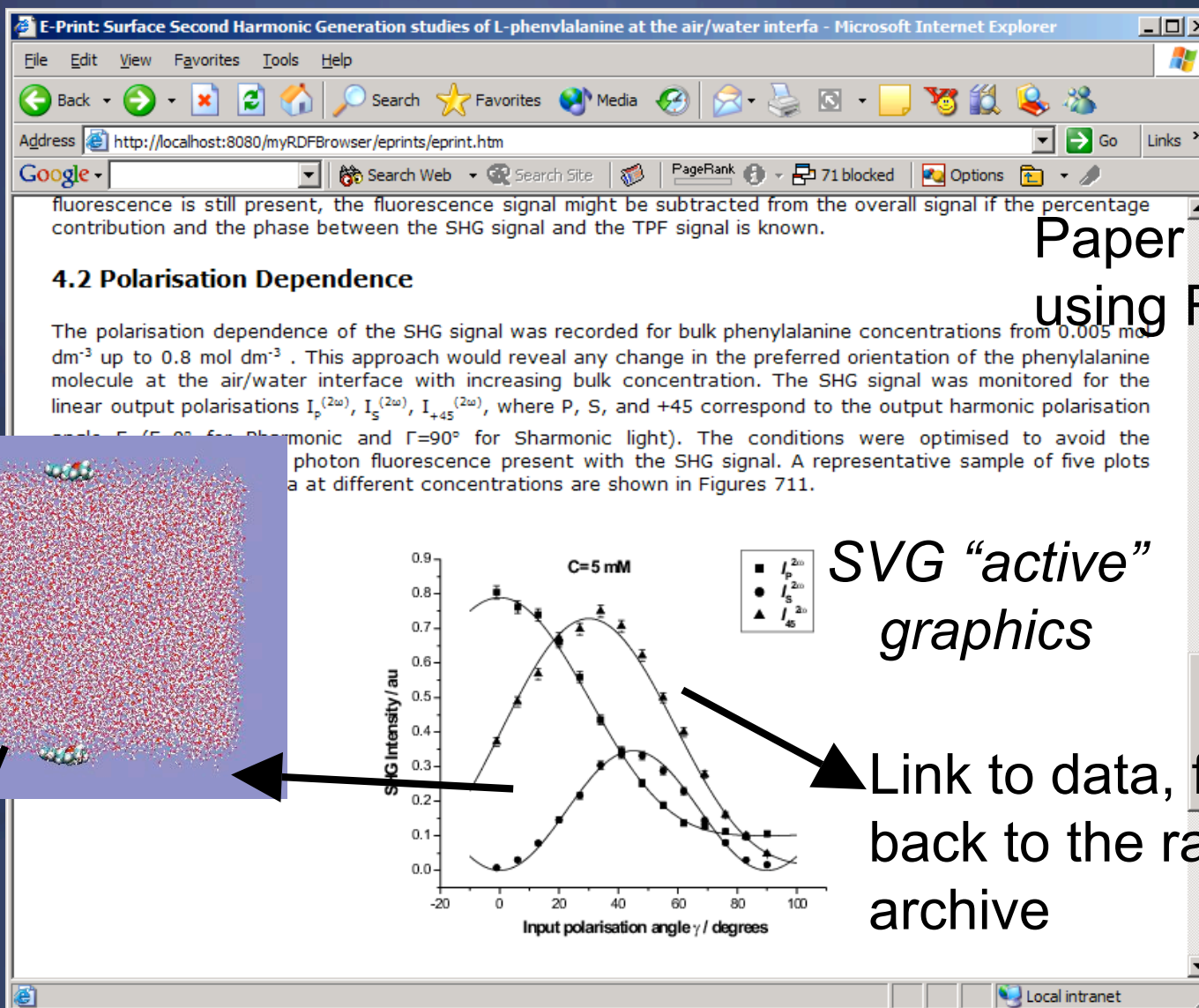
File Name	Size
02sot064.HKL	338k

**checkCIF/PLATON report (full structural check)**

**Summary report for Directory: disks\02sot082**

Unit cell

Parameter	Value
a	10.111(1)
b	10.111(1)
c	10.111(1)
alpha	90.000
beta	90.000
gamma	90.000
Volume	1011.1(1)
Z	1
Density	1.45
Calculated density	1.45
Flack parameter	0.00(1)
Goodness of fit	1.000
R factor	0.1185
Weighted R factor	0.1046
Maximum difference	0.1243



Paper organized  
using RDF

SVG "active"  
graphics

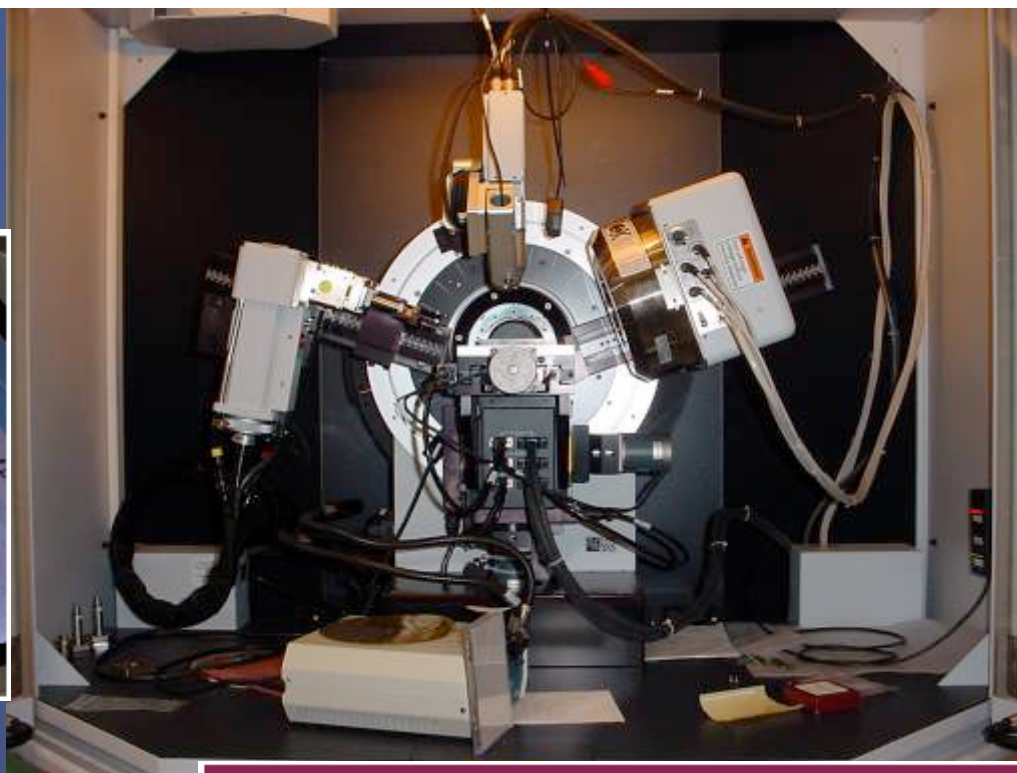
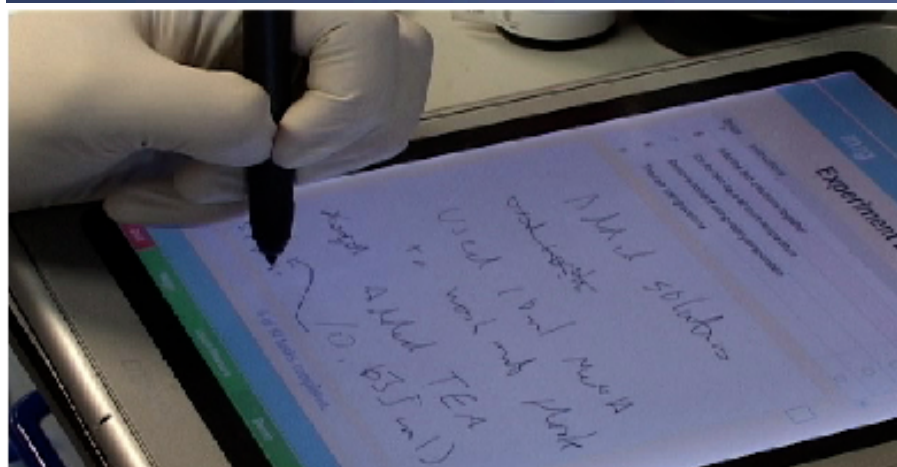
Link to data, follow links  
back to the raw data  
archive

Link to simulation, full  
simulation data archived  
in BioSimGrid

R4L



# Data capture



**R4L** Repository for  
the Laboratory



the Smart Tea Project



"I can go anywhere and its, like, this is  
me and my data. Its all there, bang."

- Chris,  
a real chemist, on using Smart Tea  
instead of a paper lab book.

Smart Tea is about improving the information environment for chemists doing chemistry - within and beyond the lab. Smart Tea is about supporting chemists in the preparation, execution, analysis and dissemination of their experimental work.

# Take Homes


- ★ Making sure other people can find, understand and re-use your data easily and with confidence (even when there is a huge amount of it!)
- ★ Whole lifecycle approach from lab to publication
- ★ Significant rollout of next generation Web technologies – Semantic DataGrid
- ★ Distinctive in e-Science for focusing on laboratory, usability and collaboration
- ★ Agent of culture shift in publishing and open access to data
- ★ Outreach including schools
- ★ Platform and agenda for future research in the Pervasive Semantic Grid

www.combechem.org

CombeChem

http://www.combechem.org/ Google

The Cloud N...on manager Theme 02 -...- homepage Examples an...s for PH331 Design of E...ments (DOE) Electronic L...eplace Paper Apple (233) Amazon eBay Yahoo! >>



# CombeChem

**Home**

**About**

**Tour**


**Events**

**Publications**

**Downloads**

**Links**

**Contact**

**CombeSearch**  
powered by 

**Site Map**

Find out [about](#) Combechem and take a [tour](#)

Visit CombeChem [events](#)

Access CombeChem [publications](#)

Download [talks and software](#)

Follow [links](#) to related projects

See [who](#) is involved and how to contact them

**News**

See the articles on [mobile phones in the bar](#) and [e-Malaria](#) on the BBC News Web Site.

*Semantic Datagrid* paper accepted at [6th IEEE/ACM International Workshop on Grid Computing](#).

**Latest CombeChem Related E-Prints**

**ECS E-Prints**

De Roure, D., Frey, J., Michaelides, D. and Page, K. (2006) [The Collaborative Semantic Grid](#). In *Proceedings of 2006 International Symposium on Collaborative Technologies and Systems*, pp. 411-418, Las Vegas, USA. Smari, W. W. and McQuay, W., Eds.

Miles, S. (2006) [Electronically Querying for the Provenance of Entities](#). In *Proceedings of Third International Provenance and Annotation Workshop* (in press), Chicago.

**University of Southampton E-Prints**

Frey, Jeremy G. (2006) [Future Lab - "Smart not Dark"](#). In, *Smart Lab: Laboratory Informatics Exchange, Jumeirah Carlton Tower, London, UK, 15-16 Feb 2006*. UK, International Quality and Productivity Center (IQPC).

Coles, Simon, Frey, Jeremy, Hursthouse, Michael, Light, Mark, Carr, Leslie, DeRoure, David, Gutteridge, Christopher, Mills, Hugo, Meacham, Ken, Surridge, Mike, Lyon, Liz, Heery, Rachel, Duke, Monica and Dav... Michael (2006) [The 'end to end' crystallographic](#)



# Questions

# Grid Innovation

- ★ CombeChem has focused on accelerating science by accelerating the process and not necessarily the computation
  - ★ Uses existing cluster and grid techniques
  - ★ Early focus on security for National Crystallographic Service
  - ★ Adopted Web Services from the outset
  - ★ Uses asynchronous message passing for integration
- ★ Semantic DataGrid

# Middleware Outputs

- ★ Security and access control developed for NCS
- ★ Software written by IT Innovation for CombeChem fed into the software distribution for the EU Grid for Industrial Applications (GRIA) Project
- ★ It forked!
  - ★ GRIA now on release 5, good adoption by industrials in EU projects (e.g. SIMDAT)
  - ★ Solutions evolved with current Grid standards

# Other Outputs

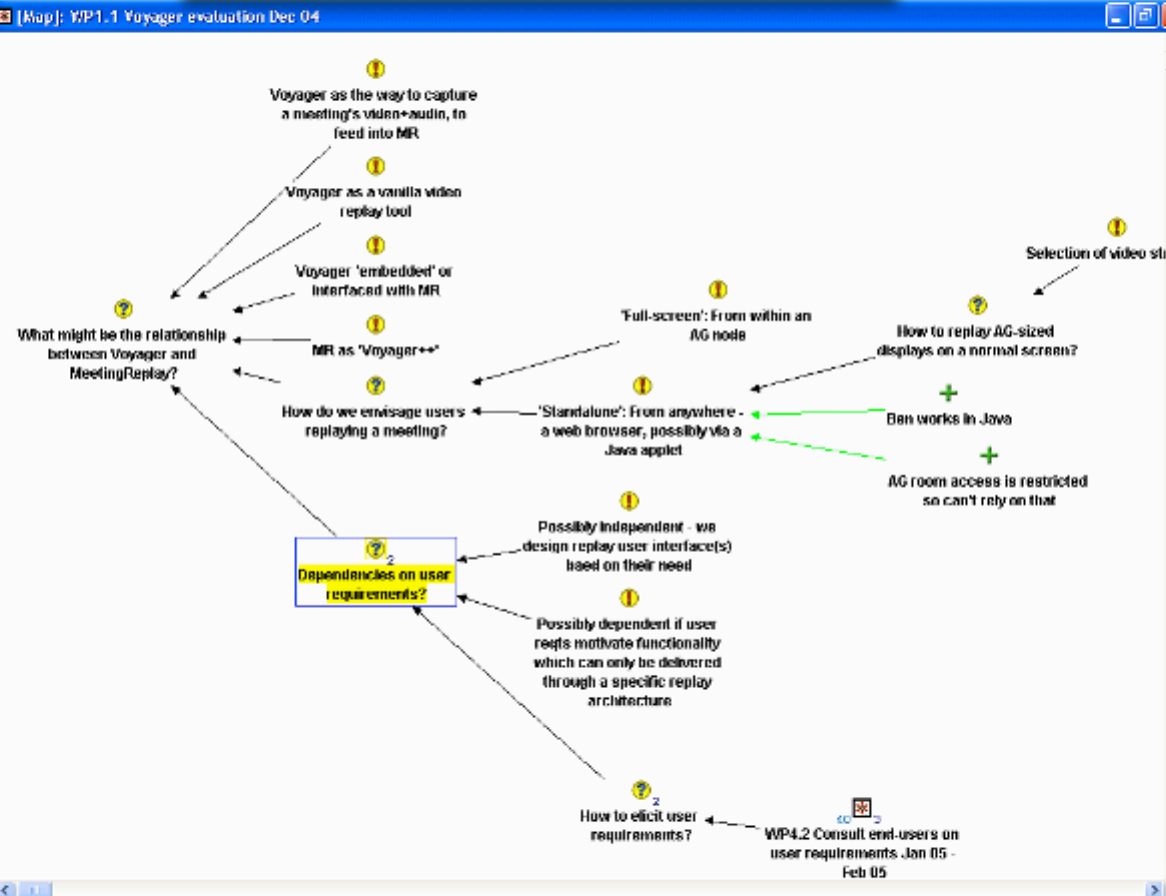
- ✦ Security and Access Control in GRIA 5
- ✦ Statistics software:
  - ✦ Design search algorithms for Generalized Linear Models
  - ✦ Design of experiments eLearning module
  - ✦ Elicitation in Chemistry Investigations (EliCIT)
- ✦ RDF streaming tools
- ✦ Units Ontology

# Staffing

- ✦ Deploy-then-research strategy
- ✦ Core team persisted through most of project and developed interdisciplinary knowledge
- ✦ Brought in additional staff for specific tasks

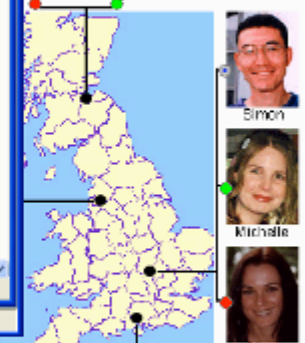
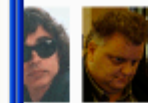


# MEMETIC

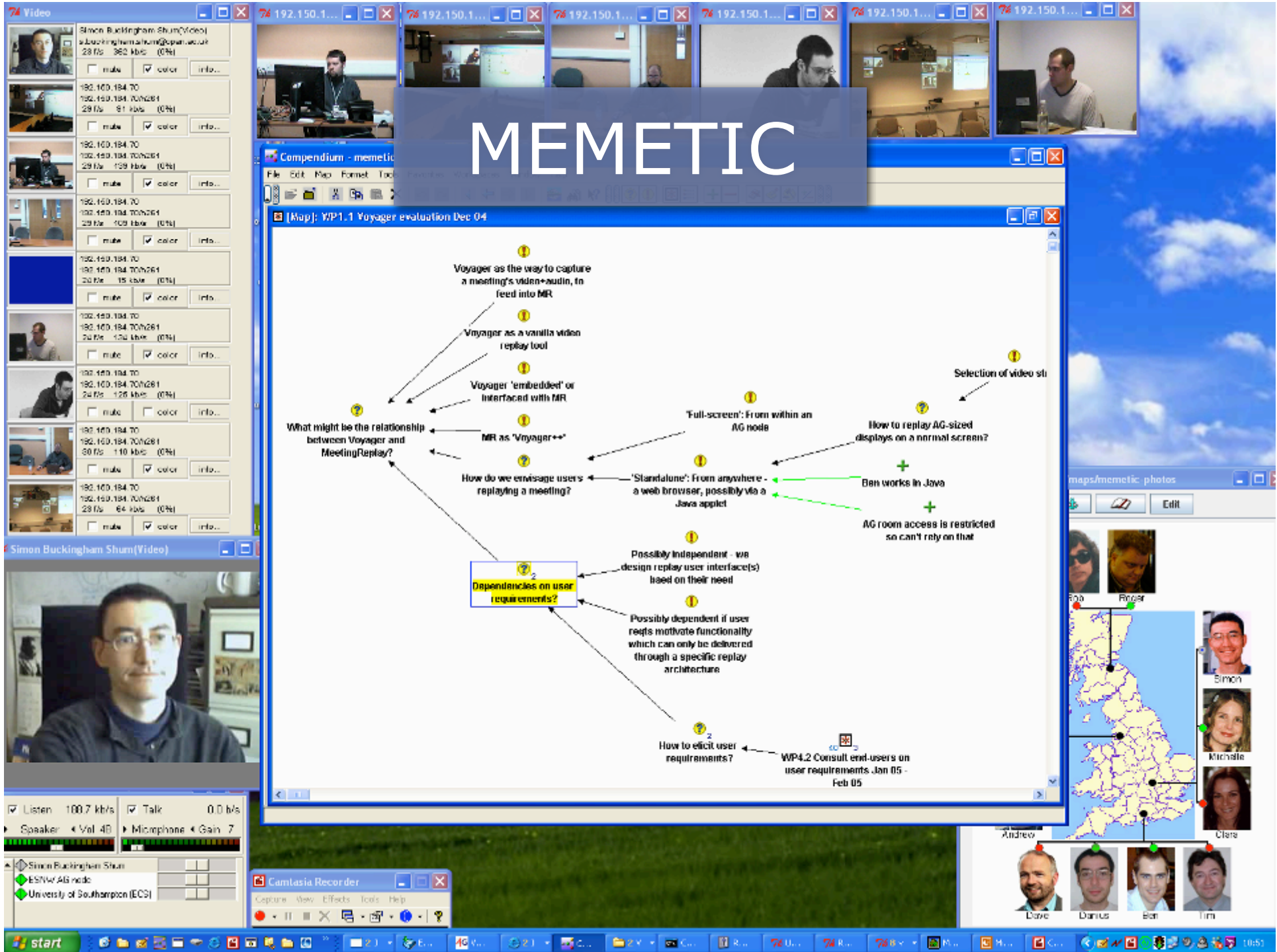
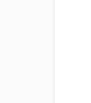


maps/memetic\_photos

Edit



Andrew

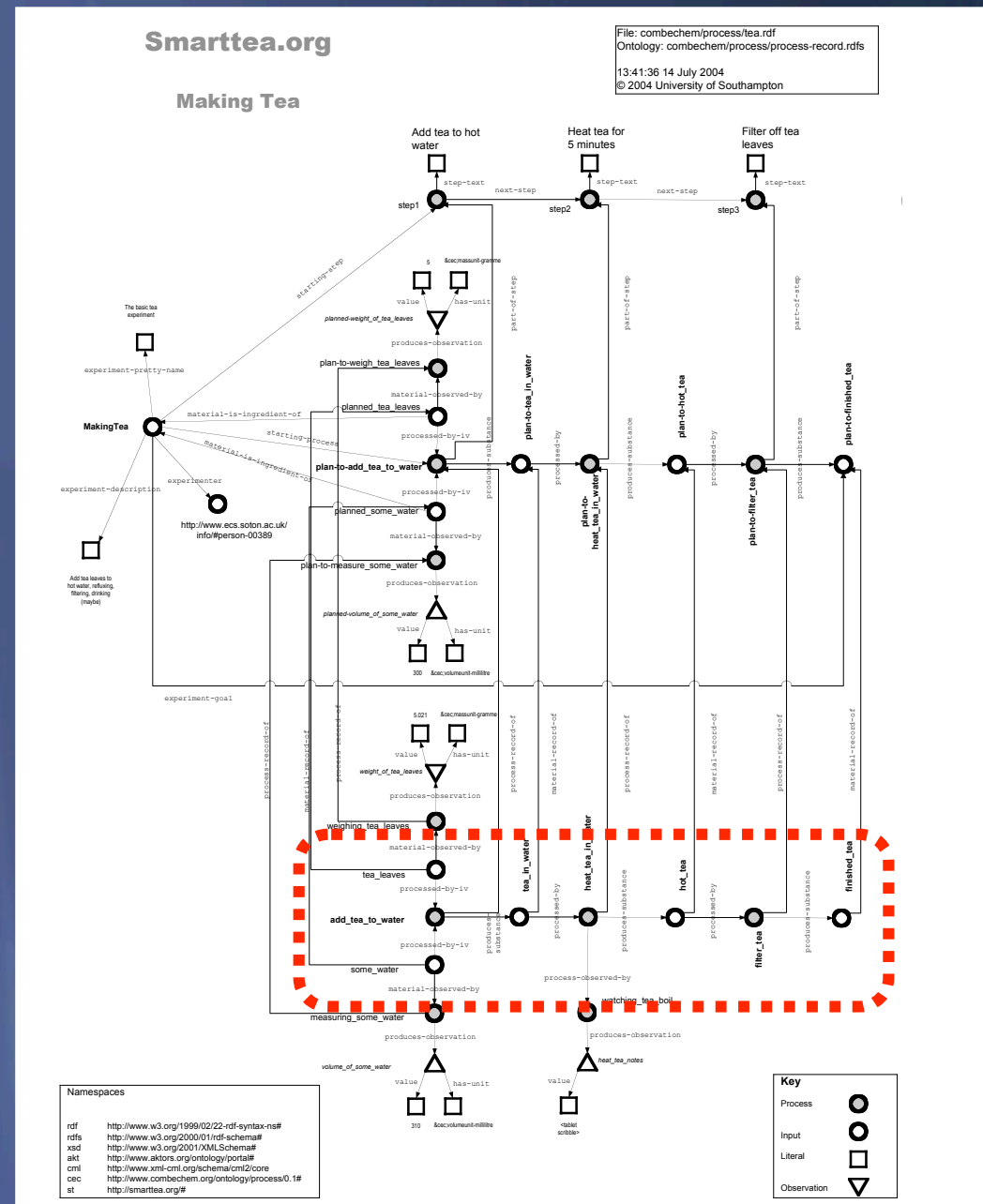




getRecord()

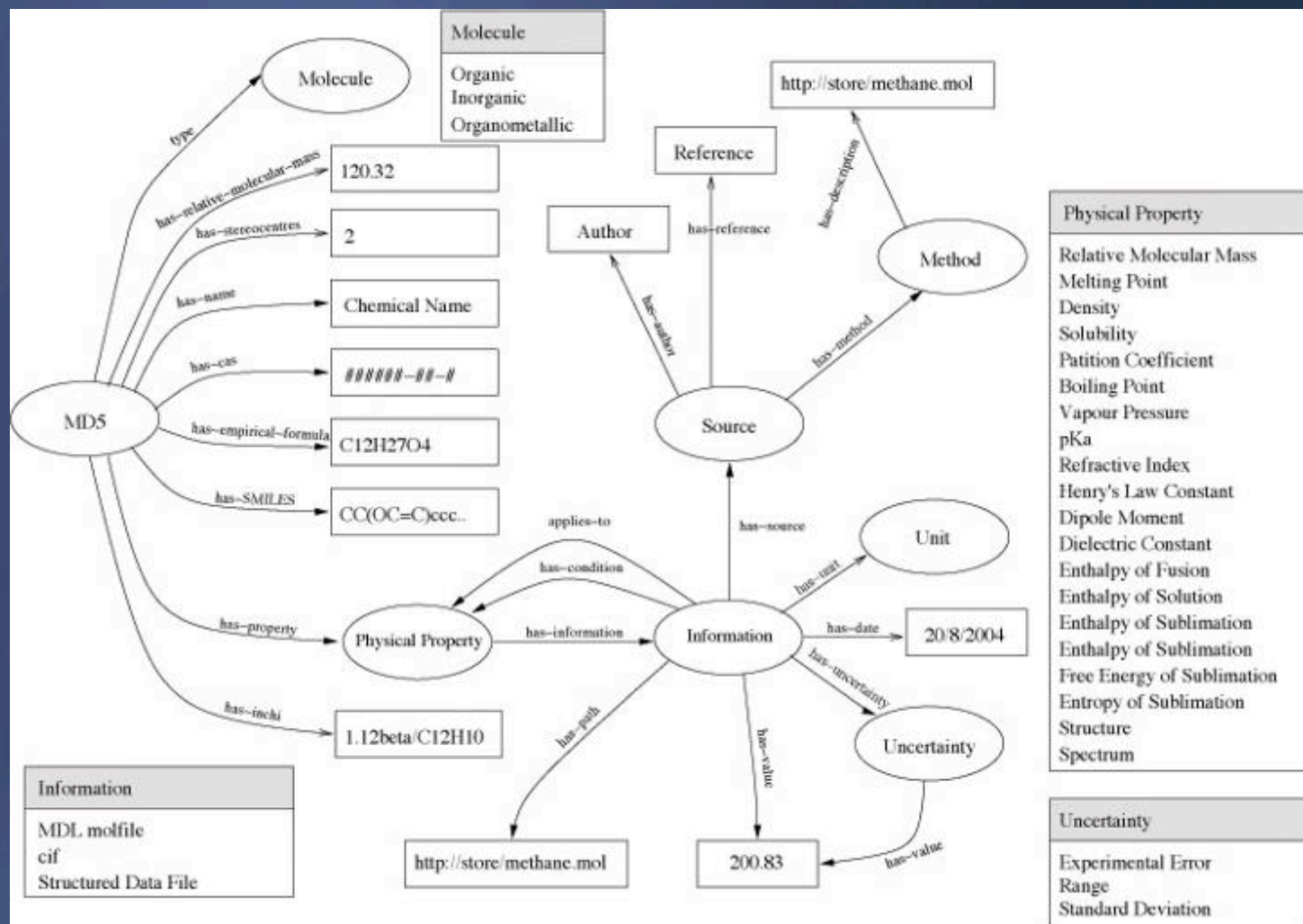
There is a potential containment problem in pulling back partial RDF graphs from the triple store.

Solved by using multiple triple stores but boundaries are a major issue for the future.



# RDF/RDFS

## High level Schema for chemical properties



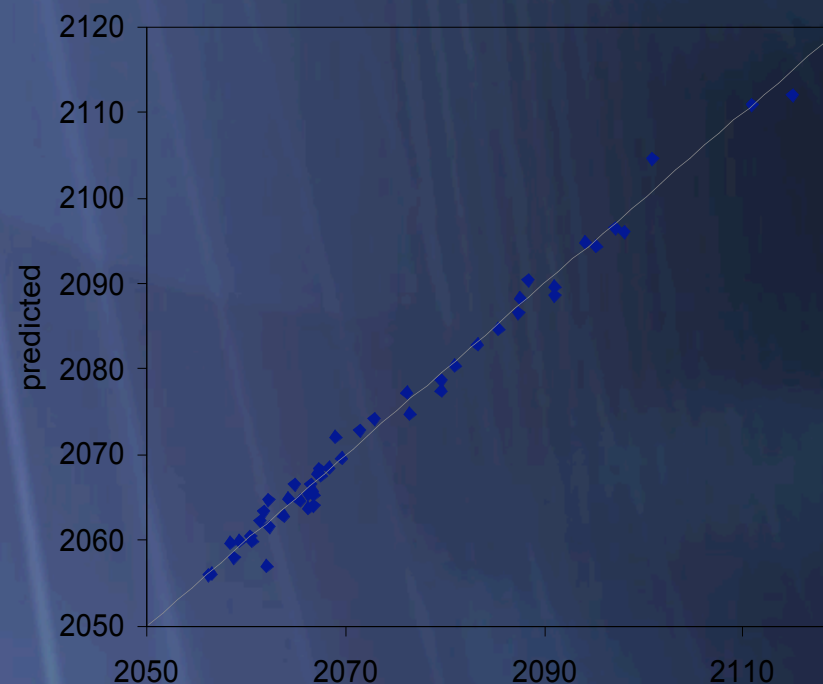
# Statistics

1. Identify and screen new catalysts *in silico* –prediction of desirable properties.
2. Direct experimental screening (high-throughput).
3. Detect and quantify ligand similarities/differences.
4. Add to chemical knowledge – interpret ligand contributions to experimental observations.

## Potential applications of ligand maps

15 Sept 06

Tolman Electronic Parameter ( $\text{cm}^{-1}$ )  
( $\nu_{\text{CO}}$  in  $\text{Ni}(\text{CO})_3\text{L}$ )

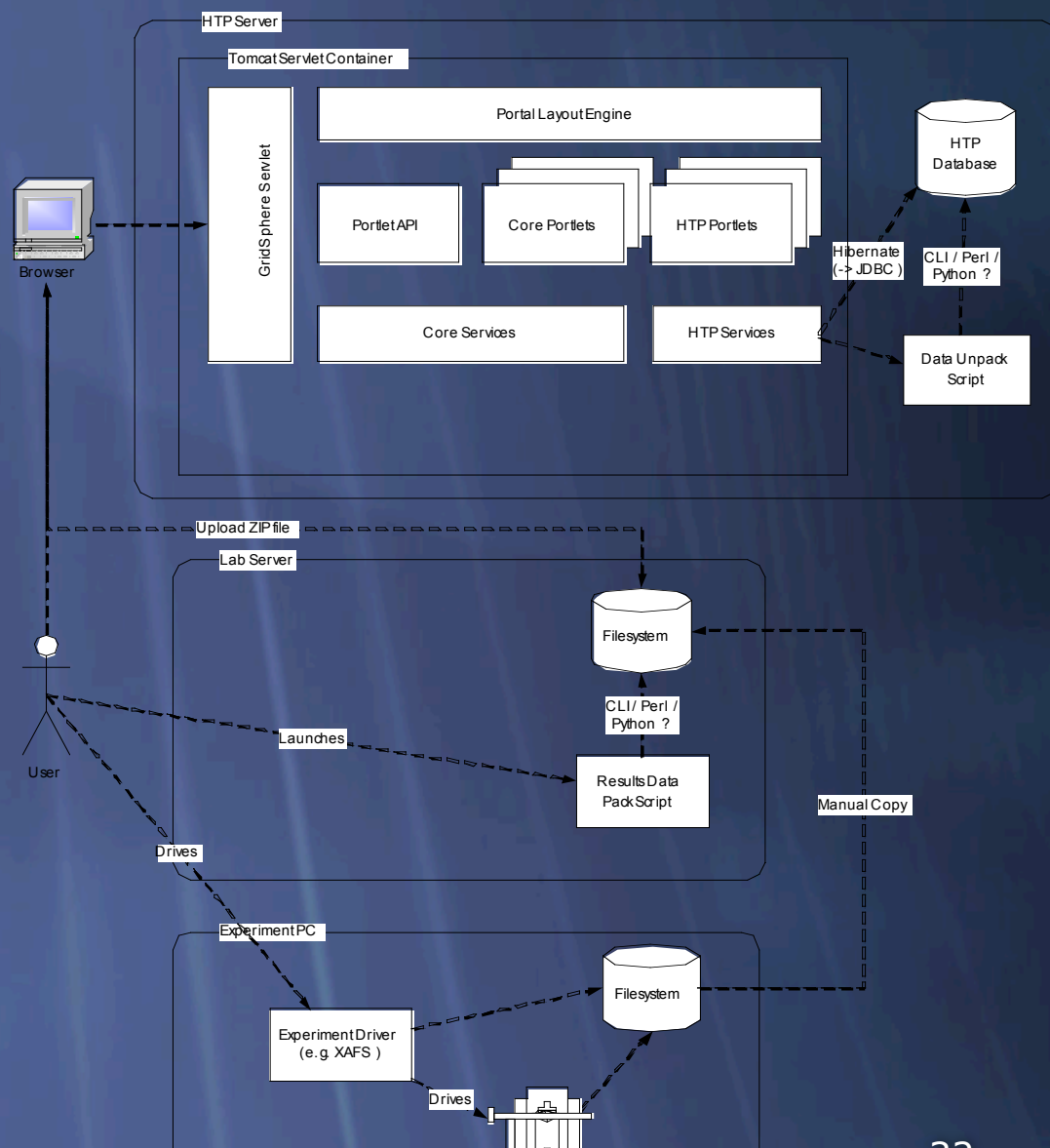


Descriptors: PA, %s, Q(Pt fragm.),  $\text{He}_8_{\text{steric}}$ , P-B, P-Pt,  $\Delta R\text{-P-R(Pd)}$

# HTP Sample Tracking

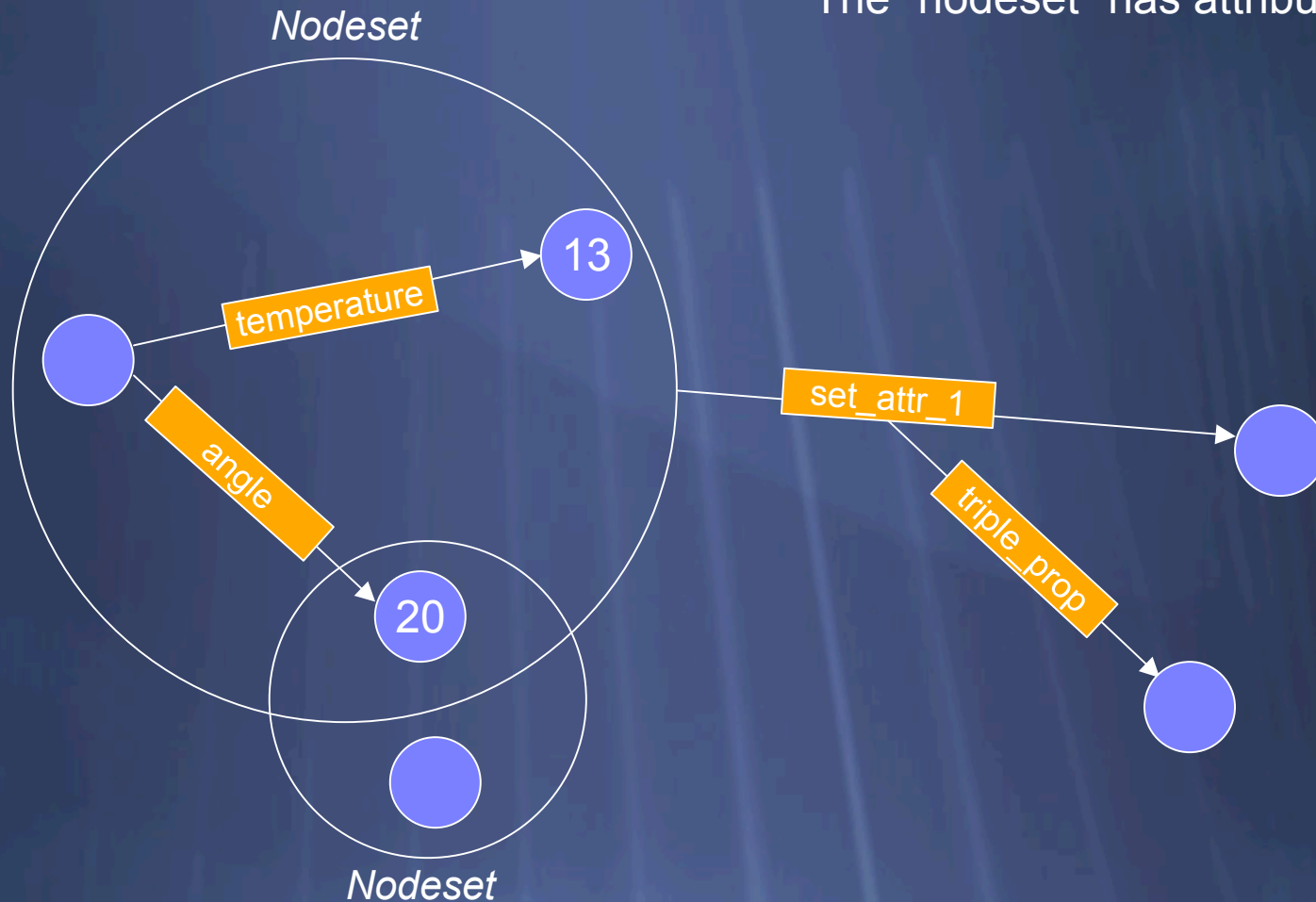
Using ideas from the NCS Grid Service we have produced a prototype for a high throughput catalyst experiment involving array samples investigated by Raman, MS, EXAFS with the samples manufactured at one site and tested at several others

HTP Architecture (First Prototype )





The “nodeset” has attributes



- ★ The edge with the attribute name `set_attr_1` is an attribute of a nodeset.
- ★ The edge with the attribute name `triple_prop` is an attribute of the above edge.

# Grid and Pervasive Computing

- ✦ Electronic Lab Notebook
- ✦ Lab Environment
- ✦ Mobile Devices
- ✦ Semantic throughout

