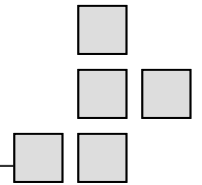# Open Discovery Workflows

# Beyond Pipelining

Jonathan Sheldon, CSO

(LifeScience and Healthcare)

InforSense Ltd
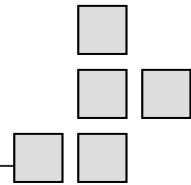
# InforSense Introduction

London (HQ)

Shanghai

European Field Office

Japanese Distributors

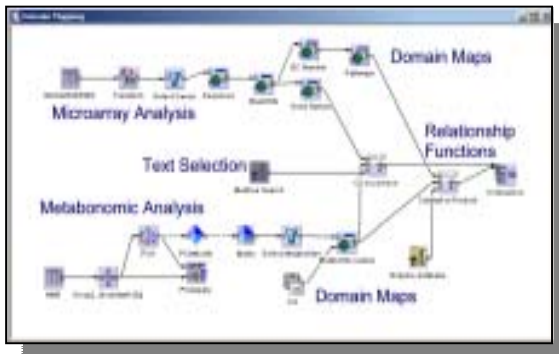Cambridge, MA

Turkey
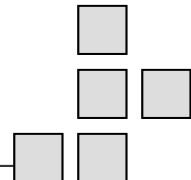Distributor

# InforSense Technology Foundation

Workflow based portal building technology:
Automatically deploying workflow composing
services to an enterprise portal
***BioWorld IT 2005 Best of the Show Award -
Knowledge Management & Collaboration***

Workflow technology :
Building applications by
composing services
***IEEE Supercomputing 2002 :
Most Innovative Award***

Web Services based informatics
infrastructure: Harness multi-
domain distributed resources
***IEEE Supercomputing 1998 :
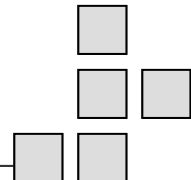Most Innovative Award***

# Workflow is not new…



with certain orders

des and

– Application Integra

# Building a workflow

**InforSense**
The Integrative Analytics Company

**1 Select:**

- Data Sources
- Data Mining / Statistics
- Data Processing / Transformation
- SQL Functions
- Interactive data visualization / reporting

**2 Connect:**

- Connect data and components in GUI
- Workflow describes complex data processing and analysis
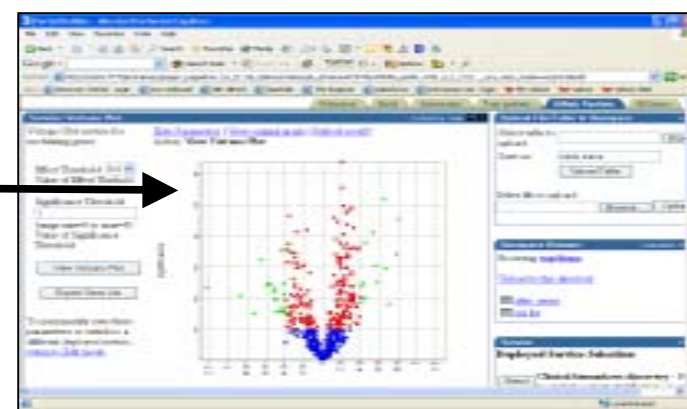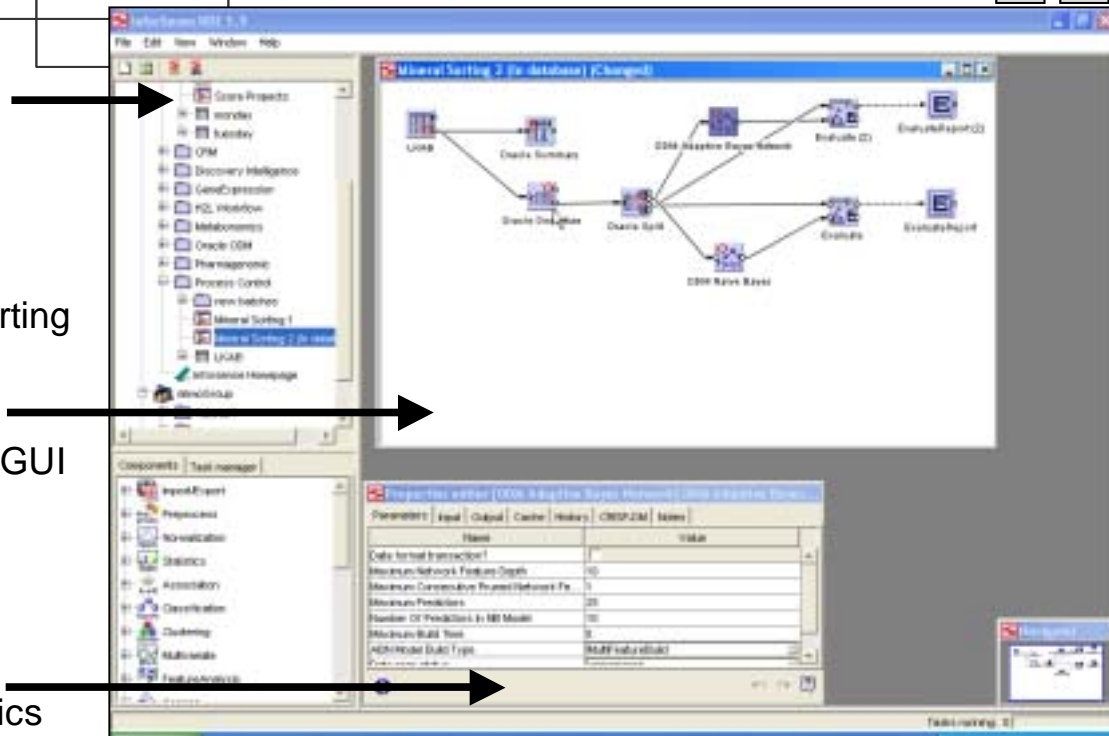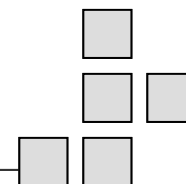
**3 Execute:**

- "In database" processing & analytics

**4 Deploy:**

- Define parameters of workflow to expose
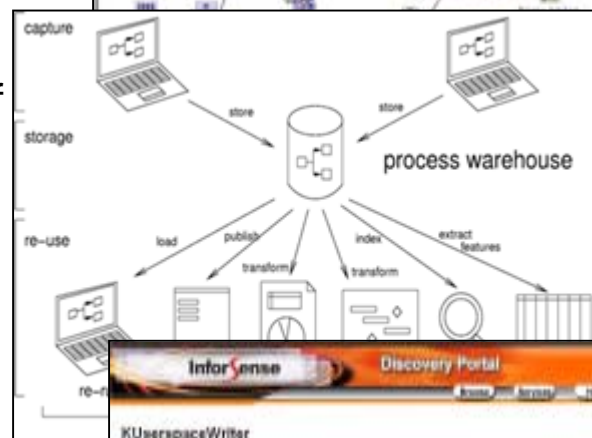- Publish as: portlet, web application, SOAP service, command line app
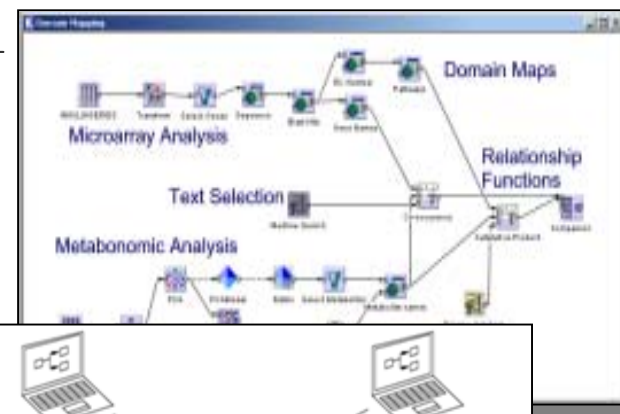
**5 Warehouse:**

- Store workflows to capture ad-hoc querying and establish best practices

- **Constructing a ubiquitous workflow:**
  **by scientists**

  – Integrate your information resources/software applications cross-domain

  – Support innovation and capture the best practice of your scientific research

- **Warehousing workflows:**
  **for scientists**

  – Manage discovery processes in your organisation

  – Construct an enterprise process knowledge bank

- **Deployment workflow:**
  **to scientists**

  – Turn your workflows into reusable applications

  – Turn every scientist into a solution builder

InforSense
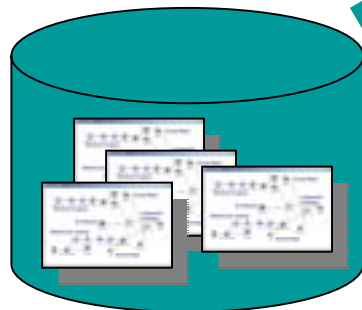The Integrative Analytics Company

Ubiquitous workflows, warehouses and deployment engines

**Domain specific tools**

**Workflows**

**Warehousing**

**Deployment**

**Mine warehouse for provenance, information and knowledge**

Experiment & design — CambridgeSoft Life Science Enterprise Solutions

Enumeration & descriptors — Optive RESEARCH

Lab & analytical data — MICROMASS MS TECHNOLOGIES / Thermo ELECTRON CORPORATION

Compound & activity data — ID Business Solutions / MDL Information Systems, Inc.

Database & documents — documentum / ORACLE / NuGenesis

Visualization — Spotfire

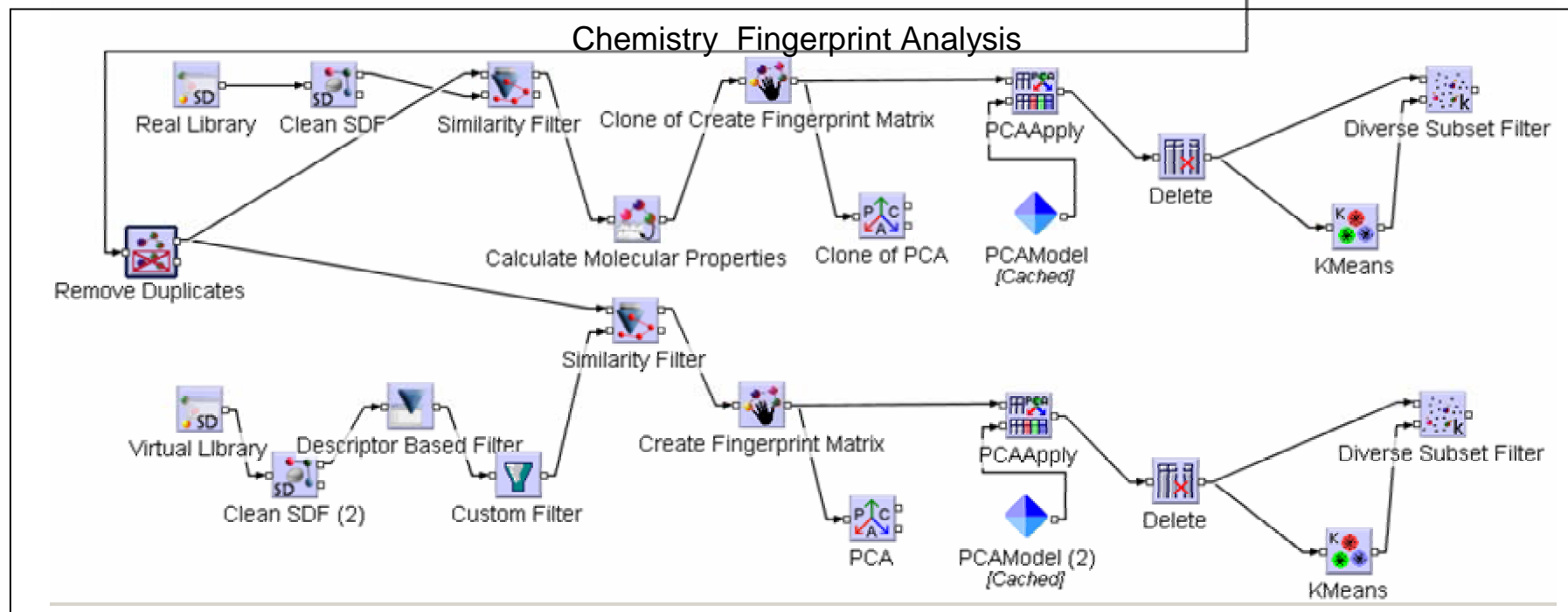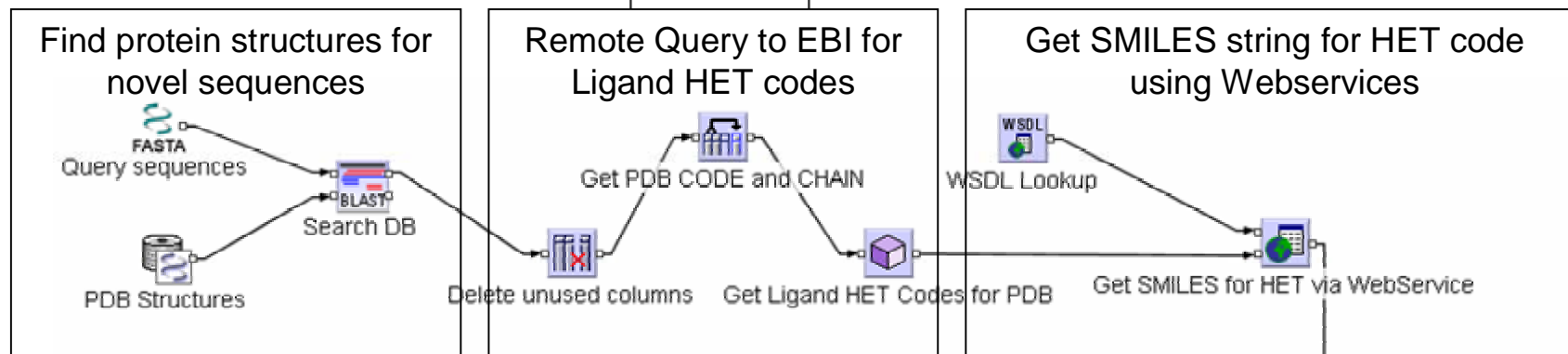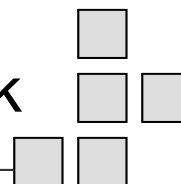**TextSense ChemScience BioScience** — InforSense analytical systems

# There is something new about the use of workflow in scientific informatics…
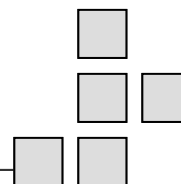
- Workflow is an "machine executable protocol" for in silico scientific activities
- Thus, a scientific workflow needs to be :
  - **Ubiquitous** : scientific data are heterogamous and scientific activities are rich in variety
  - **Integrative** : resources (data, software, instruments) need to be dynamically combined for scientific study
  - **Personal** : scientists need total freedom in defining and building their creative scientific process
  - **Interactive** : building a scientific workflow with a " trial and error" approach
  - **Open** : workflow should be able to integrate components from any vendors
  - **Reusable** : workflow should be reused by others in various form
  - **Manageable** : workflow represents process knowledge so it evolves overtime

**Ubiquitous** : Manage the complete workflow of a project in one framework

Find protein structures for novel sequences

Remote Query to EBI for Ligand HET codes

Get SMILES string for HET code using Webservices

Chemistry Fingerprint Analysis

# Integrative : Dynamic cross domain integration of databases & applications
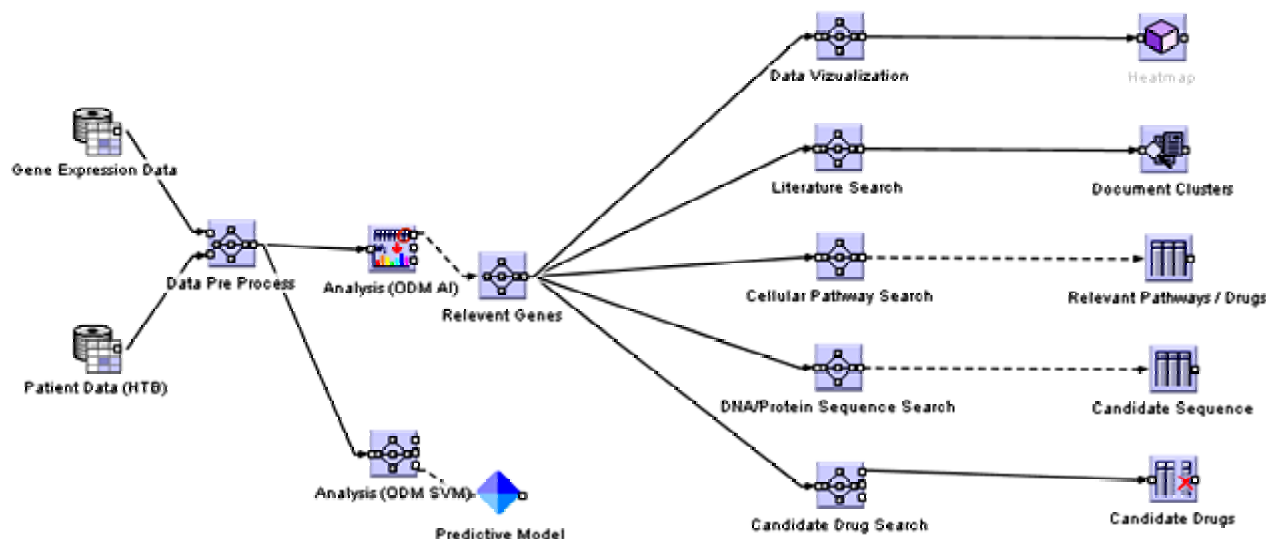
- ## Used Modules in KDE
  - 60 nodes representing 5 x DBs[1] and Oracle components[2], 3 x 3rd party applications[3], 4 x public DBs[4] and 1 web service[5]
  - Gain: 27 times faster (10 vs. 270 Man Days), reusable and sharable

## The Workflow
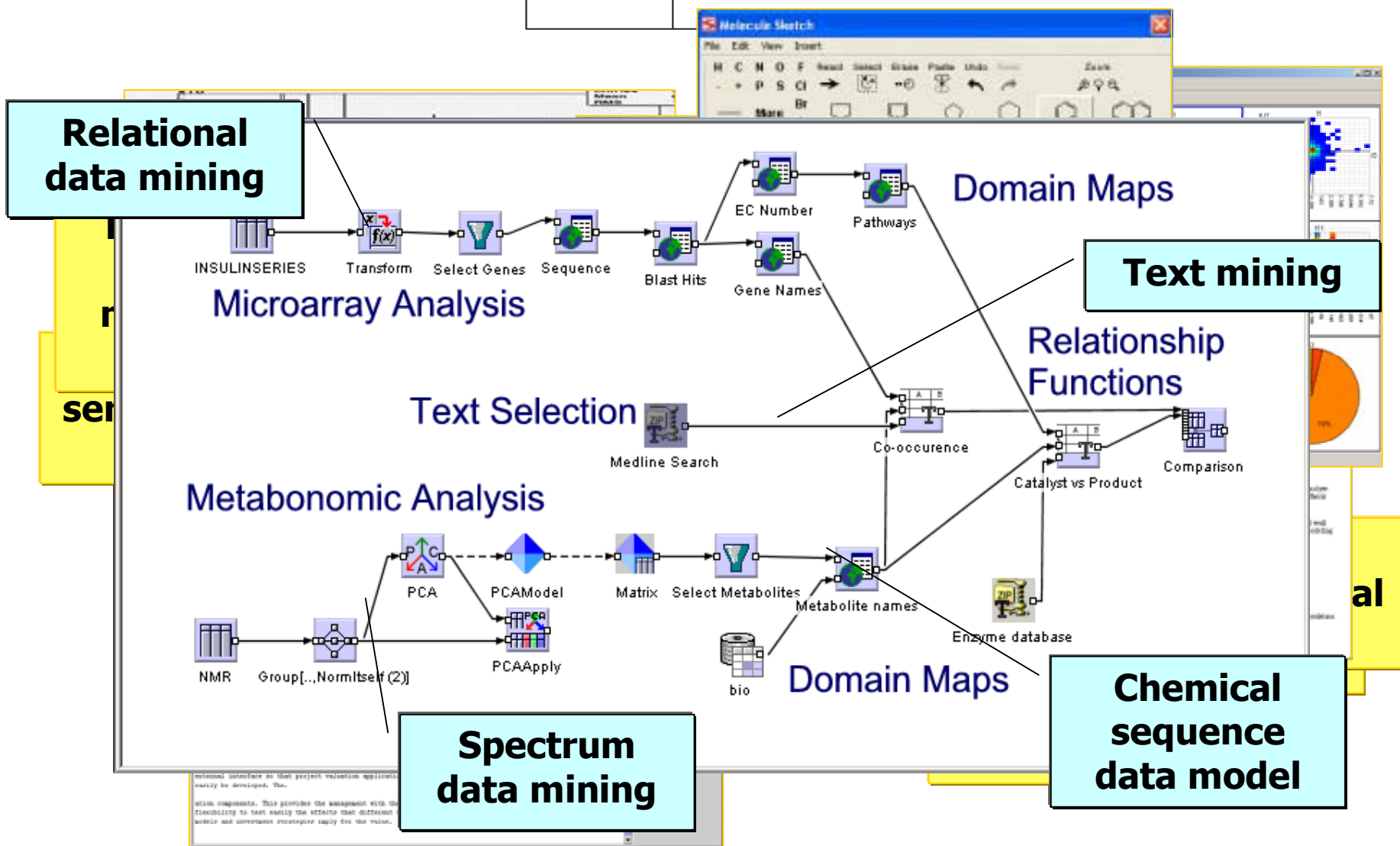
Data Acquisition                Multiple Types of Analysis                Relevant Results
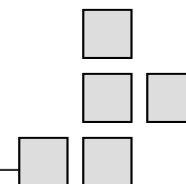


[1] Gene Expression, Patient, [2] ODM (SVM, AI, BLAST, TEXT), [3] Chemical Editor, Blast Viewer, HeatMap, [4] ChemBank, PubMed, NCBI, Beilstein, [5] KEGG

**Interactive** : Building Complex Discovery Process with Visualisation

**Relational data mining**

**Text mining**

**Spectrum data mining**

**Chemical sequence data model**

# Personal : Building and deploying personal applications

**Open** : Working with any vendor

- # Higher level orchestration of workflows
  - conditional execution, loops..

- # 'Useful' error reporting
  - exactly what has broken !

- # Cached results
  - resilience, debugging

- # Asynchronous jobs
  - long running processes

- # In database execution
  - speed, security
- # Audit trial
  - who, what, when
- # Meta data level consistency & error checking
  - construction of workflows with compatible data types
- # Support for SOA
  - web services
- # Enterprise infrastructure
  - cluster support, server admin tools

# Case Study:
# Library design

## Collaboration with GSK Cheminformatics

## Historically
- Low volume
  - 30-50 cmpds/yr/chemist: 10,000s assay wells/yr
- Low information diversity
  - scientists generally dealt with limited types of data
- Reductionist approach
  - limited information per experiment
    - Interpretation critical for next step
  - scientists required:
    - simple systems to assist in information monitoring
    - decision making resides with the scientist

Slides from their e-science presentation at the EPSRC, April 2005.

- ## **What happened in the last 5 years?**
  - ### "Industrialisation" - Application of "principles of industrialisation" to drug discovery
    - high volume
      - 10,000 cmpd/yr/chemist/100+ million wells/yr
  - ### Biology revolution
    - Human genome
      - "system biology" – holistic view and interpretation
      - high content data --- images
      - multiple result types from each experiment – bio-markers, pathways
  - ### Knowledge integration
    - scientific discipline integration
  - ### Scientists required
    - complex systems, algorithms, statistics…….
    - decision making shared between systems and scientists
    - "Informatics" essential – partnership not service

# How has pharma IT tackled the transition?

- ## Business as usual
  - ### problem centric view
    - build applications
    - integrate applications

- ## Educate scientists in the realms of IT
  - "Now I need to be an IT expert alongside chemistry, biology, genetics, robotics, engineering ......"
  - interesting time scale - generations

- ## Technology is our saviour!
  - client server, web services, java, C#, Corba, OO programming, extreme programming, grid computing, .....

# What are the results?

**Chemistry**

**screening**

**data**

**"library" design**

**samples**

**infrastructure**

- **"islands" or "silos" of process & data**
  - *complex integration problem*
    - *"spaghetti" joins our worlds - unsustainable - cost*
- *control with "IT"*
  - *mismatch in cycle time to change*
  - *engineered out serendipity*
  - *service role reversed*

# Doodling in knowledge and experiment space

**physical & information worlds merge**

Information Resources

Target List & Status

Q: - are these results real?

Q: - what do I know about these compounds?

Q: - what other data can I acquire?

**Exclusion Lists**

**IC$_{50}$ Assay**

**Structure Validation**

**Other Assay…**

Q: - what other data can I acquire?

Target Leads

- no predefined steps
  - capture what was done, don't restrict what can be done?
    - don't restrict the non-obvious

# Integrative Analytics for Life Sciences: Building an Informatics Infrastructure

- **Objective:**
  - Dynamic cheminformatics infrastructure
- **Specification:**
  - Access to world-class scientific algorithms and tools
  - Access to disparate data sources, multiple locations
  - An intuitive GUI for medicinal chemists – not just IT experts
  - Build and deploy "just-in-time" applications
  - Capture and dissemination of "best practice"
- **Solution:**
  - InforSense KDE + IOE open SOA, integrating services, data, tools
  - OWPN partner integrations
  - Spotfire collaboration -> integrative visual analytics

KDEs Visual Network orchestrates the selection process across multiple tables and multiple Spotfire visualisations



- Toolbox: scientific models, chemical handling, chemical properties, data access, statistics, data visualisation, ....
- Scientists can doodle in chemical space
  - Capture how scientists made decisions
- New algorithms, data sources added in < 1 hour

# Case Study:
# Mapping the evolution of SARS

## Collaboration with the Shanghai Center for Bioinformation Technology (SCBIT)

# SARS

- By 2003, SARS was rapidly developing into a pandemic
- The WHO needed to quickly find a way to:
  - Rapidly establish the relationship between genomic variations and the biology of SARS
  - Characterise genomic variations (deletions, single nucleotide variations) during the epidemic
  - Rapidly aggregate and analyse multiple data sources using many different tools, requiring different analysis techniques, from multiple research groups world-wide, in multiple formats on multiple platforms
  - Publish and inform research and treatment groups world-wide to enable appropriate treatment regimes for each patient

# SARS: Genomic Analysis



Examining the variations in different strains

Retrieval of publicly available knowledge

Alignment performed on the Grid

# SARS: Proteomic Analysis

Isolating interesting genomic regions

Identifying relevant protein sequences

Observing the variations in the resulting protein

# SARS: Protein Annotation



- Annotation using applications from EMBOSS package

- Execution on different physical resources

- Annotations merged

# SARS: Genotype Clustering

- **5 nt motif that characterises the phylogenetically related genotypes**

He J-F *et al* (2004)
Science **303**, 1666-1669

# SARS: Delivering the knowledge

- ▶ **Access functions via the web**

- ▶ **Rapidly deploy workflows as web applications**
  - ▪ _**No programming required!**_



Build       Encapsulate       Deploy

# SARS: Delivering the knowledge



- **Updated as new datasets and analytical techniques become available**

- **Virtual Grid based research environment**

- The ability to rapidly integrate multiple, disparate data services, available on the net

- The ability to integrate over 200 required compute services, including grid service infrastructure for compute intensive tasks

- Enable the knowledge and insights gained to be rapidly and automatically deployed as new web services and portals

- Minimise the investment required in tools and services – maximising the reuse of existing resources

One integration platform from Discovery to Clinic, and bench to boardroom

Discovery — Development — Clinical — Manufacturing & Distribution — Sales & Marketing — Patient Care

| Target ID | Lead ID | Lead Opt | Tox | Animal Models | Phase I | Phase II | Phase III | Mfg | Distrib. | Sales & Mktg | Phase IV |

Integrate knowledge

VP of Research – Portfolio decision

RA Head – Project decision

Scientist – Compound decision

**Bioinformatics**
- Target discovery
- Target validation
- Biomarker discovery
- Systems Biology

**Cheminformatics**
- Lead discovery
- Lead Optimization
- Computational Chemistry
- Library Design
- ADME Tox

**IOE and Oracle HTB**

**Personalised Medicine**

Integrate cross domain processes

**Biomarker discovery**
(Chemistry – Biology – Clinical Data)

**Chemogenomics**
(Chemistry – Gene Expression – Toxicology)

**Translational Medicine**
**Patient Lifestyle/Environment and the study of Disease progression – From bedside to bench**