IST-5-033437

# The Chemomentum Data Services
## A flexible solution for data handling in UNICORE

Katharina Rasch, Robert Schöne, Hartmut Mix - Technische Universität Dresden, ZIH
Vitaliy Ostropytskyy, Werner Dubitzky – University of Dublin
Mathilde Romberg – Forschungszentrum Jülich

# Outline

- Chemomentum project overview

- Data management features

- Technical details

- User client
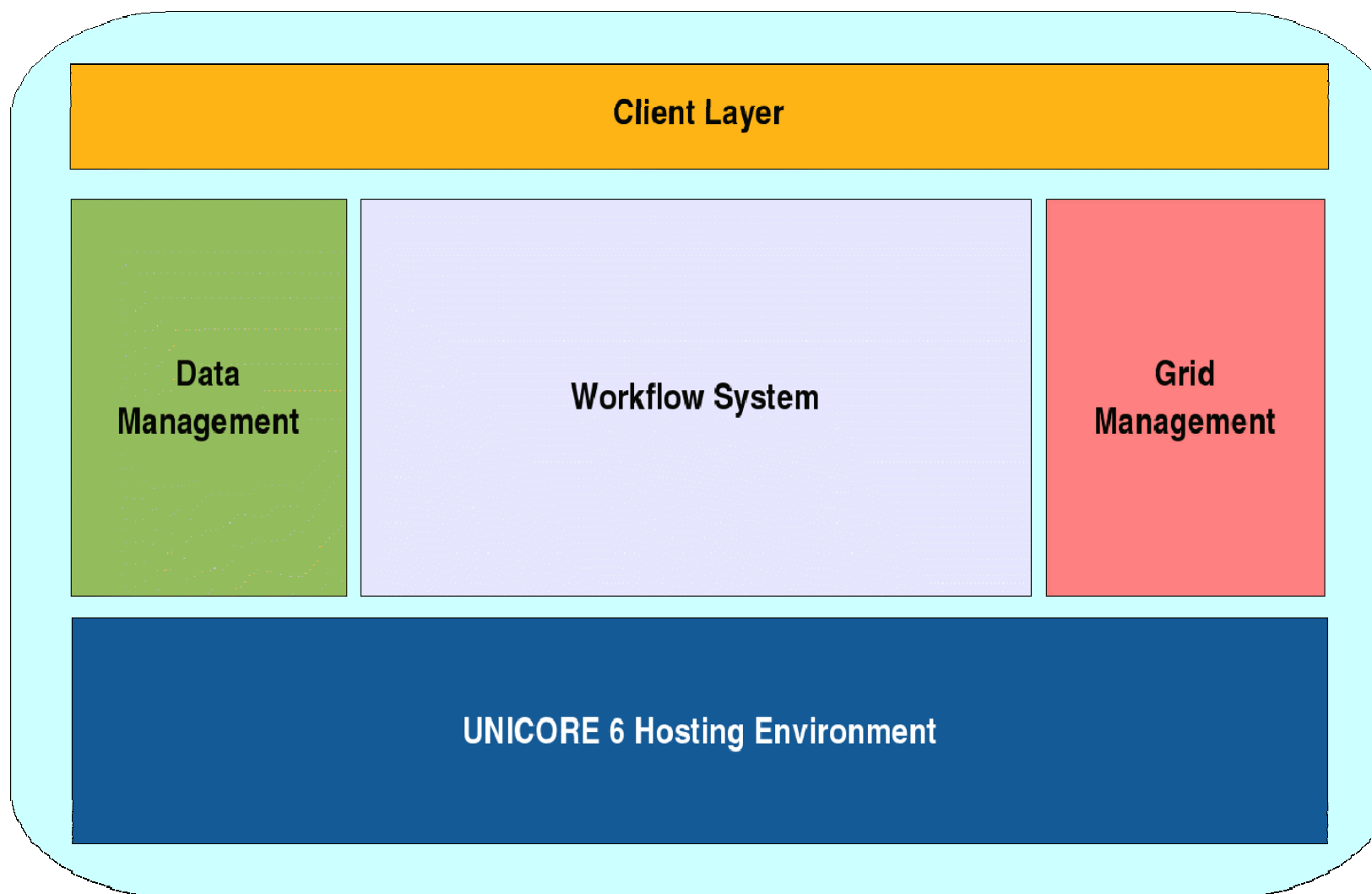
# Chemomentum project overview

- Generic, flexible system for running workflow-centric, complex applications

  e.g. computational chemistry, supply-chain management

- Deals efficiently with data and knowledge
- Focused on end users
- Use cases: drug discovery, toxicity prediction, environmental risk assessment, QSAR, protein docking
- Based on UNICORE Grid middleware
- Web site: www.chemomentum.org

# Chemomentum project overview

- 9 partners:
  - University of Warsaw, Poland (co-ordinator)
  - Research Centre Jülich, Germany
  - University of Tartu, Estonia
  - University of Technology Dresden, Germany
  - University of Ulster, United Kingdom
  - Istituto di Richerche Farmacologiche Mario Negri, Italy
  - University of Zurich, Switzerland
  - BioChemics Consulting SAS, France
  - TXT e-Solutions, Italy

- 30 month, started 01/07/2006

# The big picture

# Ambitions – Data Management

- Store data produced by workflows
  - → need metadata to retrieve data later
    - General metadata, e.g. owner, dates, applications used, workflow description
    - Domain specific metadata, e.g. chemical structures inspected
- Calculation results should be reproducable
  - → special attention to ensuring provenance of data

# Ambitions – Data Management

- Handle files and meta information produced by Chemomentum
  - Store result files and meta /provenance information
  - Browse through stored data
  - Update and delete data
- Provide access to external data sources (e.g. chemical databases)
- Use ontologies to improve search results

# Features – Data Management

- Grid storage system
  - Data identified by globally unique logical name
    - $\rightarrow$ global view of data
  - Data annotation with **extensible** meta/provenance data
  - Automatic metadata extraction
  - Distribution and replication
  - Seamless access to external data sources
  - Provide synonyms and unit conversion to improve request

# Features – Data Management

- Integrated into UNICORE/Chemomentum
  - Webservice based (using WSRFlite framework)
  - Workflow System uses data management to retrieve input files and store output files / meta information
- Integration into Chemomentum client
  - Query/browse through data and metadata
  - Manually upload/annotate/delete data and metadata
  - Administration

# Components and Interfaces

# Metadata modelling

- Scientific administrator defines metadata schema for a scientific domain

- Contains tables and attributes

- Defines metadata properties:
  - Description
  - Data type
  - Unit
  - Provenance
  - Link to other attribute
  - …

# Metadata modelling

- Metadata exchanged in domain schema format
- Automatic query building using domain knowledge
- Pluggable database handlers for DMBS support
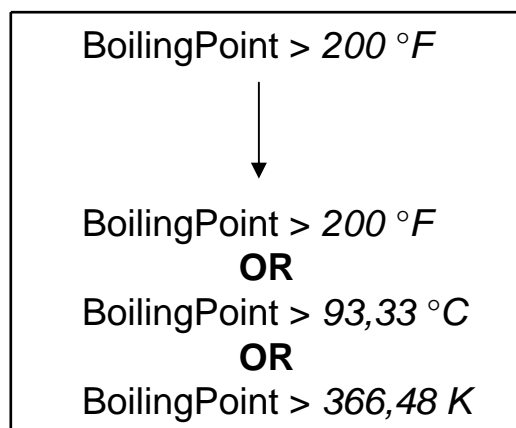
- GUI-based composition of new client views

Client (End-user Client, Workflow …)

Client API

DMS Access Service

Metadata Service

Domain knowledge

Database Handling

○ PostgreSQL

◉ MySQL

Metadata Database

# Querying data and metadata

- Seamless access to external data sources:
  SQL databases, web services,
  Excel files, web forms

→ Access to data and
  metadata regardless of
  source, e.g. in workflow
  system

# Querying data and metadata

- Automatic conversion of units in request and response
- Usage of external ontology services to broaden queries, e.g. synonyms from ChEbi

BoilingPoint > *200 °F*

↓

BoilingPoint > *200 °F*
**OR**
BoilingPoint > *93,33 °C*
**OR**
BoilingPoint > *366,48 K*

Substance = *'water'*

↓

Substance = *'water'*
**OR**
Substance = *'H2O'*
**OR**
Substance = *'aqua'*
**OR**
…

| Substance | BoilingPoint |
|-----------|--------------|
| water | 100 °C |
| arsenic | 1137,2 °F |
| helium | -268.93 °C |

# Storing files and metadata

Example: Workflow system stores result of QSAR workflow

1. Store file on UNICORE6 Storage → URL to file
2. Register file with location manager → logical name
3. Execute necessary unit conversions on metadata
4. Store metadata include logical name
5. Extract metadata from file (e.g. Structure Data Format, SDF)
6. Store extracted metadata

# Storing files and metadata

- Extract service:
    - Extraction logic in python scripts
    - Multiple extractors for single files possible
    - Uses metadata domain and file type to find matching extractors
    - Stores extracted metadata
    - e.g. create thumbnails from images, extract structure information from SDF file

# Storing files and metadata

# Security

- Uses UNICORE6 security infrastructure (X.509 certificates) to authenticate users

- XUUDB or Chemomentum VO management UVOS to authorise users

- Row-based access control lists for metadata and location information

- Metadata marked as provenance can only be modified/deleted by admin $\rightarrow$ provenance of calculation results

# Testbed installation

- Data Management System installed at TU Dresden
- Used by Workflow system to store workflow output and manage intermediate files

Chemomentum project
Testbed facilities locations

University of Tartu
Execution site

Forschungszentrum Jülich
Workflow site
Execution sites

ICM, Warsaw University
CA (pilot and demo)
Demo site
Testbed global registry
Testbed execution site

Dresden University
of Technology
Execution site
Data access
services

University of Zurich,
Execution site

IRFMN, Milan
TXT e-solutions, Milan
Execution sites

# Client

- Based on Eclipse Rich Client Platform
- Query, store, update and delete data and metadata
- Administrative functions, e.g. edit/create domain schemas
- GUI-based composition of new client views using domain knowledge, e.g. generation of query forms
- Extension points to build own interaction possibilities (e.g. integration of other views for data visualisation)

# Client: File upload

# Client: Search aquire

# Client: PDB and JMOL

# Thank you.