

MoSGrid – Molecular Simulations in a Distributed Environment

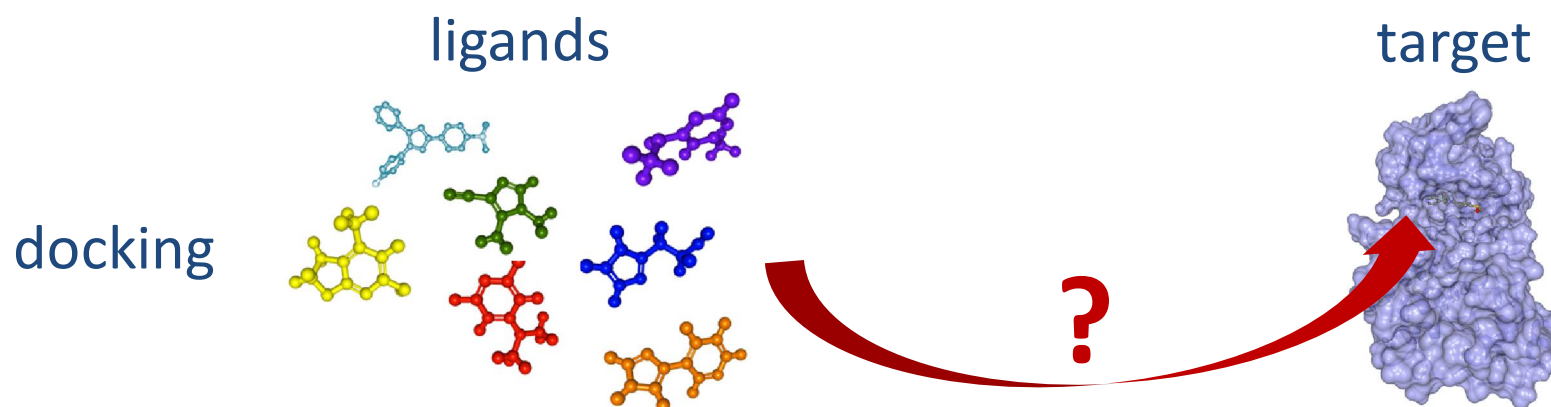
Sandra Gesing

sandra.gesing@uni-tuebingen.de

18 June 2013

Molecular Simulations and Docking

- Prediction and analysis of molecular structures
- Support by sophisticated tools and methods
- Numerous applications, e.g.
 - Materials science
 - Drug design



Molecular Simulations and Docking

- Prediction and analysis of molecular structure
- Support by sophisticated tools and methods
- Numerous applications, e.g.
 - Materials science
 - Drug design



- Data intensive and compute intensive
- Structure databases, e.g., ZINC with ~20 mio structures
- Sensitive and „expensive“ data
- Distributed data management available
- DCIs (Distributed Computing Infrastructures) available

Why do researchers not use the distributed environments on a large scale?

Open Issues

- Usability of tools often limited
- Complexity of methods
- Lack of graphical user interfaces

- Usability of tools often limited
- Complexity of methods

```
=====
| Version: 1.1
| build date: Jan 10 2012
| execution host: vomitoxin
| execution time: 2012-09-09, 16:39:43 (MST) |
=====

Available parameters are ('*' indicates mandatory parameters):
* -i <in.file>      input molecule file
* -o <out.file>     output file
  -ef <double>     error fraction; print error if fraction of invalid mols is larger
  -write_par <out.file> write xml parameter file for this tool
  -par <in.file>   read parameters from parameter-xml-file

Available flags are:
  -ri    remove invalid molecules.
  -ut    check for unique topologies
  -nc    no not check for unique conformations
  -rm    remove input file when finished
  -help  show help about parameters and flags of this program

This tool checks all molecules of the given input file for errors. Supported formats are mol2, sdf or drf (DockResultFile, xml-based).

The following checks are done for each molecule:

* bond-lengths may not be completely senseless (i.e. <0.7 or >2.5 Angstrom)
* each 'molecule' in the input file may only contain one actual molecule, i.e. there may be no unconnected atoms or fragments.
* each atom must have a valid assigned element
* the molecule must be protonated (since this is necessary for docking/(re-)scoring).
* 3D coordinates must be present (instead of 2D coordinates; also necessary for docking/(re-)scoring)
* partial charges may not contain completely senseless values (>5 or <-5).
* each conformation should appear only once within the given file, otherwise it is rejected and not written to the output file. However, if option '-ut' is used, molecules will instead be checked for unique topologies.

If option '-ri' is used, only those molecules that pass all those tests are written to the output file. If this option is not used, all molecules are written to output containing a property 'score_ligcheck' with a value of 1 if the molecule passed all tests or with a value of 0 if it did not pass them.

sshgw-bs[13] █
```

- Usability of tools often limited
- Complexity of methods
- Lack of graphical user interfaces
- **Workflows**

a sequence of connected steps in a defined order
based on their control and data dependencies

Open Issues

- Usability of tools often limited
- Complexity
- Lack of integration
- Workflow support

The image displays a collage of various bioinformatics web tools and interfaces, illustrating the complexity and lack of integration in the field. The tools shown include:

- AliBaba2.1**: A web-based tool for identifying complete gene structures in genomic DNA.
- NCBI BLAST**: The National Center for Biotechnology Information's Basic Local Alignment Search Tool.
- Signal Scan**: A web-based tool for predicting signal peptide positions in protein sequences.
- RepeatMasker**: A web-based tool for identifying and masking repetitive elements in DNA sequences.
- EMBOSS**: A collection of open-source bioinformatics software tools.
- InterProScan**: A web-based tool for protein domain and motif identification.
- EMBL-EBI**: The European Bioinformatics Institute's search and analysis tools.
- NCBI Signal Scan**: A web-based tool for predicting signal peptide positions in protein sequences.
- ABGENT**: A web-based tool for protein structure prediction.
- SUMOPLOT**: A web-based tool for protein structure prediction.
- National Center for Biotechnology Information**: A web-based tool for accessing and analyzing biological data.
- RepeatMasker Web Server**: A web-based tool for identifying and masking repetitive elements in DNA sequences.
- EMBOSS**: A collection of open-source bioinformatics software tools.
- InterProScan**: A web-based tool for protein domain and motif identification.
- EMBL-EBI**: The European Bioinformatics Institute's search and analysis tools.

Red arrows indicate workflow connections between these tools, showing how data from one tool is used as input for another, highlighting the complexity and lack of integration in the field.

Slide copied from: Stuart Owen „Workflows with Taverna“

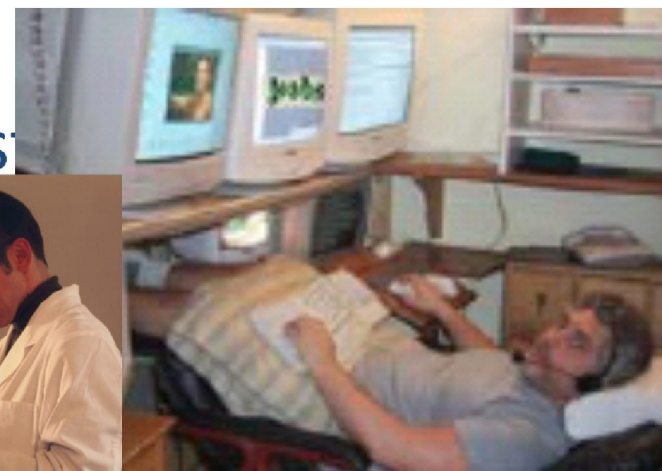
Open Issues

- Usability of tools often limited
- Complexity of methods
- Lack of graphical user interfaces
- Workflows
- **Complexity of infrastructures**
- **Users are generally not IT specialists**

Open Issues

- Usability of tools often limited
- Complexity of methods
- Lack of graphical user interfaces
- Workflows
- Complexity of infrastructures

not IT specialists



- Usability of tools often limited
- Complexity of methods
- Lack of graphical user interfaces
- Workflows
- Complexity of infrastructures
- Users are generally not IT specialists

⇒ User interfaces need to be intuitive and self-explanatory

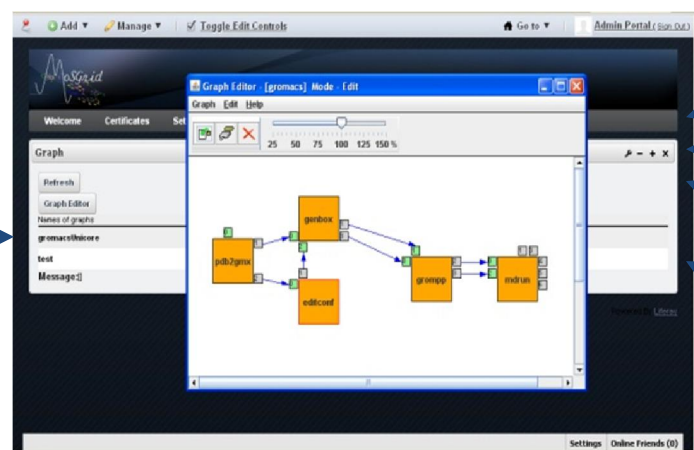
⇒ Science gateways

Science Gateways

“A Science Gateway is a community-developed set of tools, applications, and data that is integrated via a portal or a suite of applications, usually in a graphical user interface, that is further customized to meet the needs of a specific community.”

TeraGrid/XSEDE

Community



Gaussian

GROMACS FAST. FLEXIBLE. FREE.

XTREEMFS



Usability of software and data

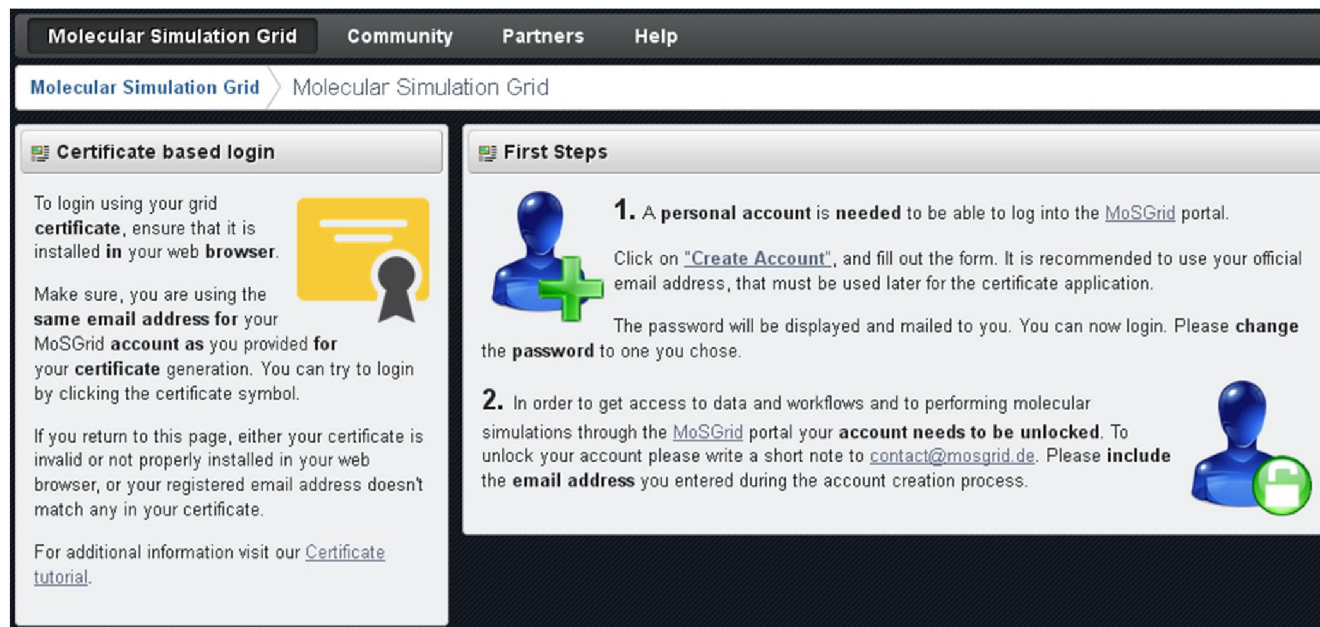
"After all, usability really just means that making sure that something works well: that a person ... can use the thing - whether it's a Web site, a fighter jet, or a revolving door - for its intended purpose without getting hopelessly frustrated."

(Steve Krug in *"Don't make me think!: A Common Sense Approach to Web Usability"*, 2005)



Molecular Simulation Grid

- Science gateway integrated with underlying compute and data management infrastructure
- Distributed workflow management
- Data repository
- Open source



The screenshot shows the MoSGrid portal interface. At the top, there is a navigation bar with "Molecular Simulation Grid", "Community", "Partners", and "Help". Below this, a breadcrumb trail shows "Molecular Simulation Grid" > "Molecular Simulation Grid". The main content area is divided into two columns. The left column is titled "Certificate based login" and contains instructions on how to use a grid certificate for login, including a note about email addresses and a link to a certificate tutorial. The right column is titled "First Steps" and contains two numbered steps: 1. A personal account is needed to log into the MoSGrid portal, with instructions on how to create an account and change the password. 2. In order to get access to data and workflows and to performing molecular simulations through the MoSGrid portal your account needs to be unlocked, with instructions on how to unlock the account. Both columns feature icons representing certificates and user accounts.

Molecular Simulation Grid

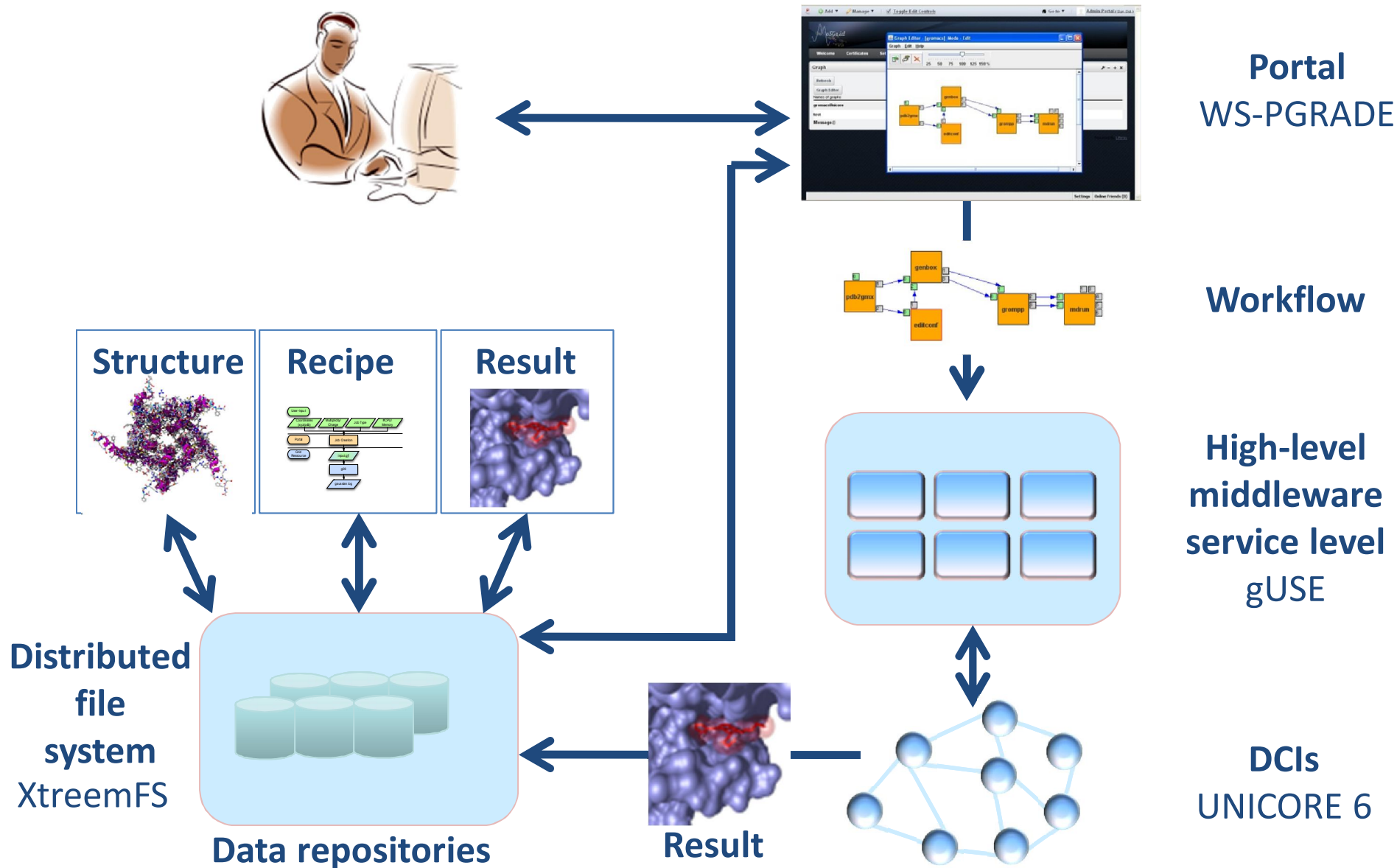
- Science gateway integrated with underlying compute and data management infrastructure
- Distributed workflow management
- Data repository
- Open source

Survey of willingness to share knowledge in the community

⇒ 90% share workflows

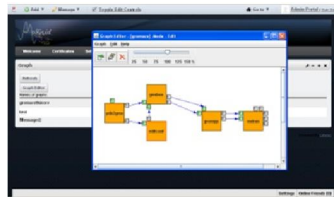
⇒ 70% share results after publication

MoSGrid in a Nutshell



Job and Workflow Management

grid User Support Environment



User Interface
WS-PGRADE

**Workflow
storage**

**Application
repository**

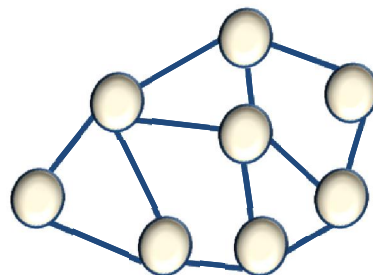
**Information
system**

**Workflow
engine**

Submitters

Logging

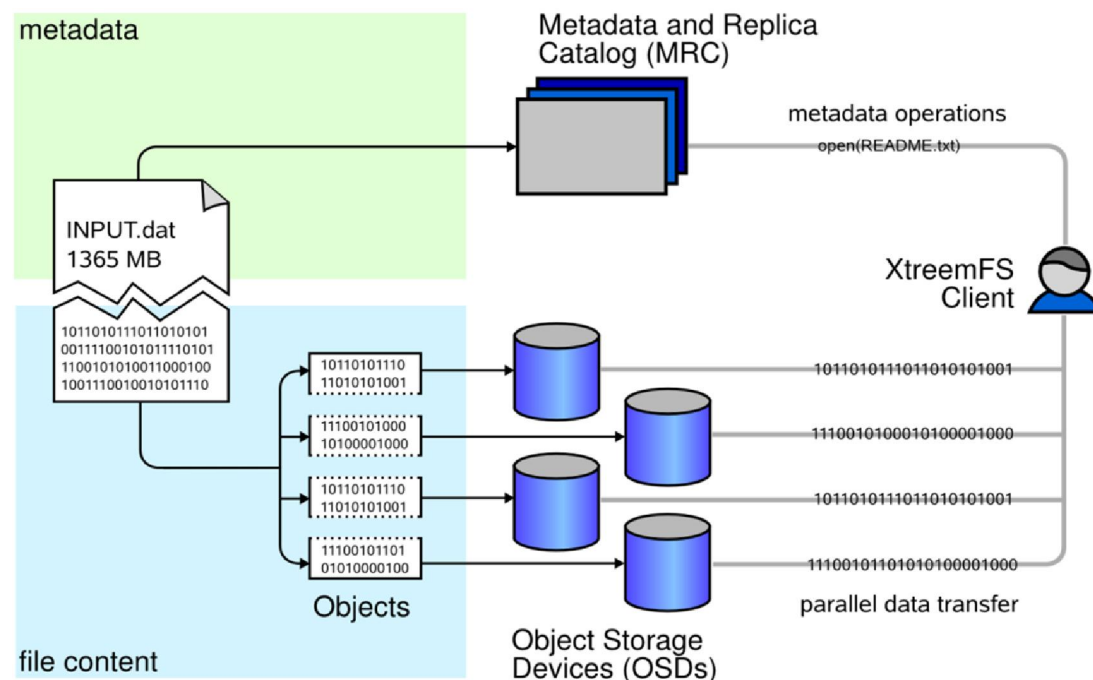
**High-Level
Middleware
Service Layer**
gUSE



**DCI Resources
Middleware Layer**
UNICORE 6

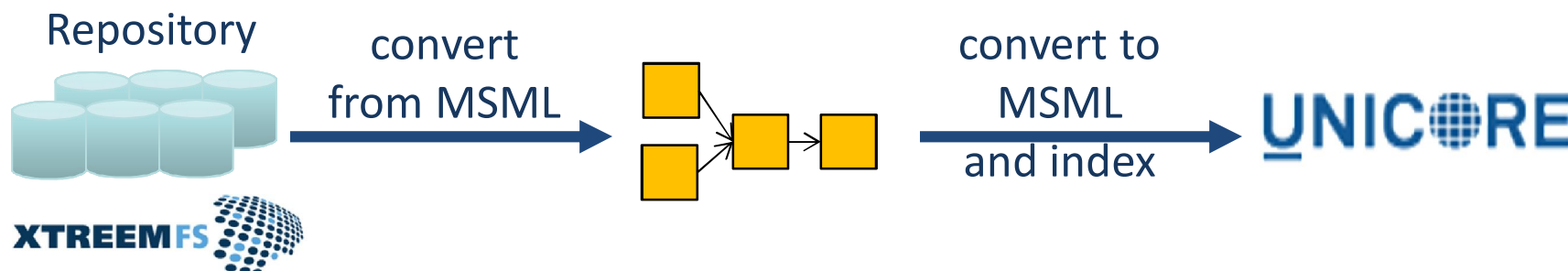
Distributed Data Management

- XtremFS is an object-based grid and cloud filesystem
- Replication for availability, locality, bandwidth, and latency
- Easy integration in heterogenous environments
- UNICORE extension with URL scheme `xtreemfs://` available



- Repository consists of data and metadata storage
- MSML (Molecular Simulation Mark-up Language)
- Subset and extension of CML (Chemical Mark-up Language)
- Unified data representation
- Used for storing structures of molecules and macromolecules, simulation descriptions, and results
- Parsers and adapters used for conversions to and from MSML

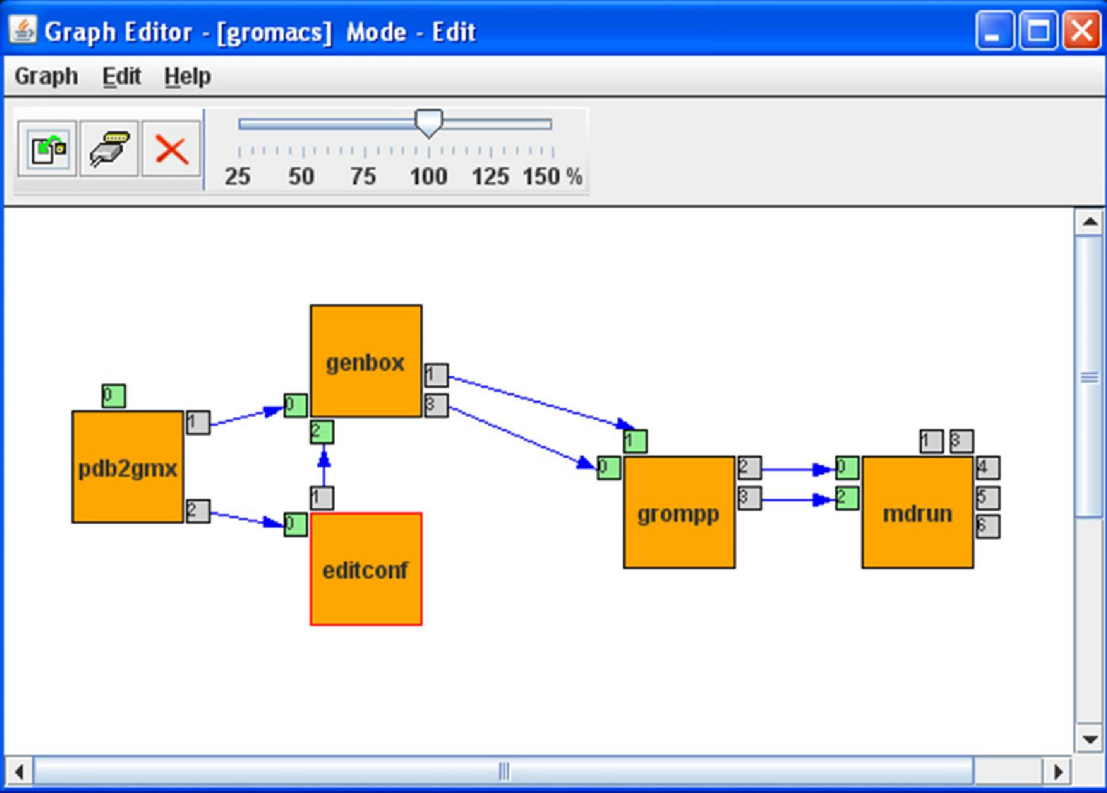
- First step in workflows: MSML to application specific input
- Results are computed
- Output converted to MSML
- Last step in workflows: MSML to JSON metadata (metadata extractor in UNICORE)
- JSON indexed for searchable results by UNICORE and LUCENE (text search engine API)



Graph Editor - [gromacs] Mode - Edit

Graph Edit Help

25 50 75 100 125 150 %



Refresh

Graph Editor

Names of graphs

gromacsUnicore

test

Message:[]

Powered By [Liferay](#)

Settings Online Friends (0)

Job Configuration

Job's name: ParserProtein

Optional note: Description of Job

[Job Executable] [Job I/O] [JDL/RSL] [History]

WorkflowService Binary ?

Type: unicare

Grid: flavus.informatik.uni-tuebingen.de:8090

Tools: Bash shell 3.1.16

Execute parser:

Replicate settings in all Jobs:

Copy job names to tools: ?

Kind of binary: Sequential Java MPI

MPI Node Number:

Executable code of binary: Recently stored:

Parameter: genparser.sh ProteinProc

Job Configuration

Job's name: PDBCutter

Optional note: Description of Job

[Job Executable] [Job I/O] [JDL/RSL] [History]

WorkflowService Binary

Type: unicore

Grid: flavus.informatik.uni-tuebingen.de:8090

Tools: PDBCutter 1.0.0

Execute parser:

Replicate settings in all Jobs:

Copy job names to tools:

Kind of binary:

MPI Node Number:

Executable code of binary:

Parameter:

- ModelCreator 1.0.0
- MolCombine 1.0.0
- MolDepict 1.0.0
- MolFilter 1.0.0
- MolPredictor 1.0.0
- nwchem 6.1
- obabel (OpenBabel) 2.3.1
- PartialChargesCopy 1.0.0
- pdb2gmx 4.5.5
- PDBCutter 1.0.0**
- PDBDownload 1.0.0
- Perl 5.8.8
- PocketDetector 1.0.0
- POVRay 3.5
- Predictor 1.0.0
- PropertyModifier 1.0.0
- PropertyPlotter 1.0.0
- ProteinCheck 1.0.0
- ProteinProtonator 1.0.0
- Python_Script 2.4.2

Remote File Configuration

2012-11-21

Job's name: ParserProtein

Optional note: Description of Job

[Job Executable] [Job I/O] [JDL/RSL] [History]

Port Number:0 Port Name: genparser Description of Port

Input Port's Internal File Name: genparser.jar

Port dependent condition allowing the run of the job: View Hide

Source of input directed to this port:
xtreemfs://test/genparser.jar
 Copy to WN:

Parametric Input details: View Hide

Port Number:1 Port Name: startscript Description of Port

Molecular Dynamics

- Study and simulation of molecular motion

Quantum Chemistry

- Study and simulation of molecular electronic behavior relative to their chemical reactivity

Docking

- Main focus on evaluation of ligand-receptor interactions (e.g., for drug design)

Molecular Dynamics Simulation Portlet

Welcome Workflow Submission Workflow

Get workflow status

Workflows list


	NAME
●	birke/2011-14-1/11:59:42/EQ_
●	birke/2011-14-1/12:15:31/SIN
●	birke/2011-14-1/12:19:28/EQ_
●	birke/2011-14-1/12:52:56/EQ_
●	jk/2011-14-1/13:41:38/EQ_Gr

<https://unicore6-bisgrid.uni-paderborn.de>

list files Abort workflow Des

Files to get:

Jmol



Docking Portlet

Docking Portlet

Import Standard

Select an input file

Import Standard

Please fill in the following fields

PDB

Filename: 1DX6.pdb

PDB Model: Model

Chain A Chain B

Chain name: A

Name: GNT

Protein: []

Welcome Monitoring Debug

Workflows:

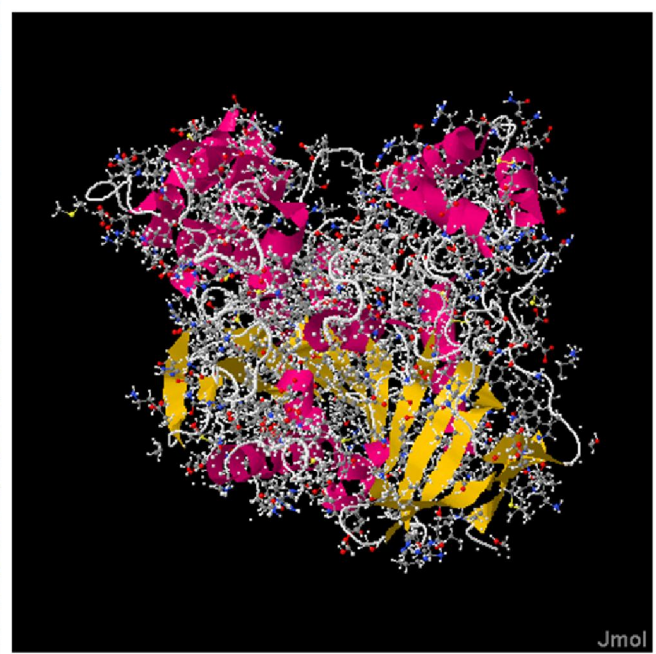
- TEST20111025
- TEST20111025
- 1EVE
 - Docking
 - recH.pdb
 - results.sorted.sdf

1EVE/ STATUS: FINISHED / selected			
ATOM	1	N	SER A
ATOM	2	CA	SER A
ATOM	3	C	SER A
ATOM	4	O	SER A
ATOM	5	CB	SER A
ATOM	6	OG	SER A
ATOM	7	HN1	SER A
ATOM	8	HN2	SER A
ATOM	9	HN3	SER A
ATOM	10	HA	SER A
ATOM	11	HB1	SER A
ATOM	12	HB2	SER A
ATOM	13	HG	SER A
ATOM	14	N	GLU A
ATOM	15	CA	GLU A
ATOM	16	C	GLU A
ATOM	17	O	GLU A
ATOM	18	CB	GLU A
ATOM	19	CG	GLU A
ATOM	20	CD	GLU A
ATOM	21	OE1	GLU A
ATOM	22	OE2	GLU A
ATOM	23	HN	GLU A
ATOM	24	HA	GLU A
ATOM	25	HB1	GLU A
ATOM	26	HB2	GLU A
ATOM	27	HG1	GLU A
ATOM	28	HG2	GLU A

Update Delete Download View in Jmol

Jmol

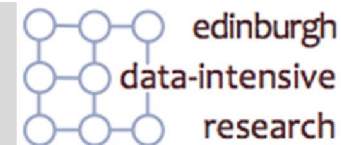
132040062594recH.pdb



Jmol

-7.010 86.637 39.078 1.00 0.00 H

The MoSGrid Science Gateway

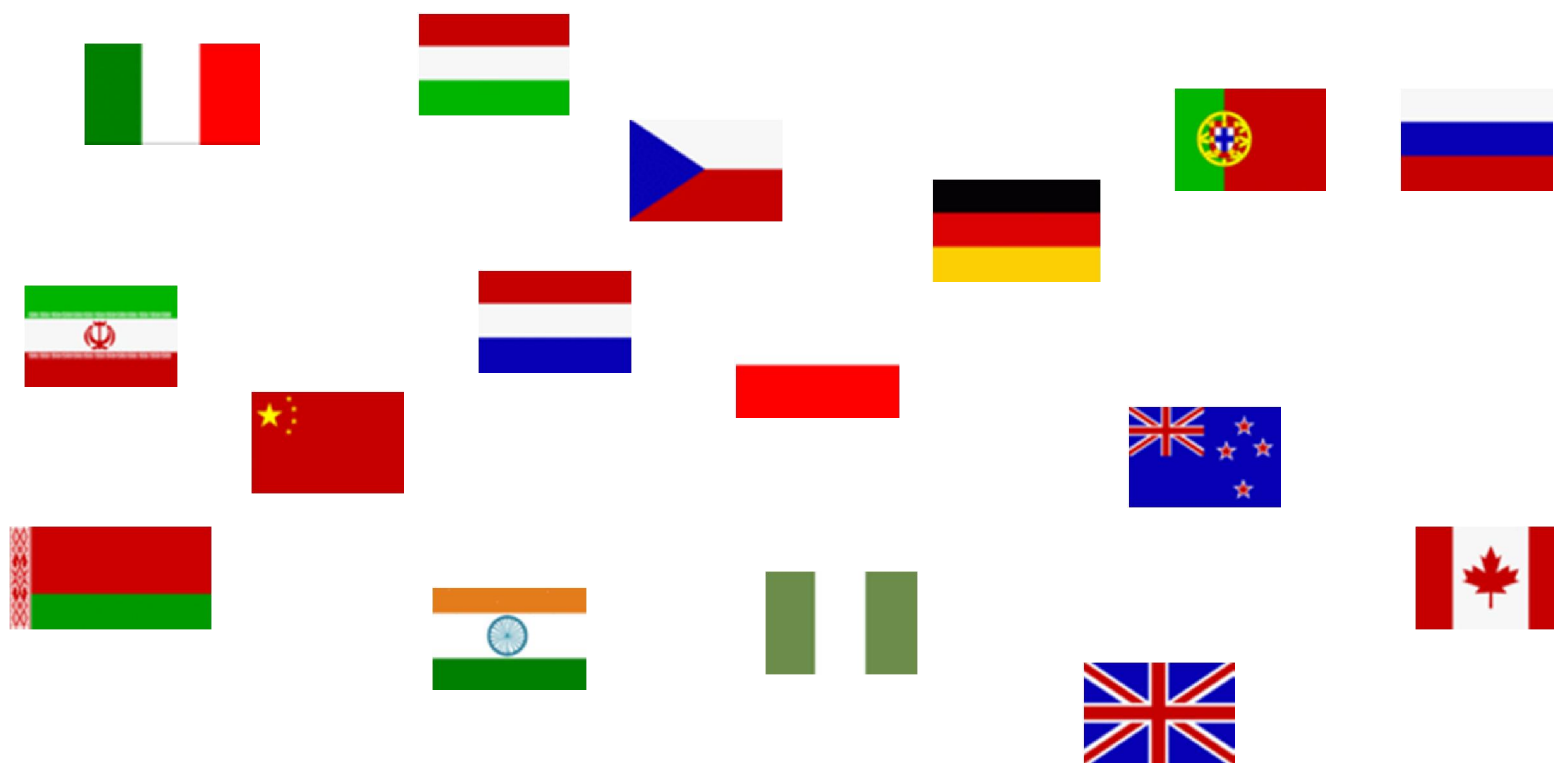


Built by 17 institutes and companies



The MoSGrid Science Gateway

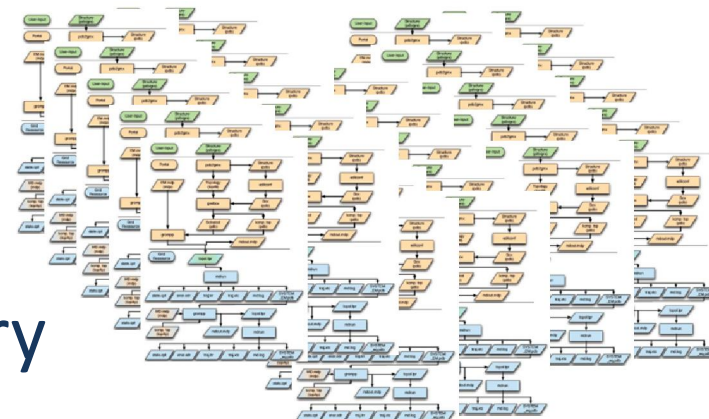
Used by 125 user groups in 16 countries



The MoSGrid Science Gateway

Currently supports

- 65 workflows in repositories
- 90 applications
- 85 GB data in central repository
- ~50,000 structures (1 – 5 GB) in user repositories



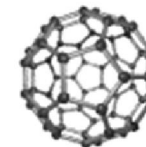
AutoDock Vina

GROMACS FAST.
FLEXIBLE.
FREE.

FlexX

Gaussian

CADDSuite



MoSGrid ended 31.12.2012 but partners participate in

SCI-BUS (SCientific gateway Based User Support)

- EU project 01.10.2011 – 30.09.2014
- Extension of the MoSGrid portal with an interactive molecule editor based on WebGL and a semantic search

ER-flow (Building an European Research Community through Interoperable Workflows and Data)

- EU project 01.10.2012 – 30.09.2014
- Integration of applications in SHIWA simulation platform
- Study of data exchange between workflow systems
- Community management



sandra.gesing@uni-tuebingen.de