

From the Desktop to the Grid: Conversion of KNIME Workflows to gUSE

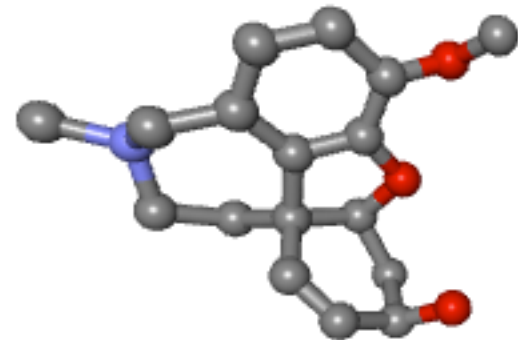
Luis de la Garza, Jens Krüger, Charlotta Schärfe, Marc Röttig,
Stephan Aiche, Knut Reinert, Oliver Kohlbacher

*Department of Applied Bioinformatics
University of Tübingen*

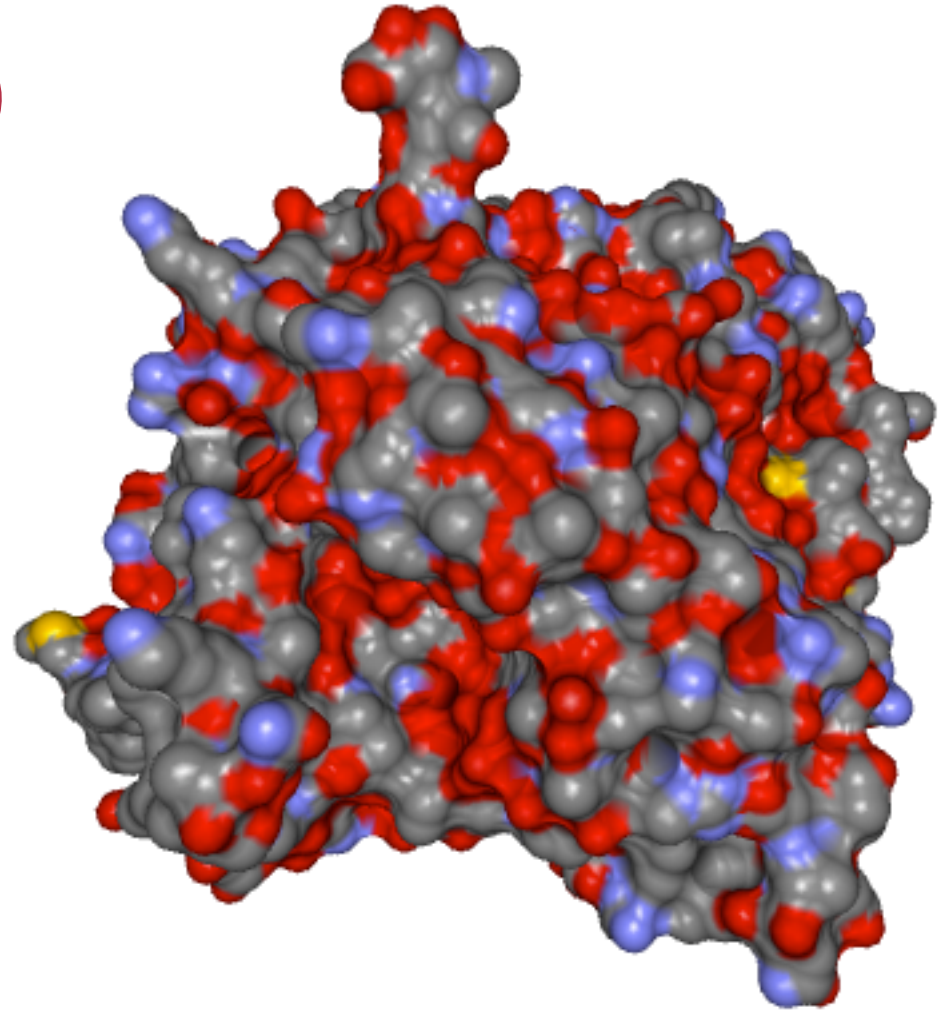
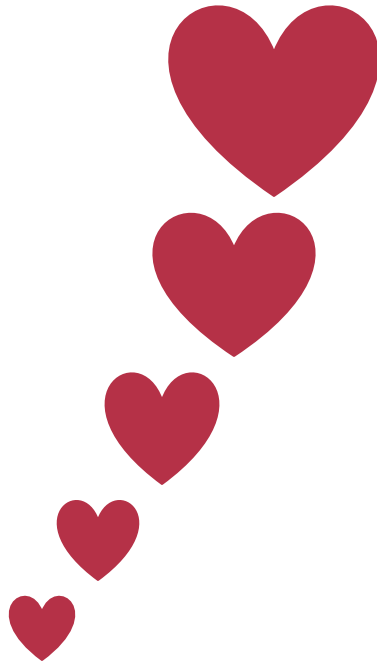
EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN



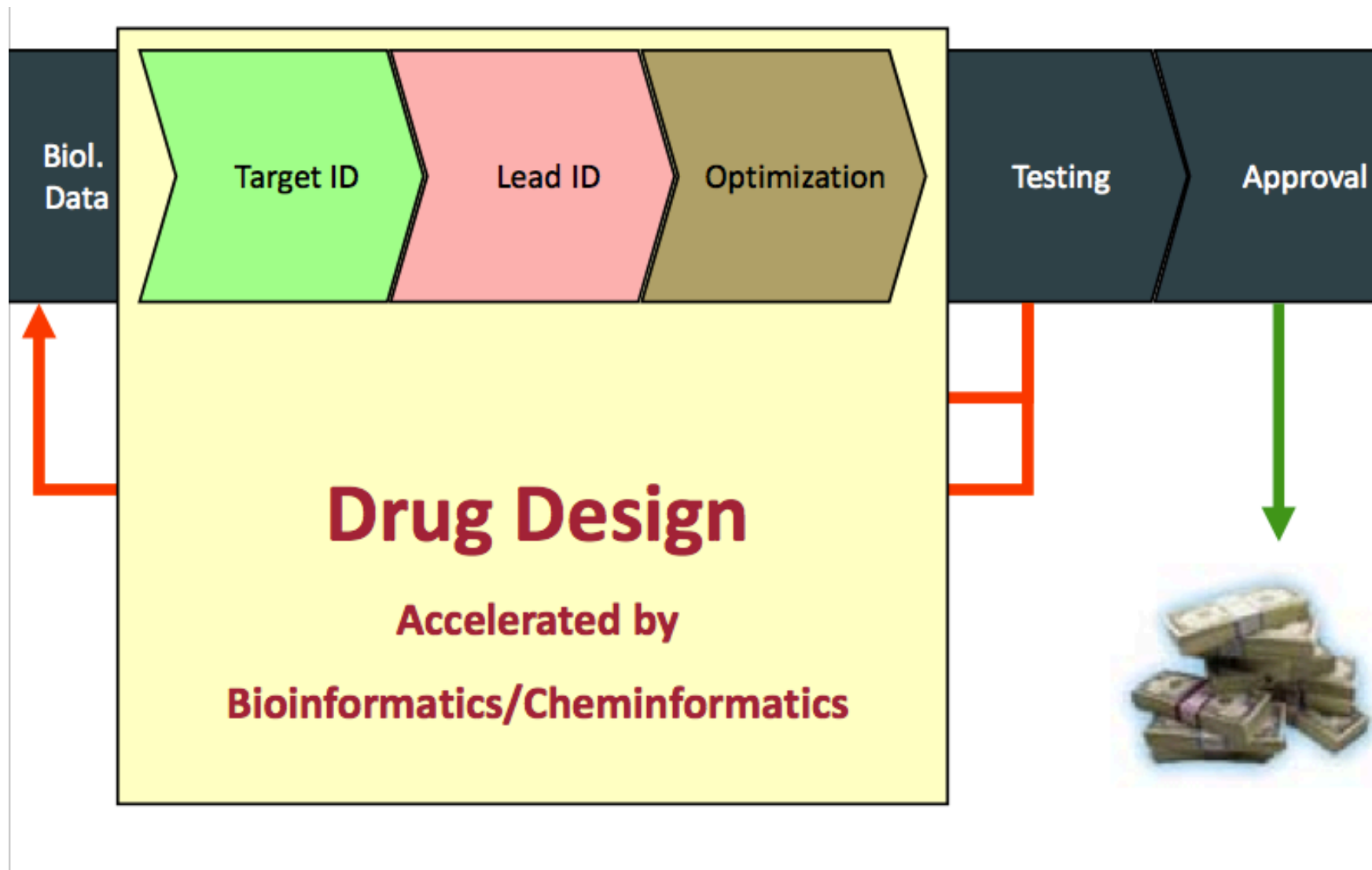
What are we researching?



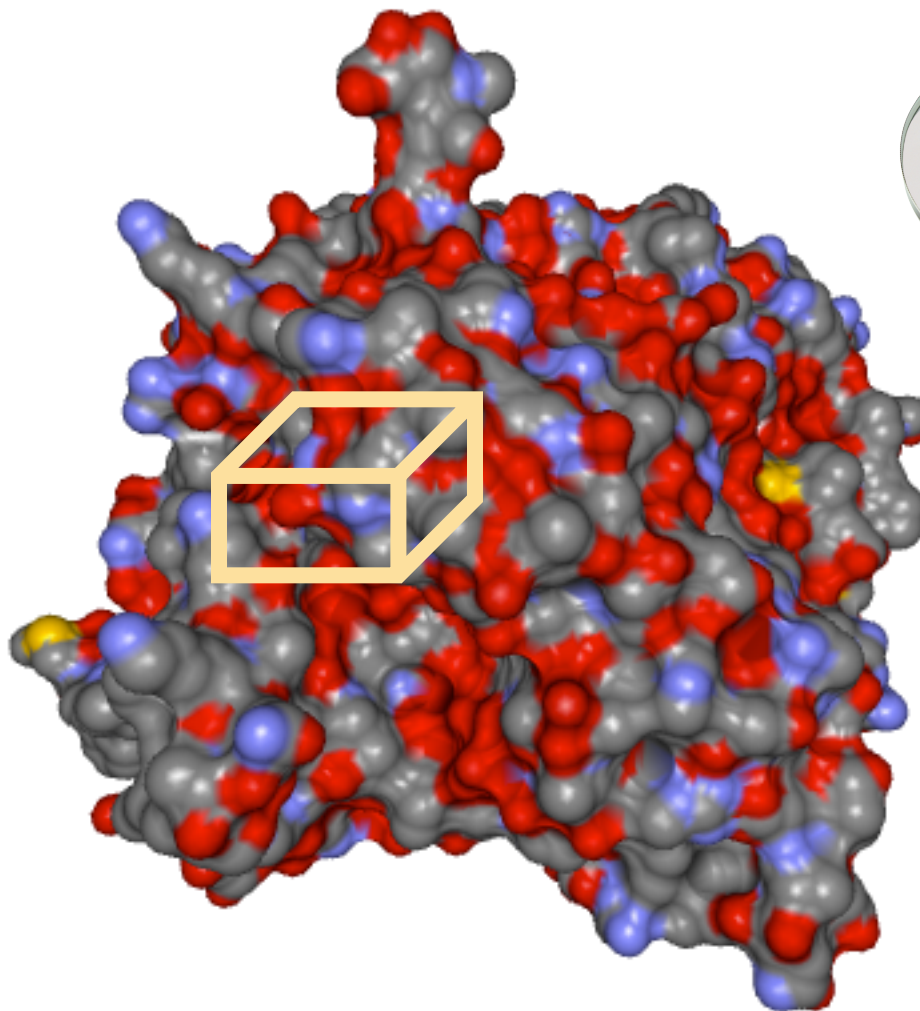
Hi there,
pretty lady



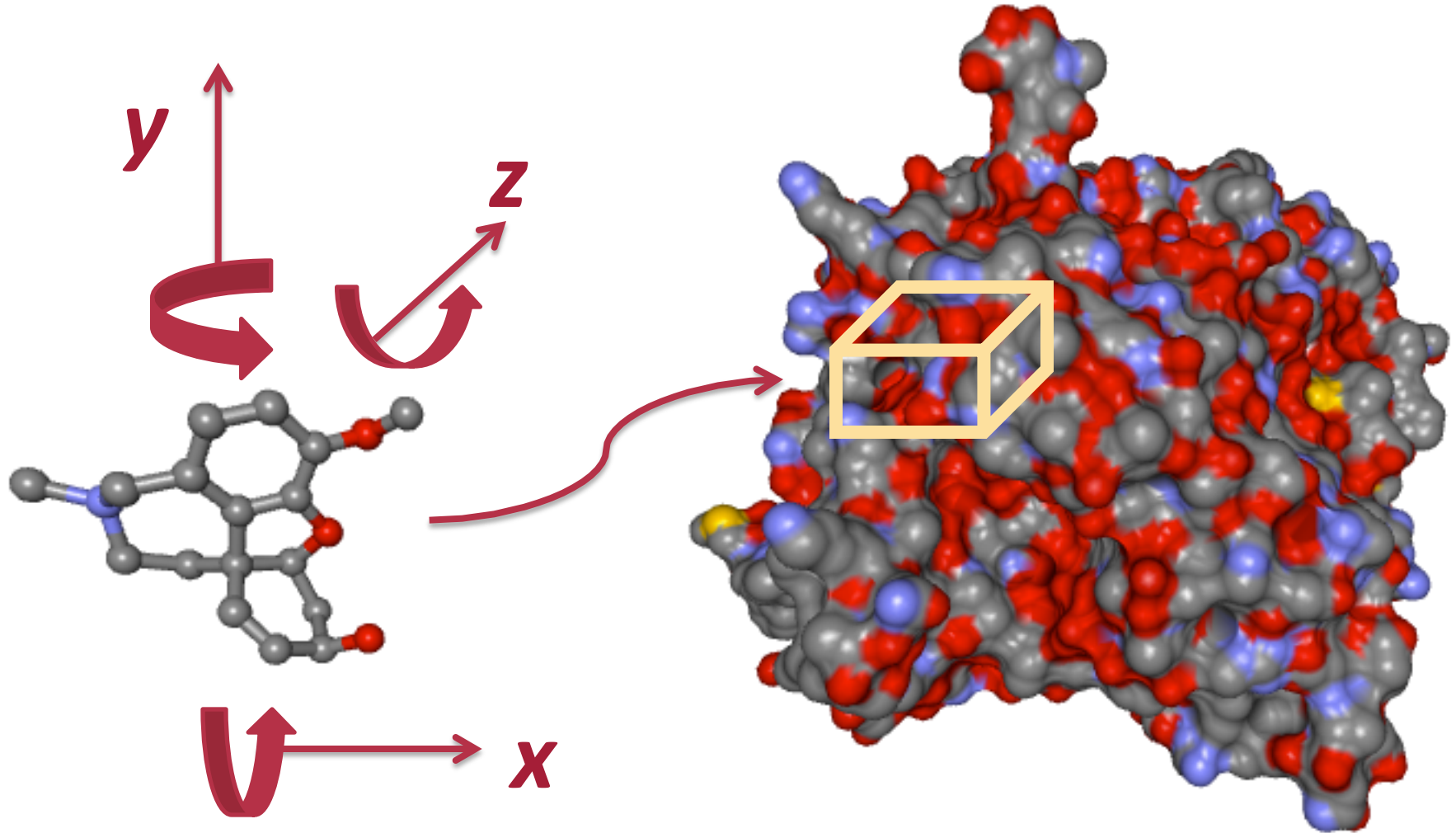
Why is Docking important?



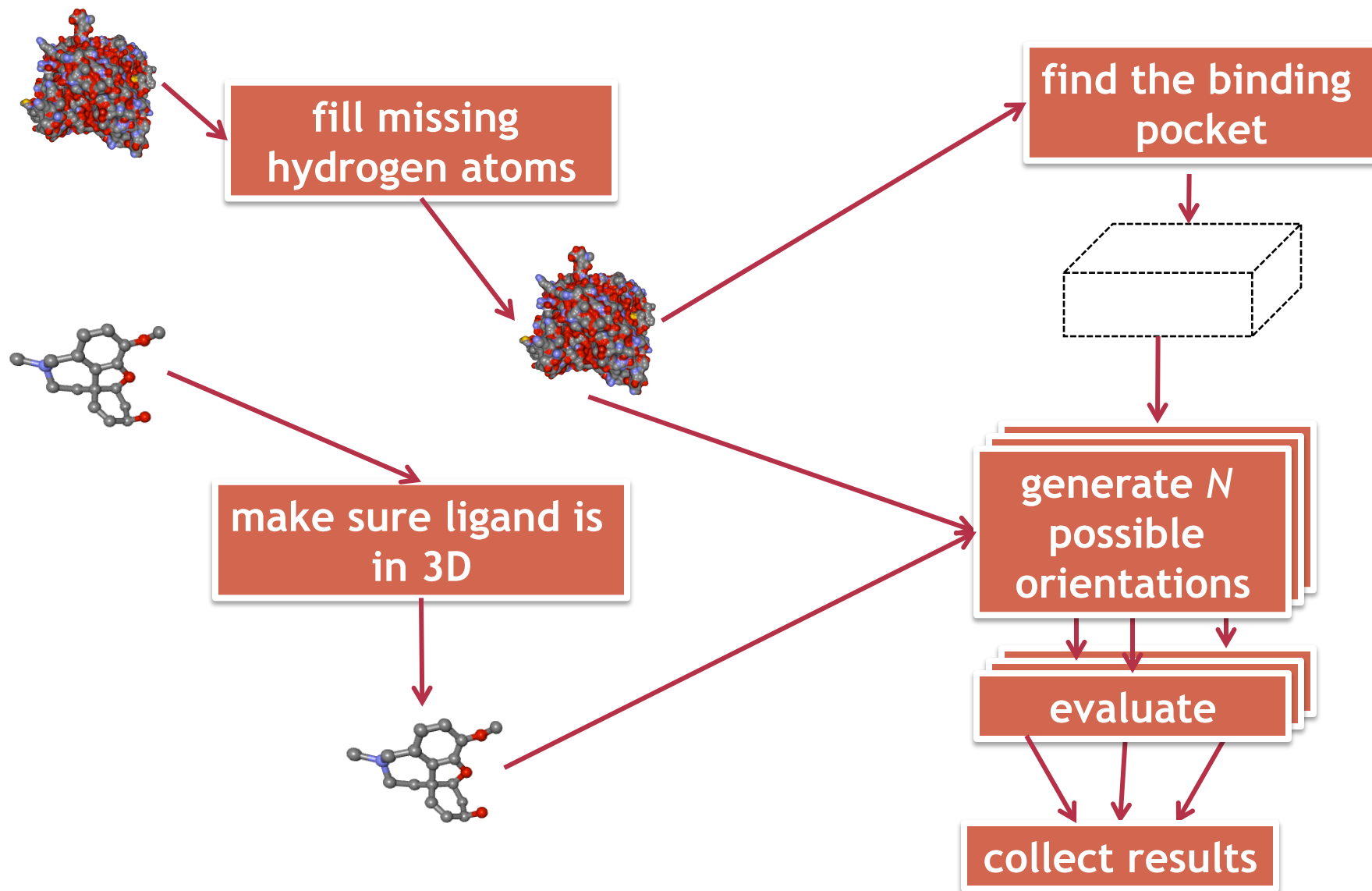
A Docking Recipe



Evaluation Orientations



It's begging to be a Workflow



Why do we use Workflows?



- Docking can be broken down as a series of small tasks; some of them can be executed in parallel
- We need access to resources offering High Performance Computing (HPC)
- We want to store intermediate results for further analysis (*i.e.*, binding pocket computation)
- We need flexibility – we would like to test different docking algorithms without making lots of changes

- Workflows have applications in other areas of bioinformatics and in other fields, such as:
 - Data mining
 - Business process automation
 - Customer relationship management
 - Business intelligence
- If you can split a process in small automated tasks, then it can be put into a workflow

- Several needs have produced different workflow managers



Taverna



Konstanz Information Miner



The screenshot shows the KNIME Desktop application window. The title bar includes the Apple logo, the text 'KNIME', and system status icons. The menu bar shows 'KNIME'. The main interface is divided into several panes:

- KNIME Explorer:** Contains 'Personal favorite nodes', 'Most frequently used nodes', and 'Last used nodes'.
- Node Repository:** A list of node categories including IO, Database, Data Manipulation, Data Views, Statistics, Mining, Chemistry, ChemAxon / Infocom, Distance Matrix, Meta, Flow Control, Misc, KNIME Labs, Time Series, Quick Form, R, Reporting, Testing, Weka, and XML.
- New Wizard Dialog:** A modal window titled 'New' with the subtitle 'Select a wizard'. It features a search bar with the placeholder text 'type filter text'. Below the search bar, a list of wizards is displayed, including 'New KNIME Workflow' and 'New KNIME Workflow Group'. A checkbox labeled 'Show All Wizards.' is at the bottom left. Navigation buttons '< Back', 'Next >', 'Cancel', and 'Finish' are at the bottom right.
- Node Description:** A pane on the right showing details for the 'Interactive Table' node. It includes a title 'Interactive Table', a description, and sections for 'Ports' and 'Views'.
- Console:** A pane at the bottom showing 'No operations to display at this time.'

Interactive Table

Displays data in a table view. If the number of rows is unknown, the view counts the number of rows when opened. Furthermore, rows can be selected and highlighted.

Ports

Input Ports

0 Input table to display.

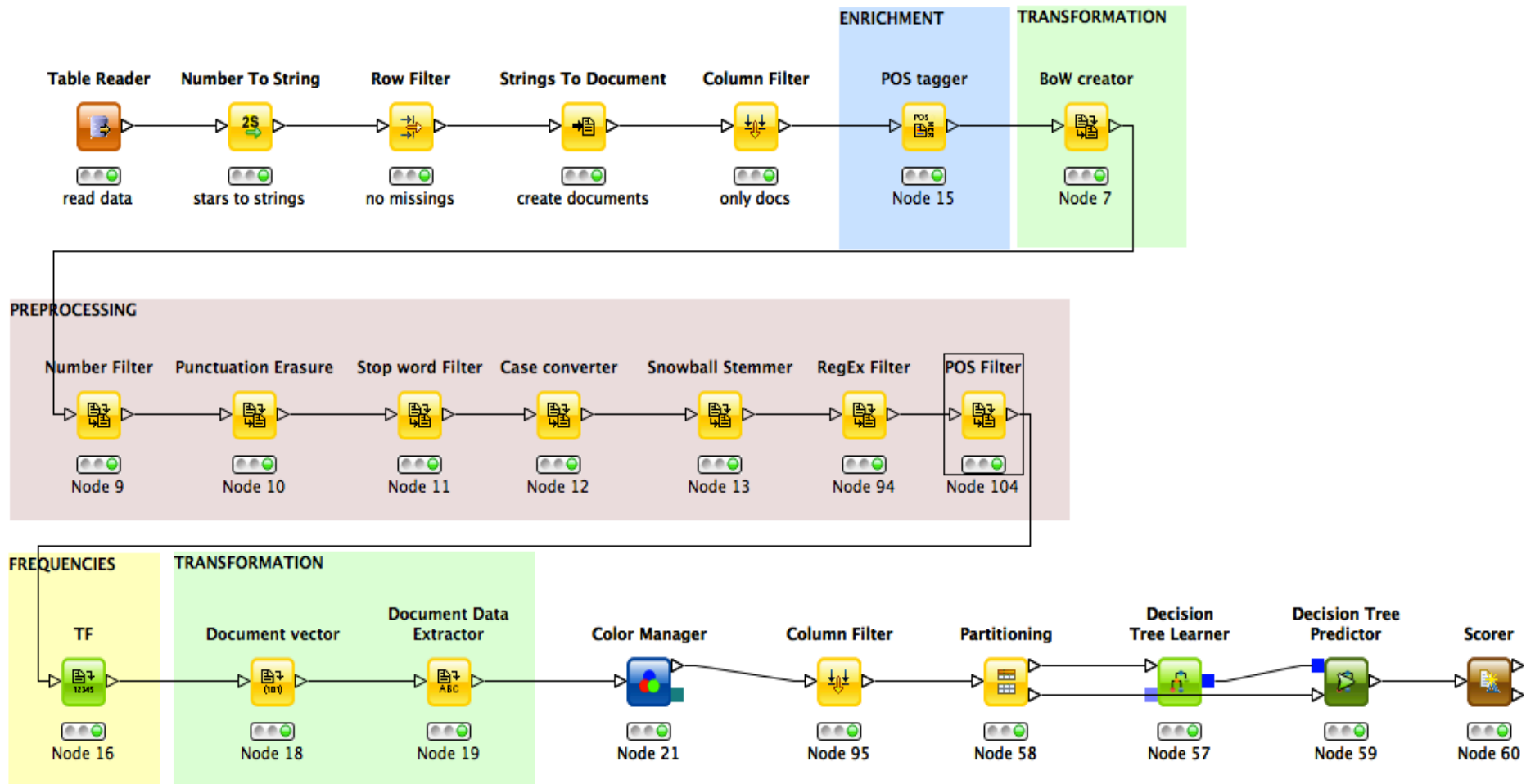
Views

Table View

Displays the data in a table view.

This node is contained in KNIME Base Nodes provided by KNIME GmbH, Konstanz, Germany.

A more elaborated Workflow



Workflow created during training at the KNIME User Meeting 2013, Zurich



- Create and execute workflows on your desktop computer; free, as in “*free beer*”, open source, available under the GPL license
- Easy-to-use interface (if you know how to *click*, you know how to KNIME)
- Already one of the most commonly used workflow management systems in the field of e-Science systems



- Lots of KNIME Community Nodes available (*e.g.* R, NGS, Schrödinger, etc.)

SCHRÖDINGER.



- Users can easily extend KNIME by creating new nodes (Java API) and extensions (based on Eclipse)





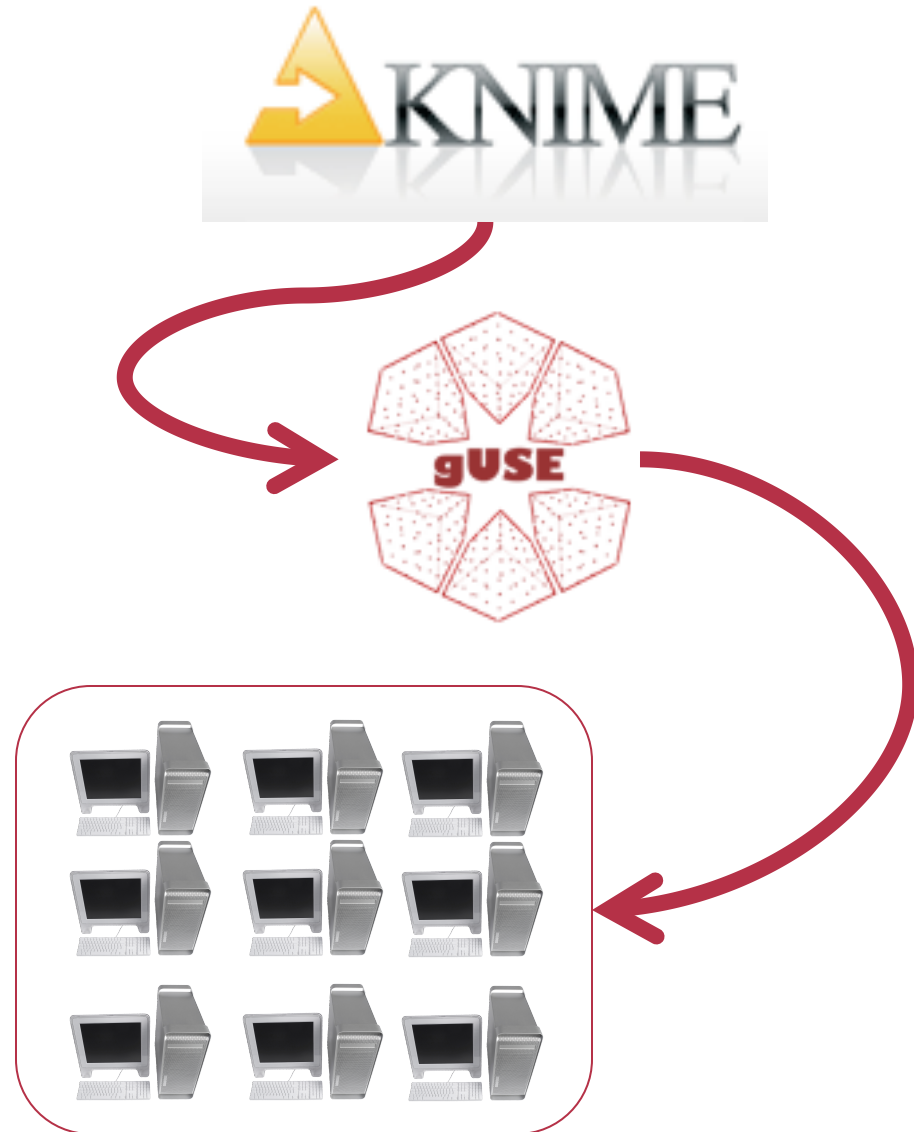
- Offers the capability to run KNIME Workflows on a cluster
- Works together with KNIME Server and KNIME Desktop to offer all features
- Not free, as in “*VISA or MasterCard?*”
- Needs to be installed as a cluster resource



Choices, choices...



- We want to use KNIME as a workflow editor
- We want access to computing resources
- We want gUSE to execute our workflows



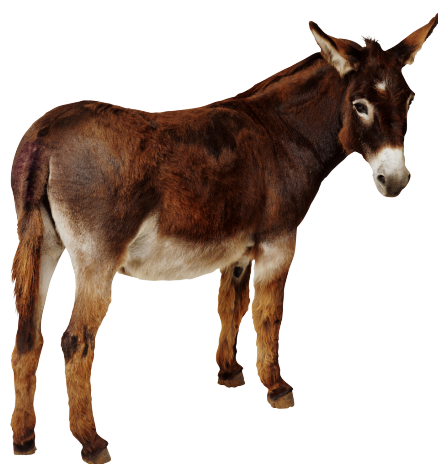


- Wouldn't it be great if arbitrary external command line tools could be integrated into KNIME?





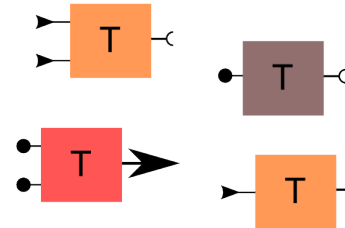
- KNIME is based on data tables, while most command line tools are based on files
- It would be possible, yet time consuming and error prone, to manually generate a KNIME node for each tool to be integrated



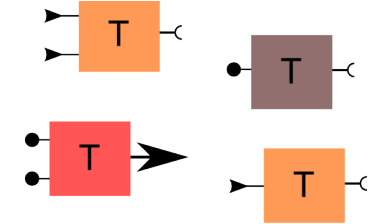


- We have developed several tool suites, each containing several tools
- Each tool should be usable as a single job in a workflow

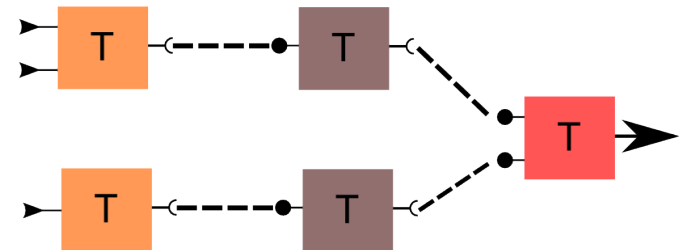
Tool Suite X



Tool Suite Y

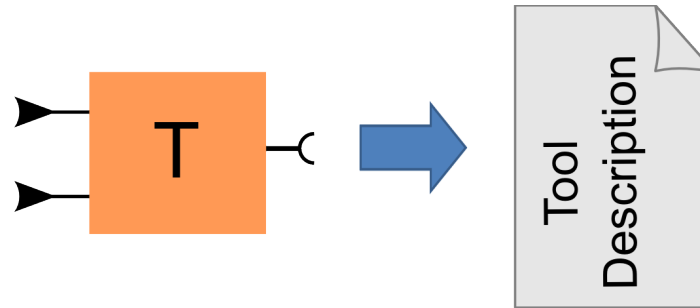


Workflow





- Give every tool a self-describing output format: semantic annotation of inputs, outputs and parameters



- Common Tool Description (CTD) was initially developed by the Open-Source Framework for Mass Spectrometry (OpenMS) team



- Each tool can “tell” its requirements and options, thus allowing easy integration and a more robust approach than “classical” tool stubs (*e.g.*, scripting)
- XML Documents – easy to understand, parse and generate



Common Tool Description



```
<!-- pdbcutter_execution.xml -->
<tool status="internal">
  <name>PDBCutter</name>
  <PARAMETERS version="1.3">
    <NODE>
      <ITEM name="i" tags="input file" value="input.pdb"/>
      <ITEM name="rec" tags="output file" value="receptor.pdb"/>
      <ITEM name="lig" tags="output file" value="ligand.pdb"/>
      <ITEM name="lig_chain" tags="required" value="A"/>
      <ITEM name="lig_name" tags="required" value="GNT"/>
    </NODE>
  </PARAMETERS>
</tool>
```

```
# using PDBCutter with a CTD as input
$ PDBCutter -par pdbcutter_execution.xml
```

```
# classic invocation
$ PDBCutter -i input.pdb -rec receptor.pdb -lig ligand.pdb
               -lig_name GNT -lig_chain A
```

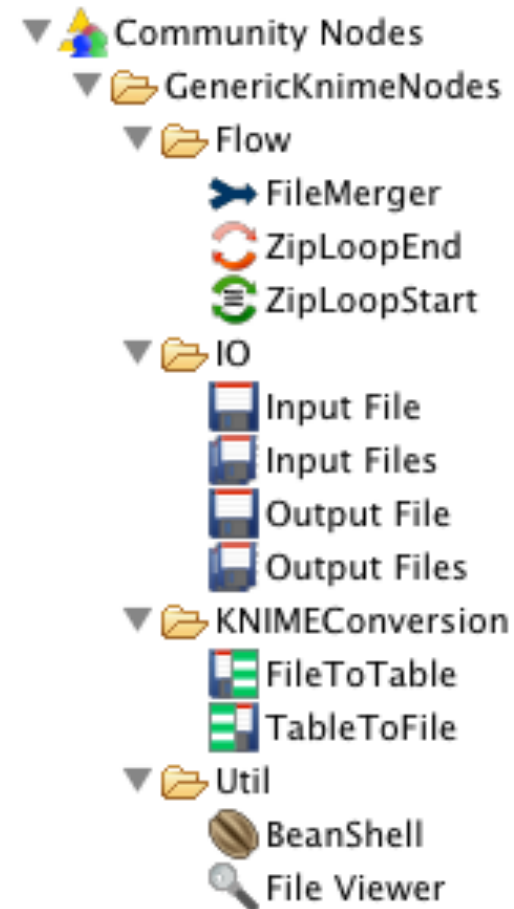
- Leaning on CTDs, external tools can be integrated into KNIME
- Generic KNIME Nodes (GKN) are 100% compatible with other KNIME nodes
- Any tool described by a CTD can be used in KNIME



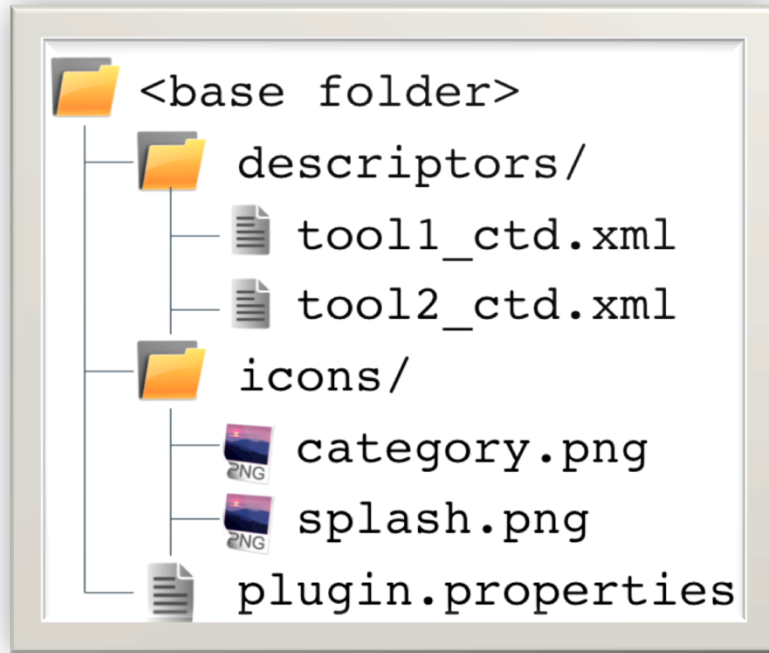


- Remember that KNIME relies on data tables and most command line tools rely on files?

GKN bridges this gap



Generic KNIME Nodes



new_nodes.jar

install as KNIME
plugin

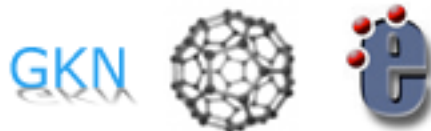


Generic KNIME Nodes



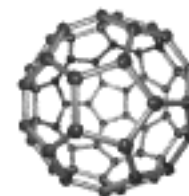
Version 2.6.0
(Build July 31, 2012)

Installed Extensions:



Loading Workbench

Copyright, 2003 - 2012, KNIME GmbH, Germany, <http://www.knime.org/>, contact@knime.org



Computer Aided Drug Design
Suite (CADD Suite)



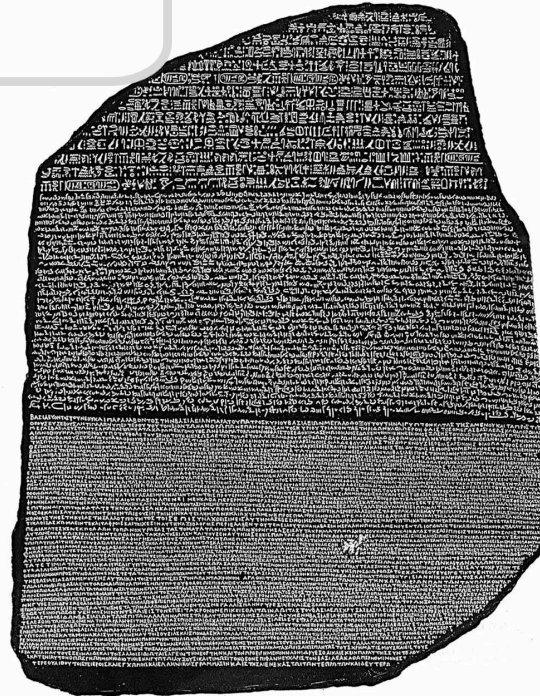
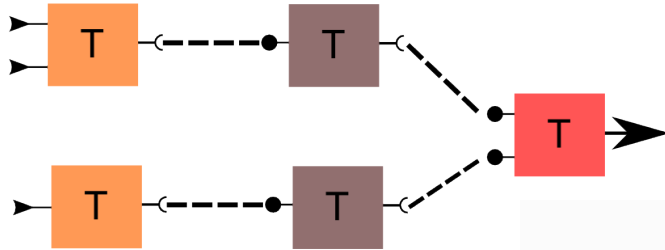
The European Molecular
Biology Open Software Suite
(EMBOSS)

CADD Suite and EMBOSS running as extensions in KNIME

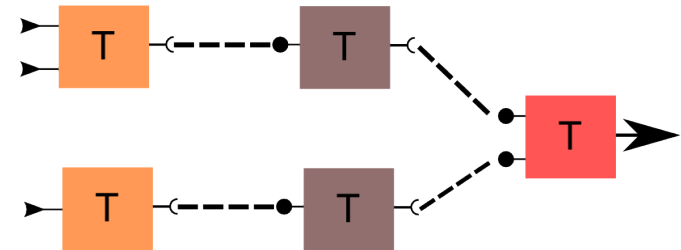
Rosetta Stone still missing



gUSE Workflow



KNIME Workflow



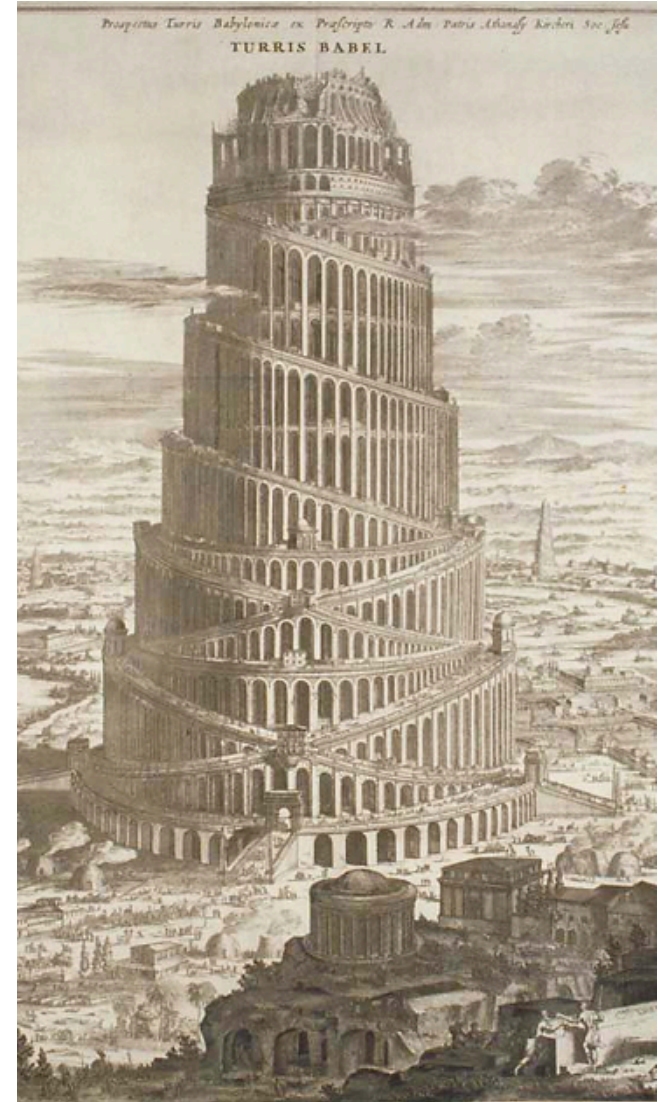
- In general, workflows are not compatible across different workflow managers
- It is not trivial to convert workflows between workflow managers



What's the Problem?



- Topology has to be converted
- Jobs have their own configuration settings and must be converted too
- Parameter sweep implementation might vary across managers
- Architecture on which jobs are executed might differ



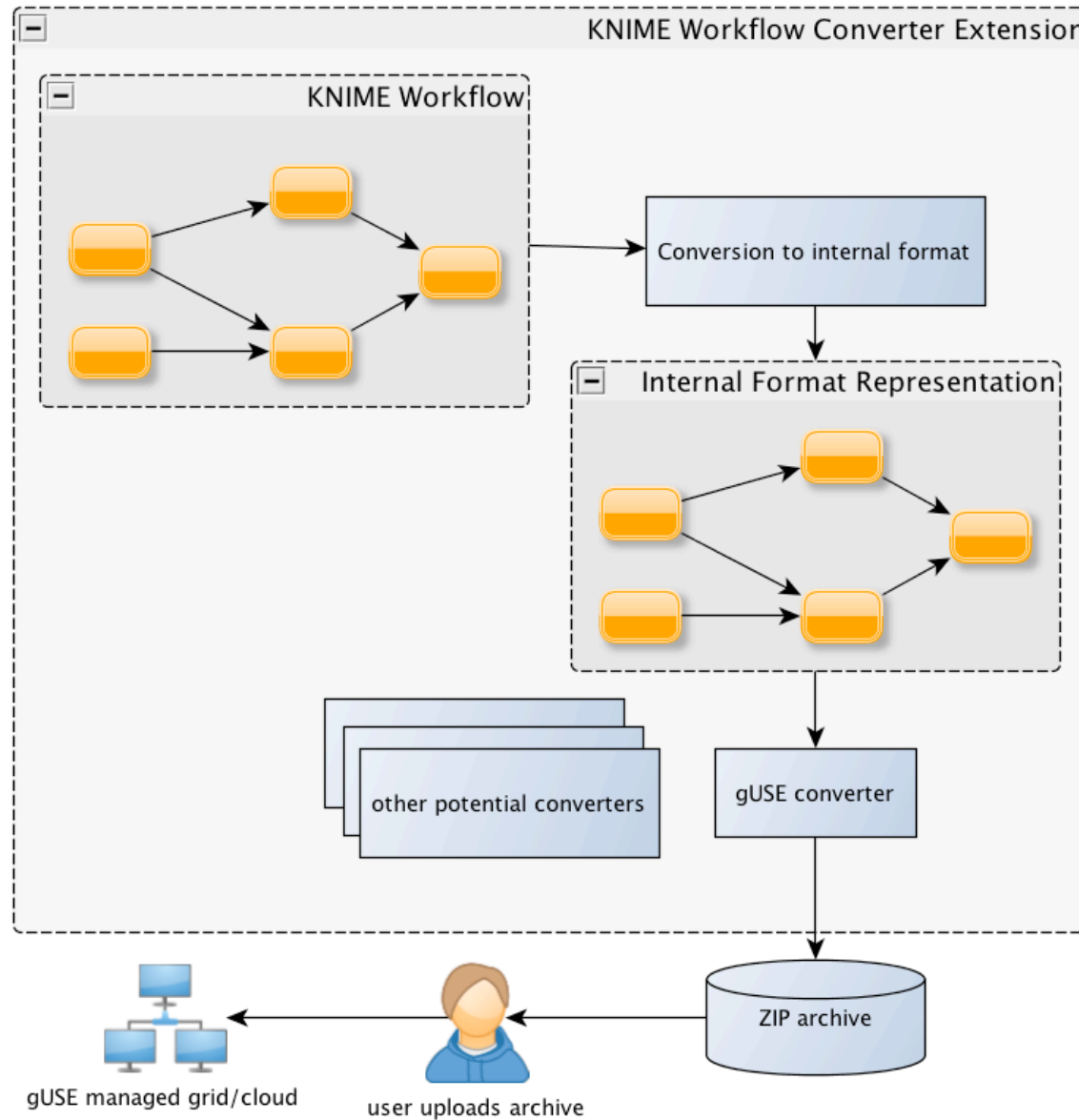
- Sharing Interoperable Workflows for large-scale scientific Simulations on available Distributed Computing Interfaces (SHIWA)
- Allows workflows to be reused across different workflow languages





- Creating and editing workflows in KNIME has a gentle learning curve
- However, desktop computers don't offer too much computing power for our needs
- **If only it were possible to upload a KNIME workflow to gUSE...**

From KNIME to gUSE





- KNIME has its own development cycles
- An internal workflow format is used in order to shield converters from any possible changes in KNIME's API
- Separation of concerns

- There is a disparity between the architecture of the desktop computer used to create the workflow and the architecture of the machines on the grid
- Different middlewares (*e.g.*, UNICORE, LSF) require different job configuration – solved using conversion tables



- Nodes imported via GKN use CTDs, therefore, generating the appropriate command line for the gUSE format is not complicated





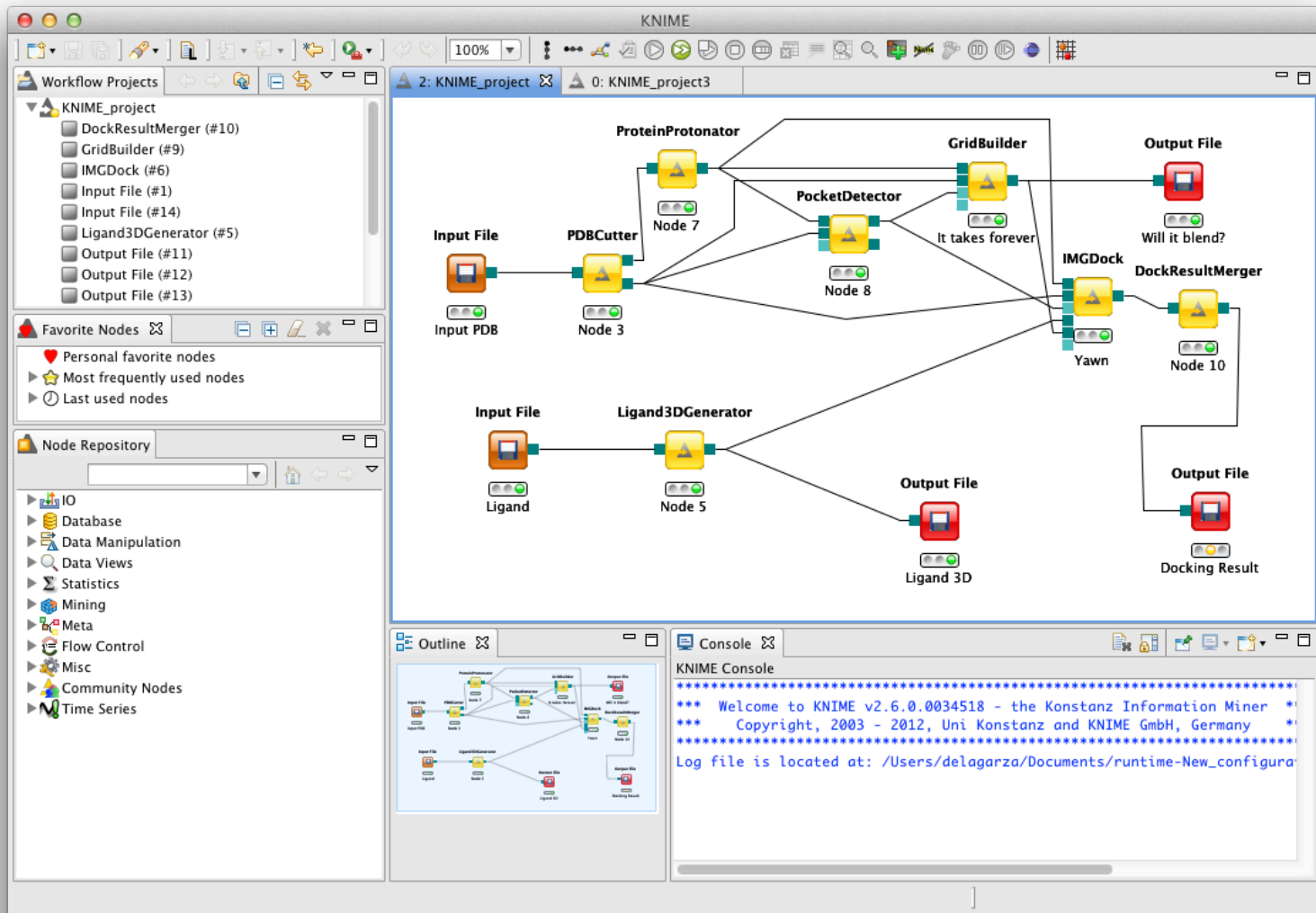
- Non-Generic KNIME Nodes are not command line tools, they are Java objects inside KNIME
- There is a way to execute KNIME via command line – KNIME Desktop must be also installed on the grid
- Workflows containing both types of nodes must be split into homogeneous sub-workflows and handled accordingly*

*work in progress

Docking on KNIME, then on the Grid

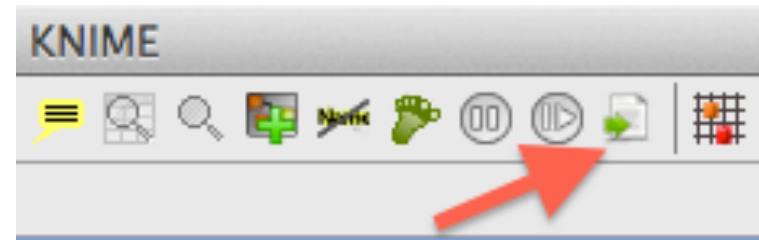
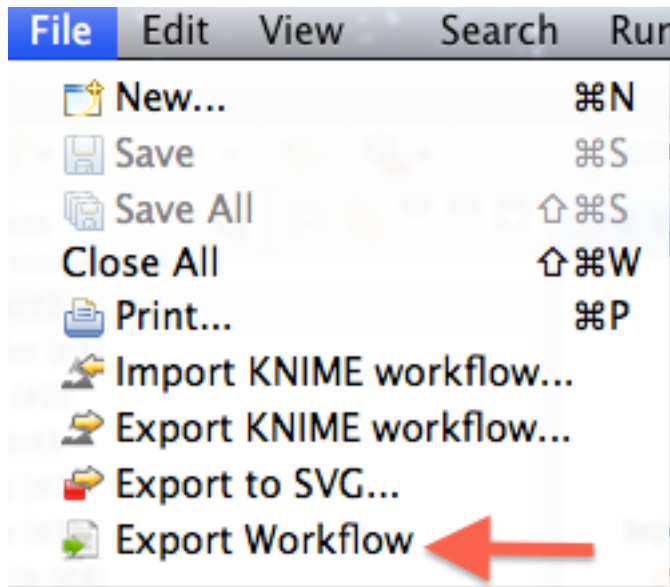


- It all starts with a workflow in KNIME



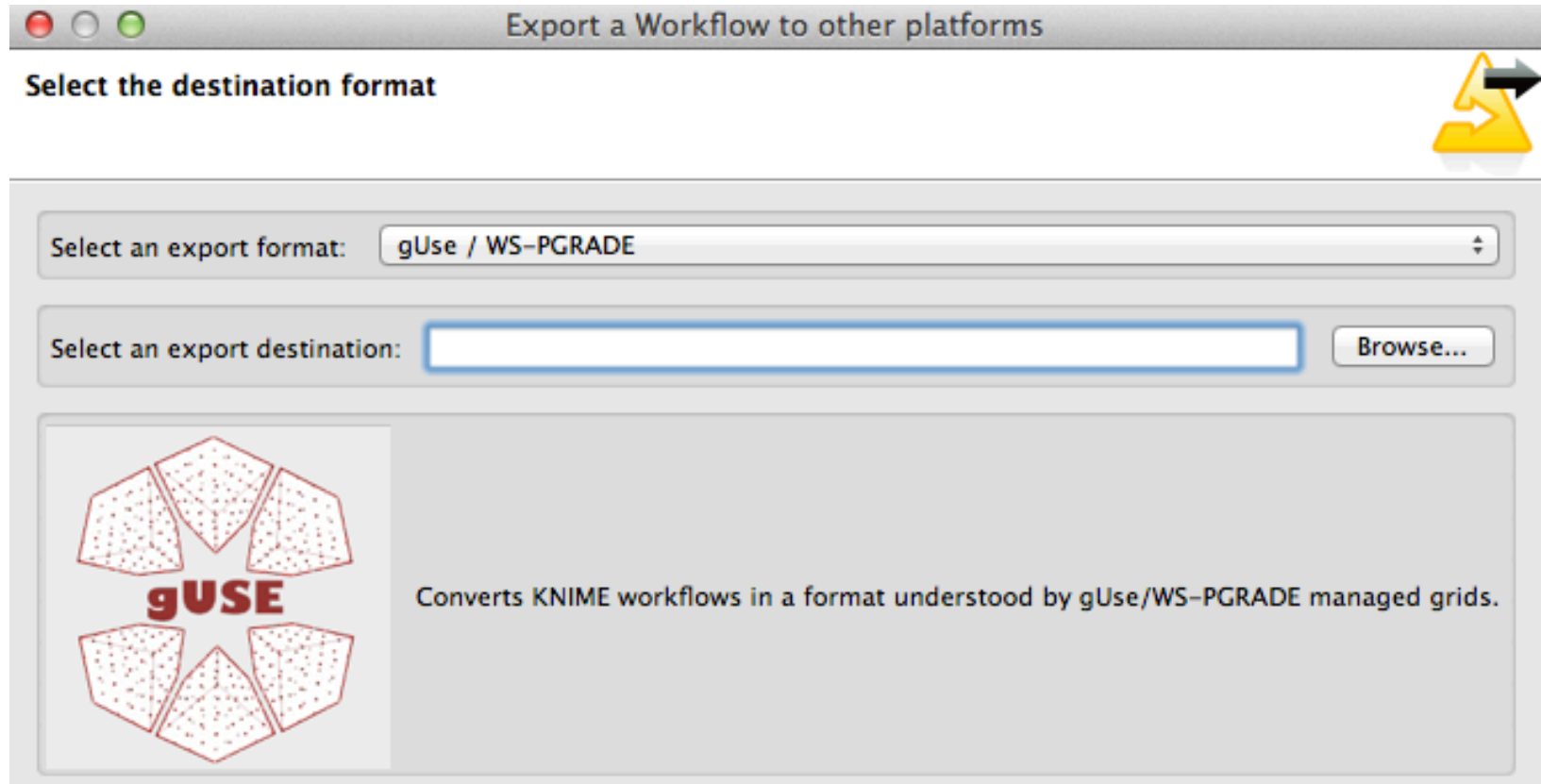


- The extension offers a simple user interface integrated into KNIME



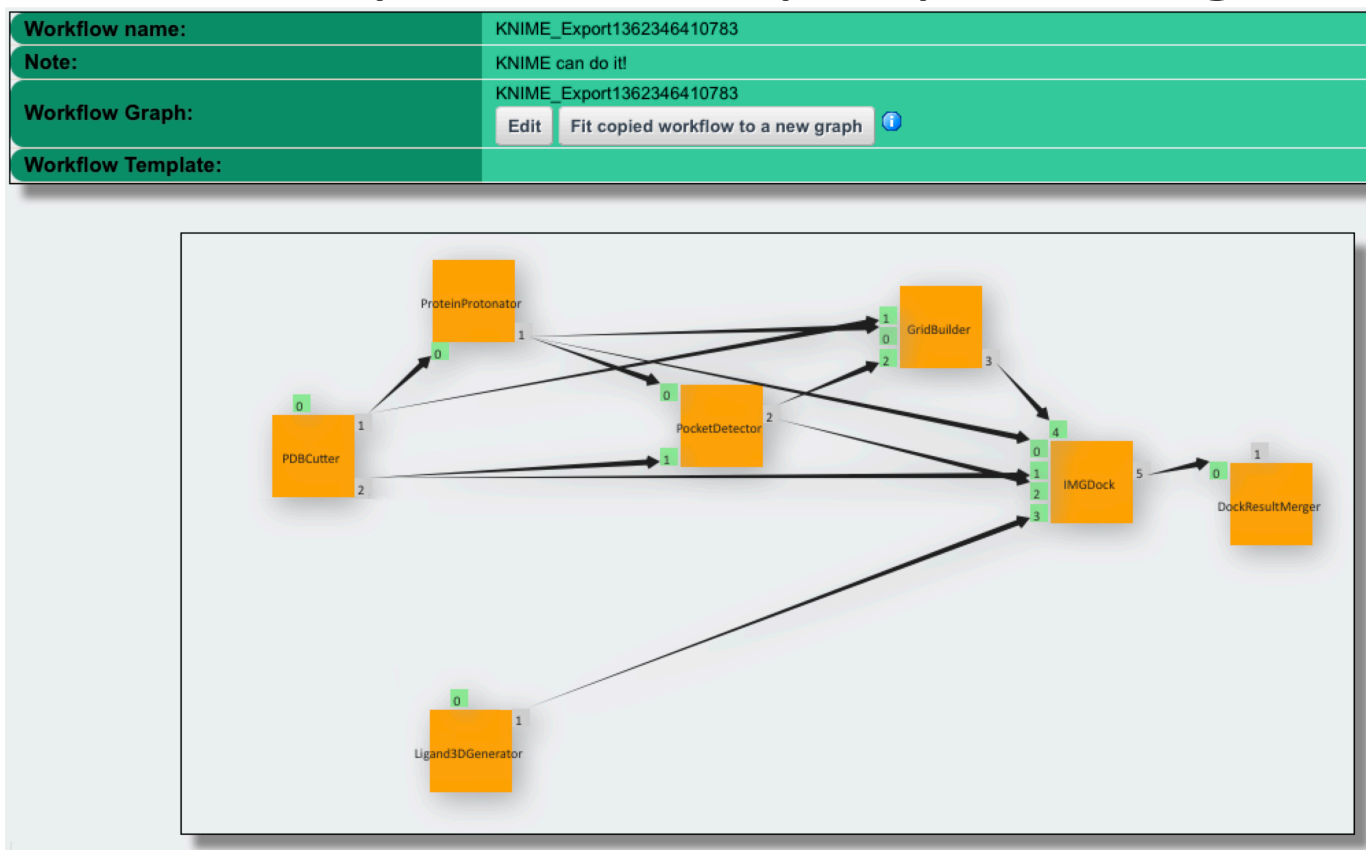


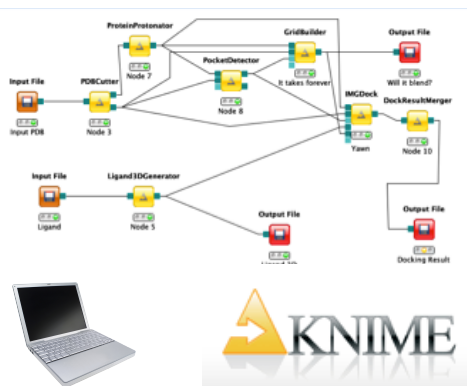
- The converter will generate a file in a specific format (*e.g.*, gUSE format) for a specific middleware (*e.g.* UNICORE)





- The generated file can be imported into WSPGRADE; input and output nodes in KNIME are converted to input and output ports in gUSE

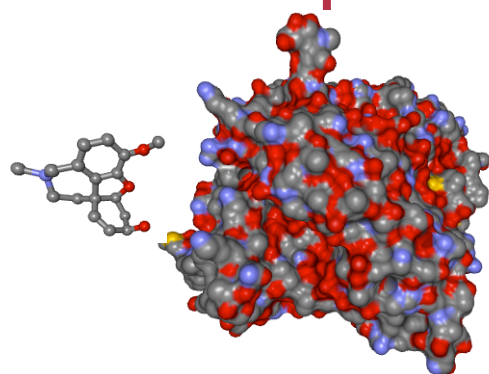




upload to gUSE
using novel converter



execute on a grid



scientific problem (docking)

Special Thanks to...

University of Tübingen

Jens Krüger

Charlotta Schärfe

Marc Röttig

Sven Nahnsen

Oliver Kohlbacher

MTA SZTAKI

Zoltán Farkas

István Márton

Ákos Balaskó

Peter Kacsuk

ETH Zürich

Béla Hullár

Peter Kunszt

University of Konstanz / KNIME.com

Peter Ohl

Throsten Meinl

Thomas Gabriel

Michael Berthold

Freie Universität Berlin

Stephan Aiche

Kurt Reinert

The MoSGrid Project (BMBF 01IG09006)

The SCI-BUS Project (EU 283481)

Questions

