

Fast Exploration of the QSAR Model Space with e-Science Central and Windows Azure

Simon Woodman

Jacek Cala

Hugo Hiden

Paul Watson

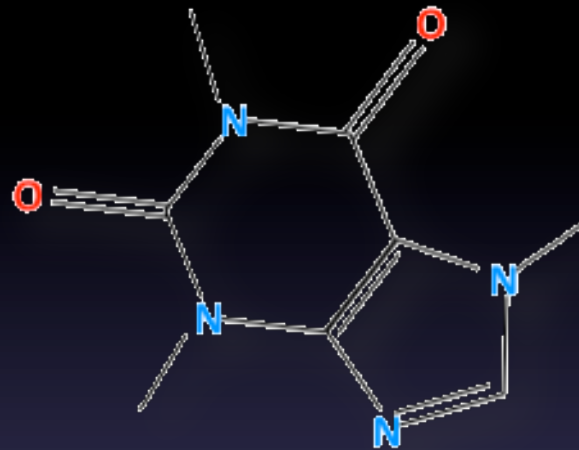
Newcastle University

The Problem

What are the properties of this molecule?

Toxicity

Solubility



Biological Activity

Perform experiments



Time consuming

Expensive

Ethical constraints

QSAR

Quantitative Structure Activity Relationship

$$\text{Activity} \approx f(\text{Structure})$$


More accurately, Activity related to a *quantifiable* structural attribute

$$\text{Activity} \approx f(\log P, \text{number of atoms}, \text{shape} \dots)$$



Currently > 3,000 recognised attributes

<http://www.qsarworld.com/>

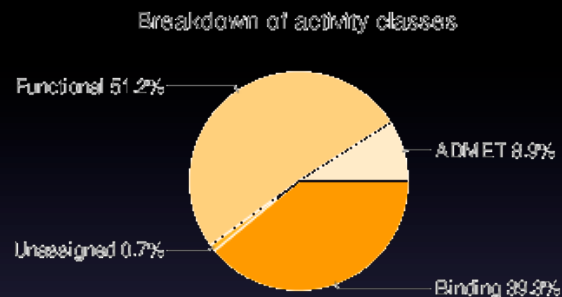
The alternative to Experiments

Predict likely properties based on **similar** molecules

CHEMBL Database: data on **622,824** compounds,
collected from **33,956** publications

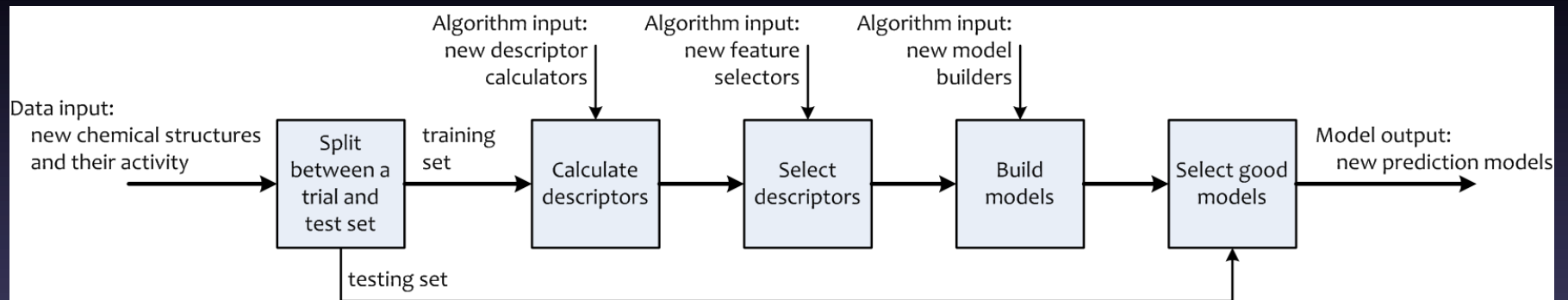
WOMBAT Database: data on **251,560** structures,
for over **1,966** targets

WOMBAT-PK Database: data on **1230** compounds,
for over **13,000** clinical measurements

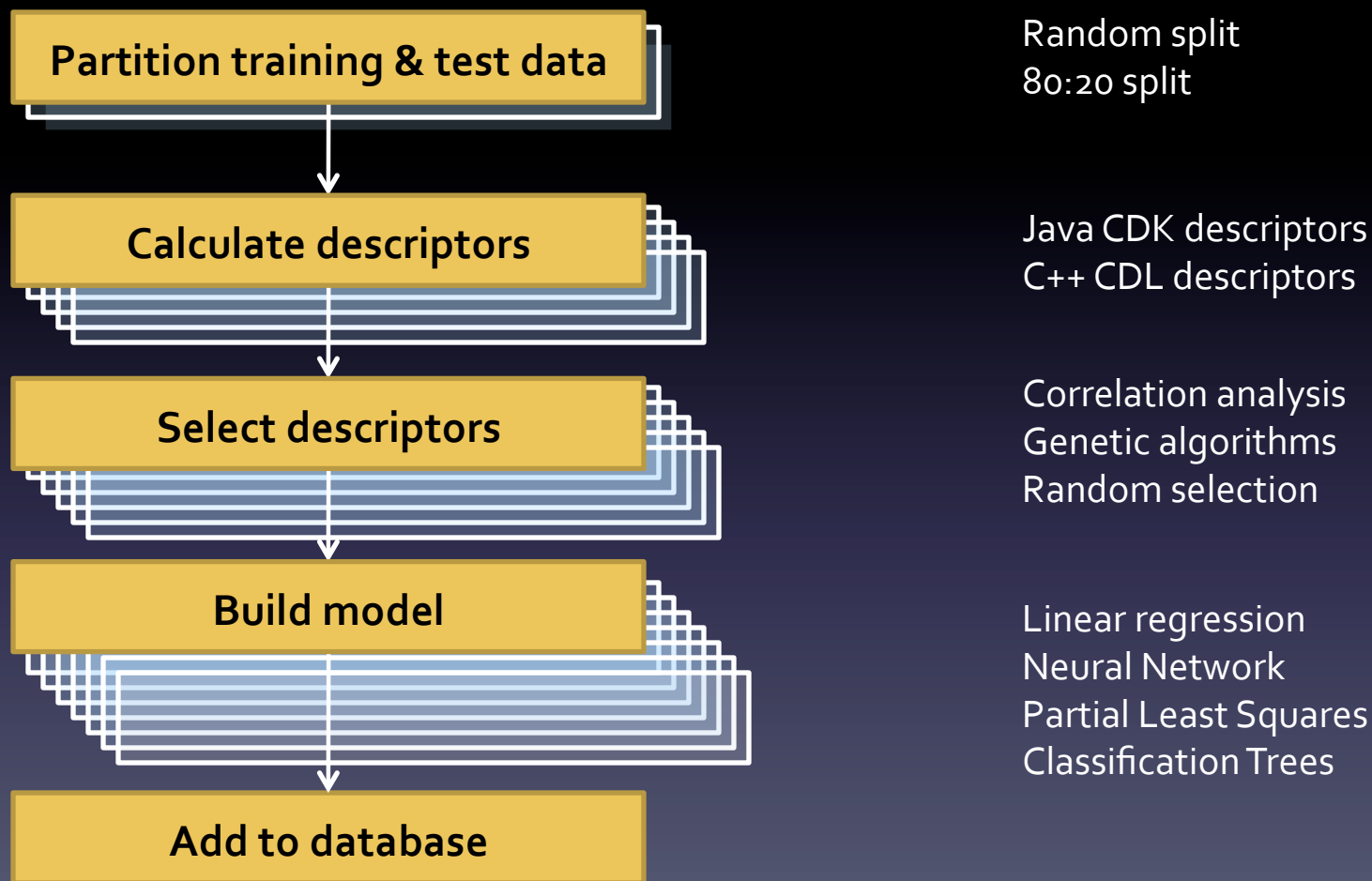


All these databases contain **structure** information and **numerical** activity data

Method

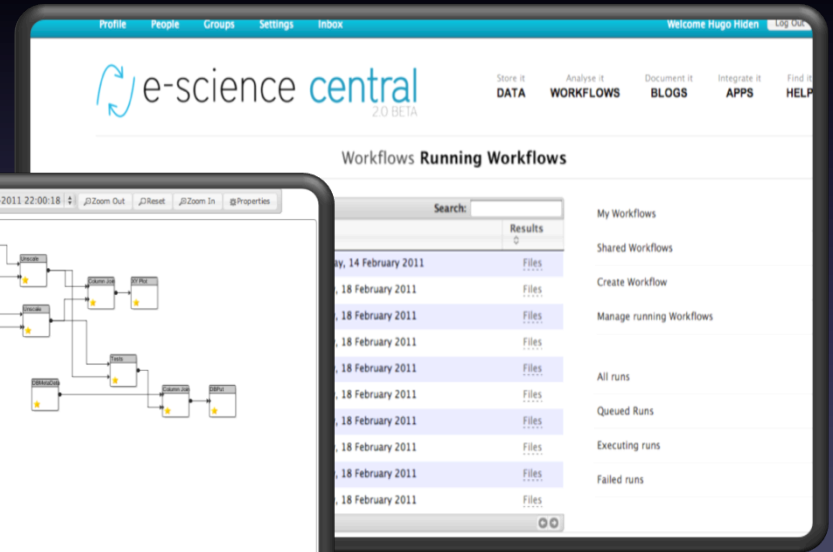
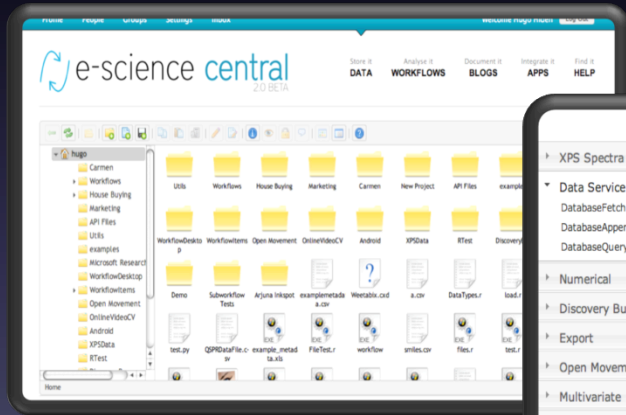


Branching Workflows



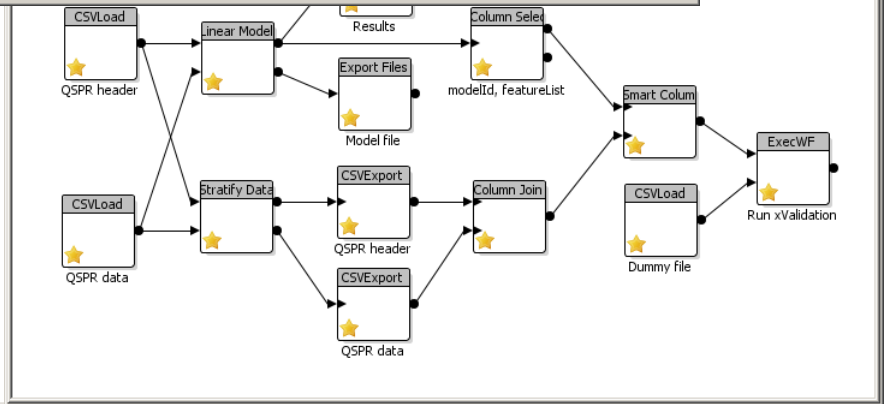
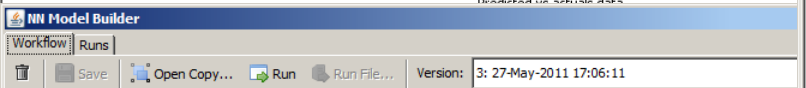
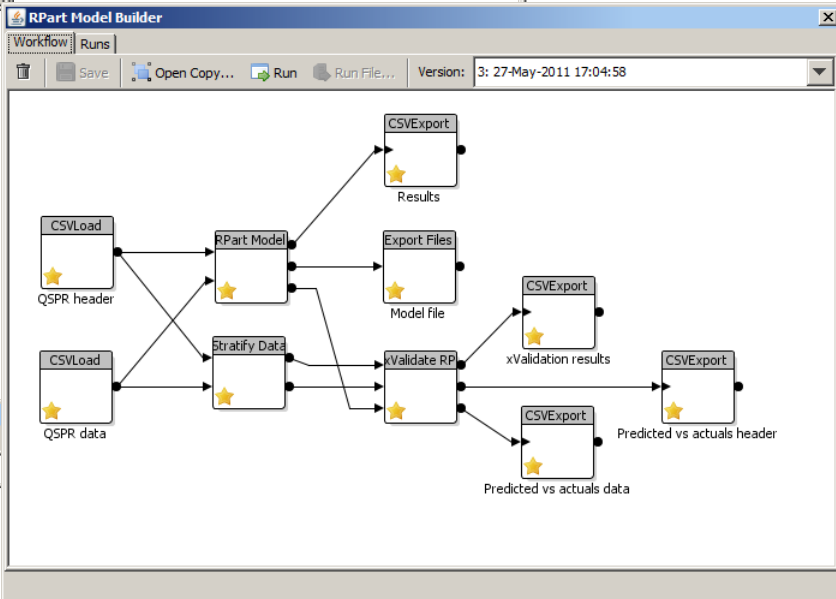
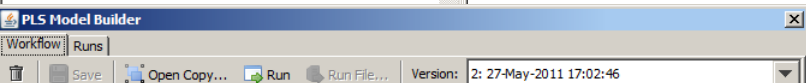
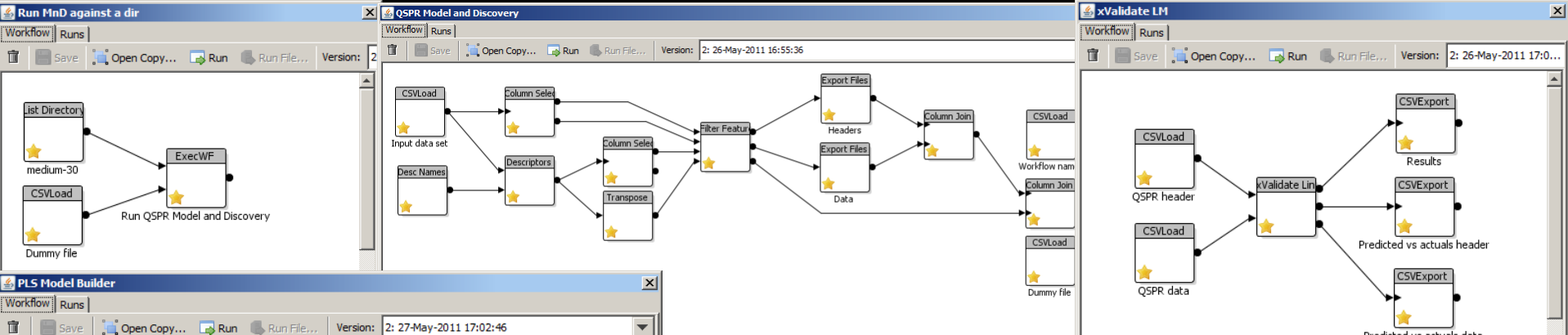
e-Science Central

Platform for cloud based data analysis

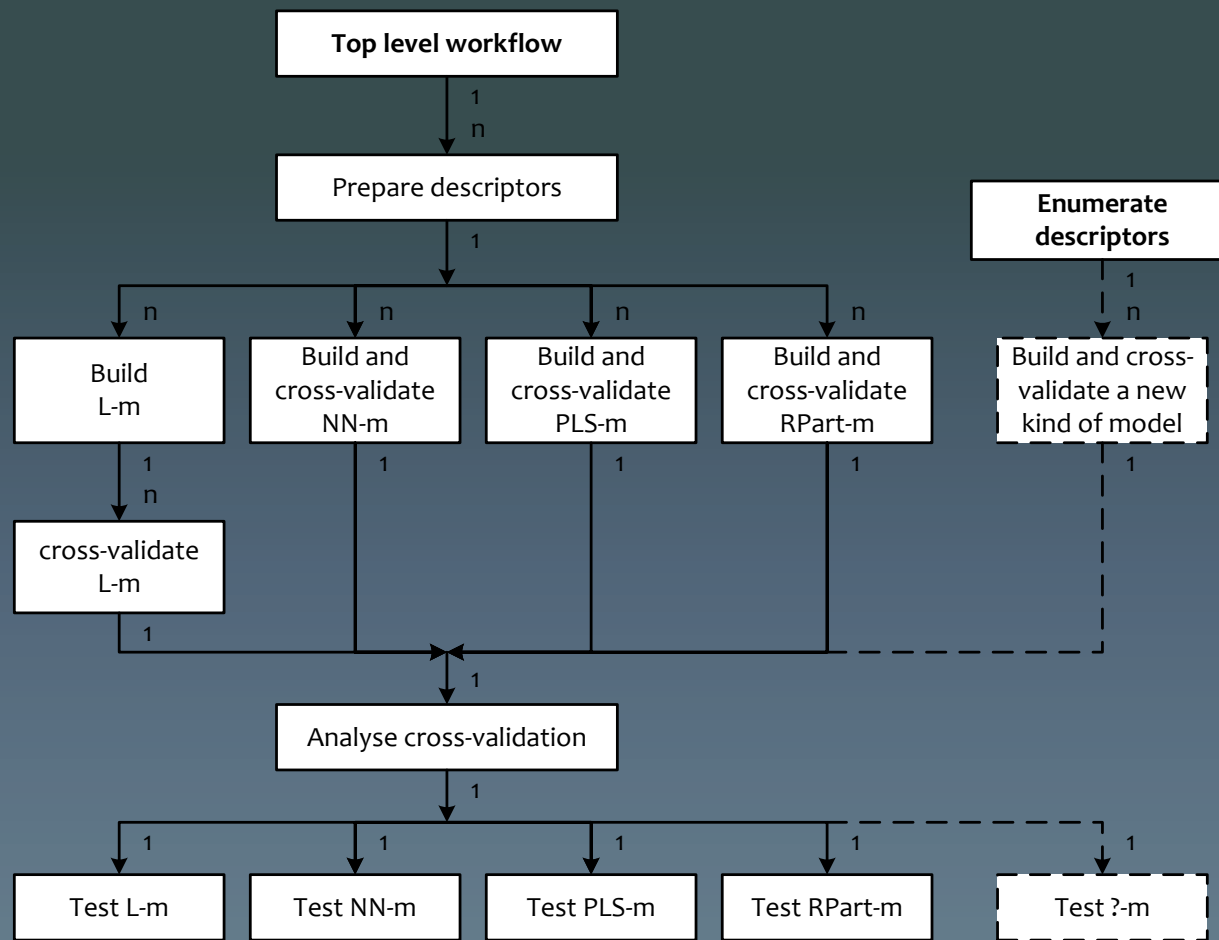


Azure
EC2
On Premise

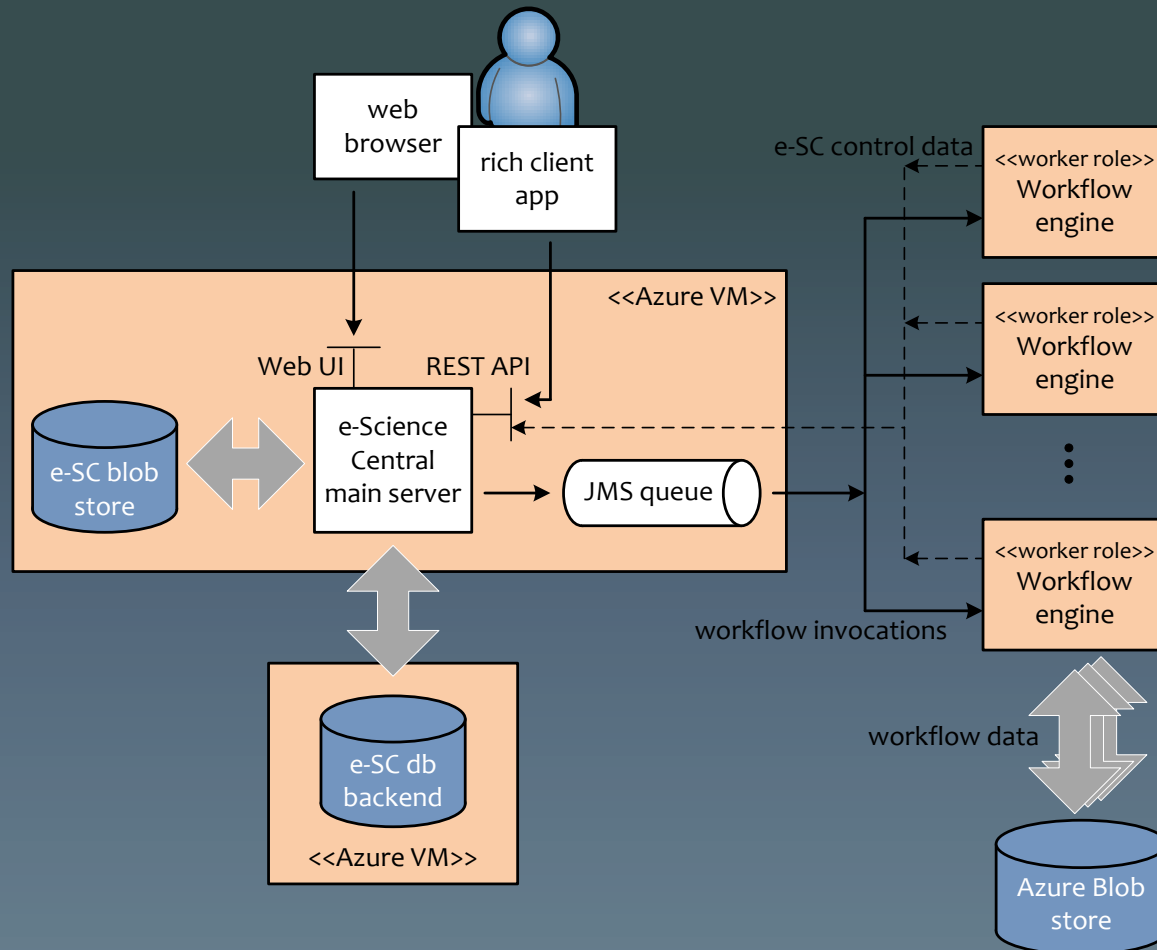
Java
R
Octave
Javascript



QSAR Implementation

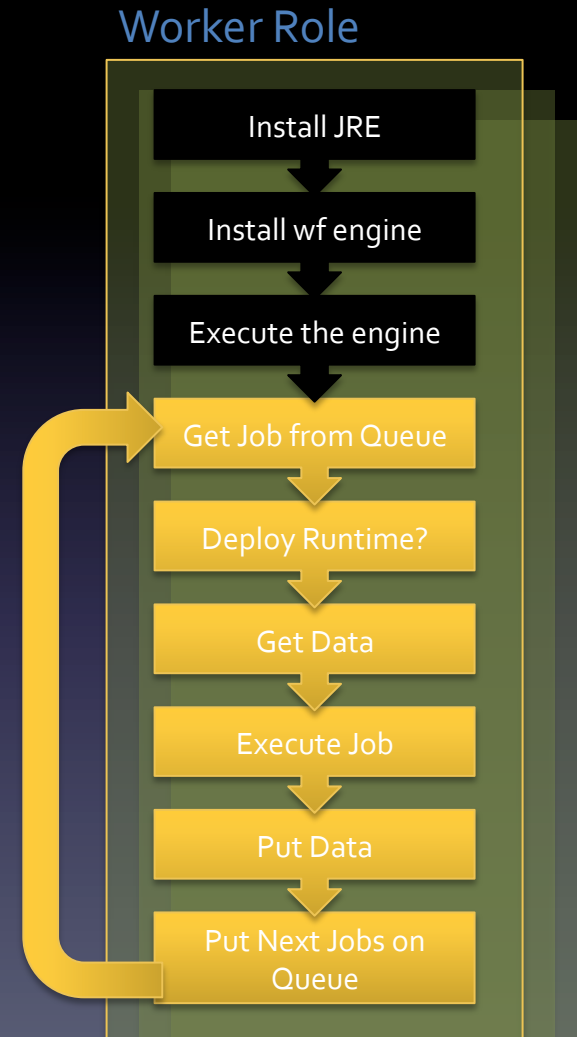


Azure e-Science Central



Workflow Architecture

- Single Message Queue
 - Worker Failure Semantics
 - Elasticity
- Runtime Environments
 - R
 - Octave
 - Java
- Deployed only once

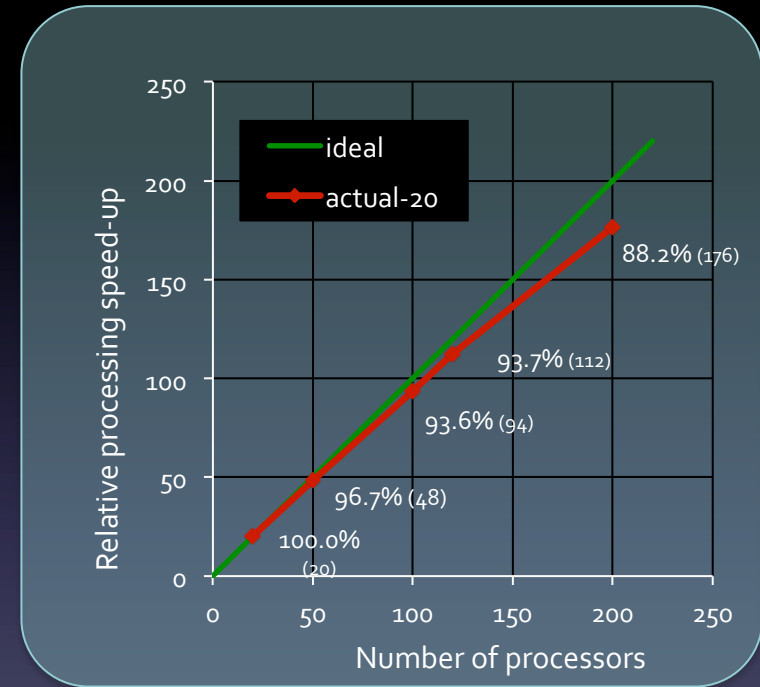


Results

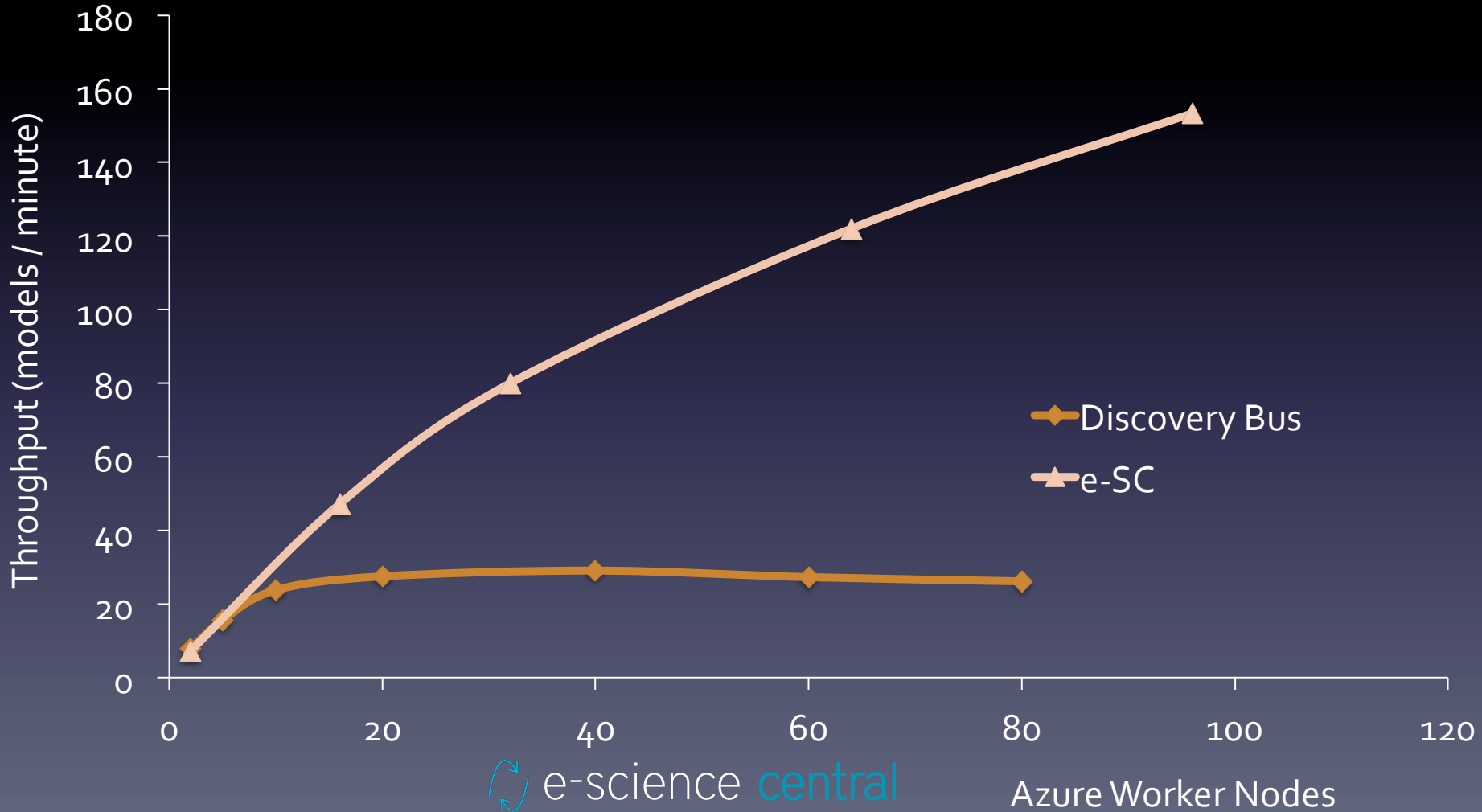
- 460K workflow executions
- 4.4M service calls
- 250k models
 - Linear Regression
 - PLS
 - Rpartitioning
 - Neural Net
- QSAR Explorer
 - Browse
 - Search
 - Get Predictions

Evaluation

Number of cores	1 + server	20 + server	200 + server
Response time	11d 20h 03m 40s	13h 00 m 11s	1h 28m 28s
Speed-up	1	20	176
Cost	\$62.64	\$40.32	\$51.84



Evaluation



Cloud Applicability

- Bursty
 - ChEMBLdb updates (delta 10%)
 - New Modelling Methods (???)
- Performance depends on how *chatty* the problem is
 - Deploy (incl download) dependencies once
 - Avoid storage bottlenecks

Performance is great but ...

Drug Development requires us to capture the
data and the **process**

Provenance Requirements

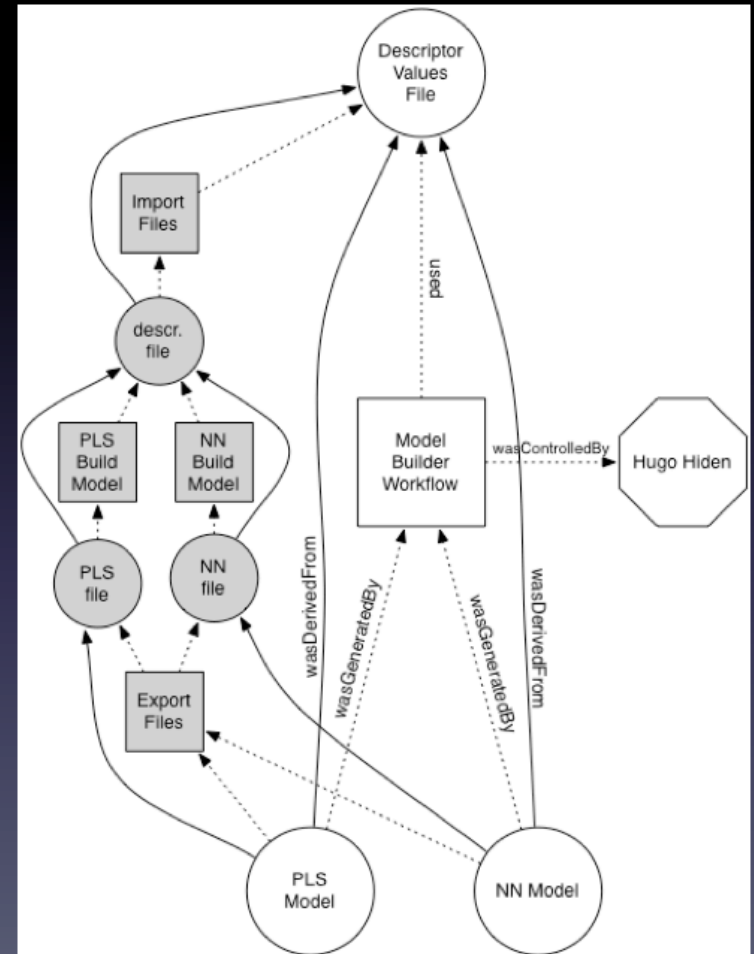
- How was a model generated?
 - What algorithm?
 - What descriptors
- Are these results reproducible?
- How have bugs manifested?
 - Which models affected
 - How do we regenerate affected models?
- Performance Characteristics
- How do we deal with new data?

Storing Provenance

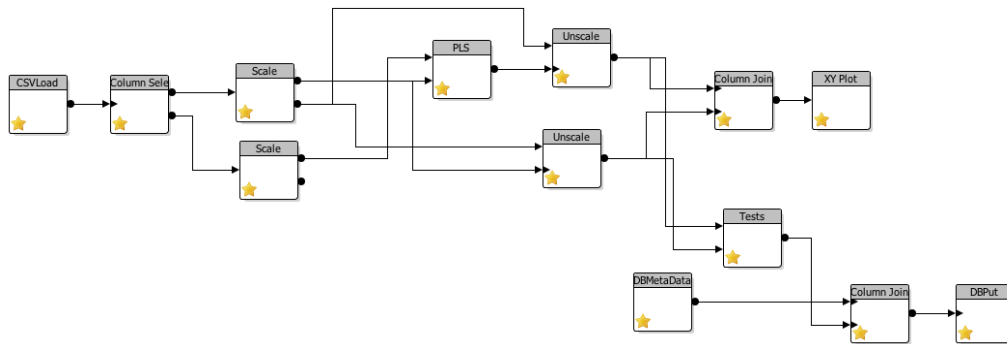
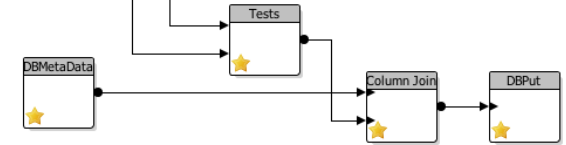
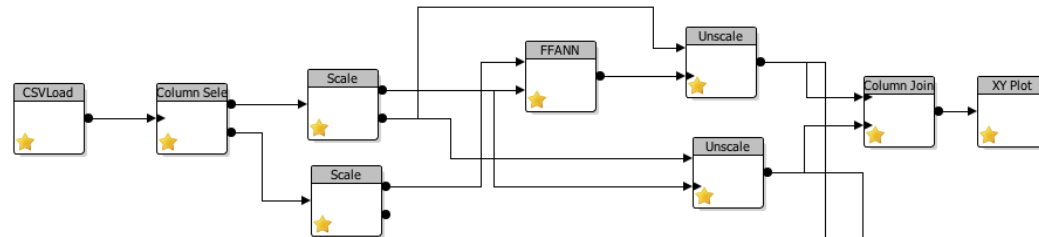
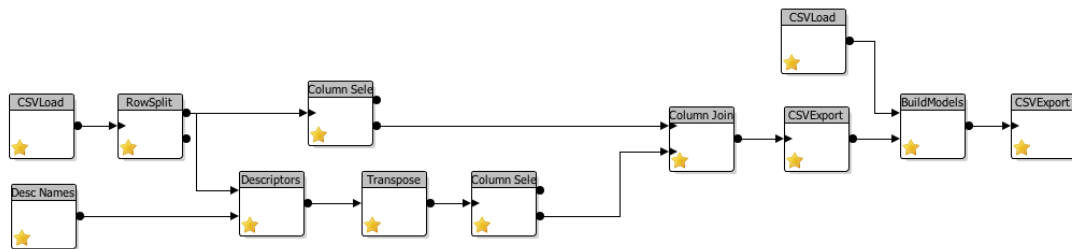
- Neo4j
 - Open Source Graph Database
 - Nodes/Relationships + properties
 - Querying/traversing
- Access
 - Java lib for OPM
 - e-SC library built on top of OPM lib
 - REST interface
- Options for HA and Sharding for performance

Provenance Model

- Based on OPM
 - Processes, Artifacts, Agents
- Directed Graph
- Multiple views of provenance
 - Dependent on security privileges

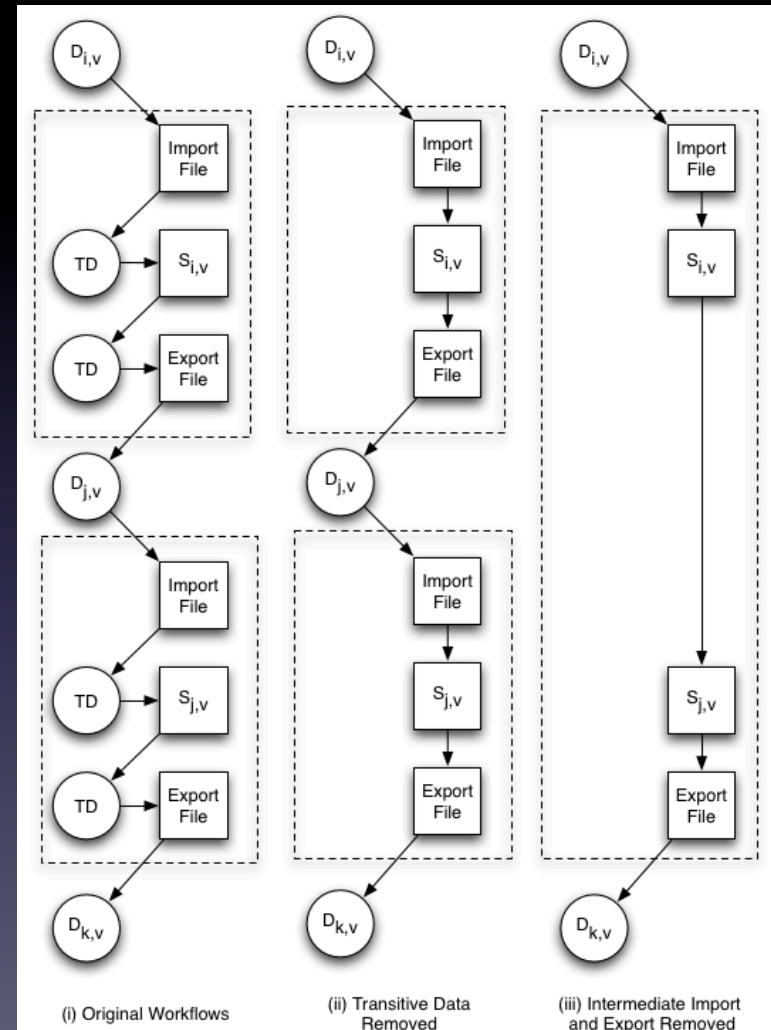


Multiple Linked Workflows



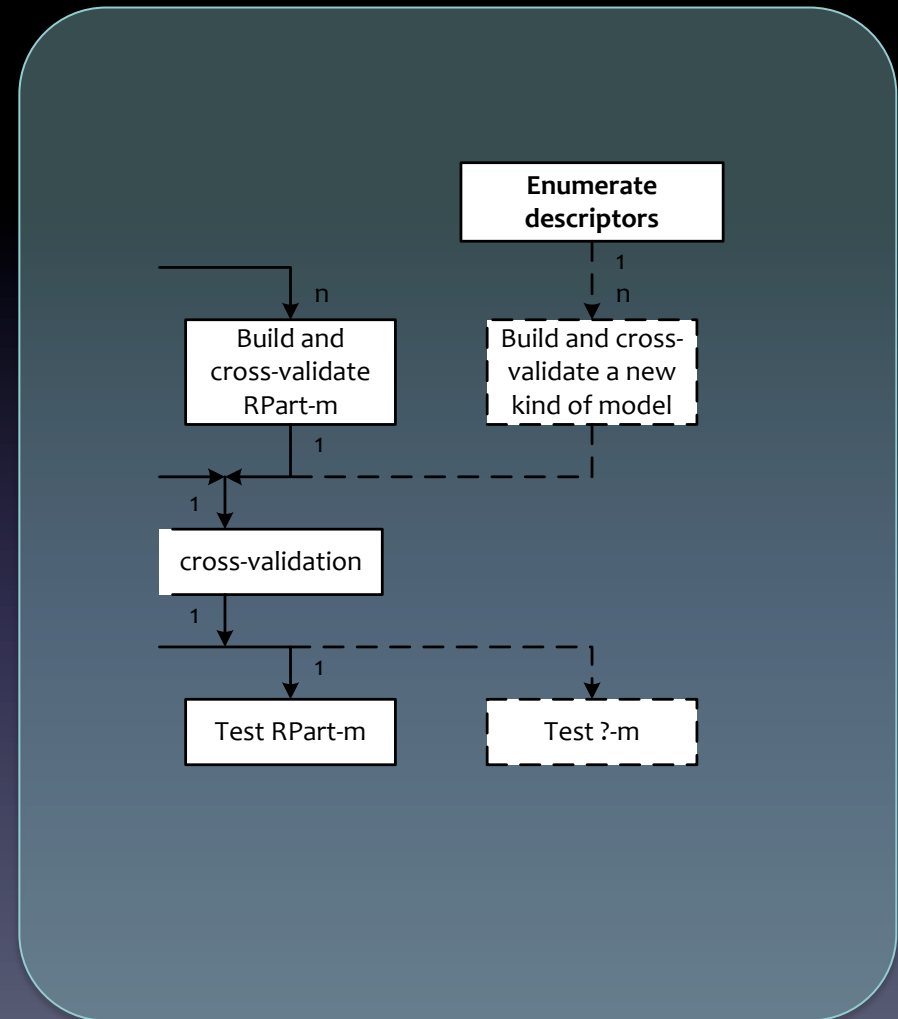
Spanning Multiple Workflows

- Regenerate Data with services updated
 - Not everything in a single workflow
- Create a single “Virtual Workflow” spanning multiple physical workflows
- Work performed on data sets by different people over a period of time



Adding new model builders

1. Add new block
 2. Mine the provenance
 3. Dynamically create virtual workflows
 4. One invocation per data set
- Work in progress...



Future Work

- Scale 200+ nodes
 - Database is bottleneck
- Provenance visualization
- Meta-QSAR
 - Provenance Mining

Questions?

- Thank you to our generous funders
 - EU FP7 - VENUS-C (RI-261565)
 - RCUK – SiDE (EP/Go66019/1)