# ChEMBL & RDKit

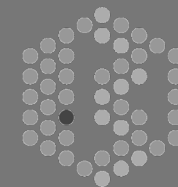## Marrying Open Data with Open Tools

**Rodrigo Ochoa**

Mark Davies

Francis Atkinson

George Papadatos

John Overington

EMBL-EBI

# Outline

- Introduction to EBI and ChEMBL
- Loading ChEMBL (PostgreSQL + RDKit)
- Comparative studies between RDKit and Symyx
- Web interface

# INTRODUCTION

# EBI structure

**Genomes**
Ensembl
Ensembl Genomes
EGA

**Nucleotide sequence**
ENA

**Functional genomics**
ArrayExpress
Expression Atlas

**Protein activity**
IntAct , PRIDE

**Protein Sequences**
UniProt

- ChEMBL database
  - Curation
  - Interface
  - Research group
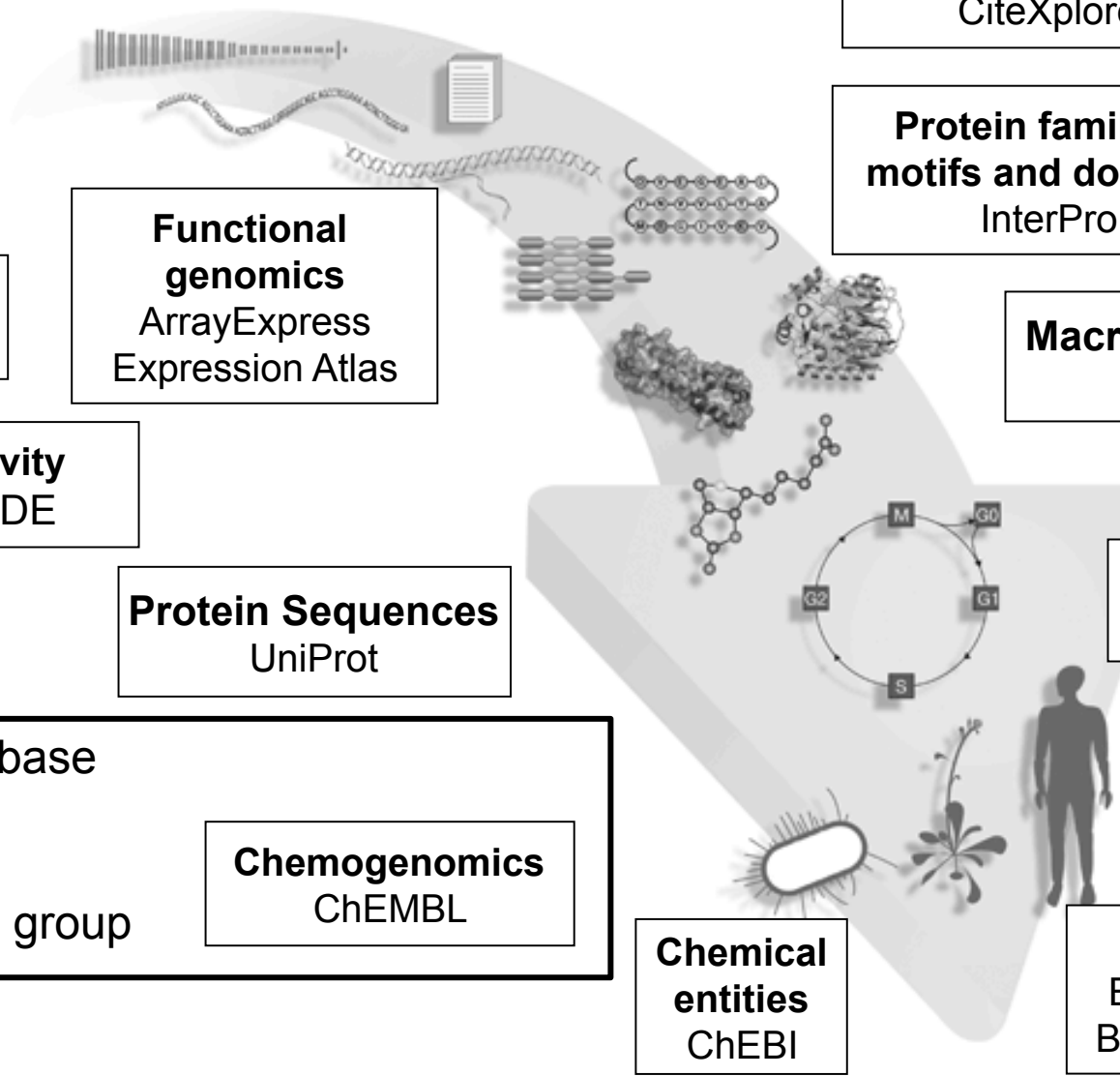
**Chemogenomics**
ChEMBL

**Chemical entities**
ChEBI

**Literature and ontologies**
CiteXplore, GO

**Protein families, motifs and domains**
InterPro

**Macromolecular**
PDBe

**Pathways**
Reactome

**Systems**
BioModels
BioSamples

EMBL-EBI

# What is the ChEMBL database?

- A freely-available, curated source of small molecules, targets, assays and bioactivity data
- Core data is from the primary Med Chem literature
  - J. Med. Chem., Bioorg. Med. Chem. Lett., J. Nat. Prod.

- Information extracted
  - Compounds tested
  - Assays performed
  - Biological targets of assays
  - Activities of compounds in assays

- Structures and data curated in-house to ensure quality

**Compounds**

**Bioactivities**

**Targets**

>Thrombin
MAHVRGLQLPGCLALAALCSLVHSQHVFLAPQQARSLLQRVRRANTFLEEVRKGNLERECVEETCS
YEEAFEALESSTATDVFWAKYTACETARTPRDKLAACLEGNCAEGLGTNYRGHVNITRSGIECQLW
RSRYPHKPEINSTTHPGADLQENFCRNPDSSTTGPWCYTTDPTVRRQECSIPVCGQDQVTVAMTPR
SEGSSVNLSPPLEQCVPDRGQQYQGRLAVTTHGLPCLAWASAQAKALSKHQDFNSAVQLVENFCRN
PDGDEEGVWCYVAGKPGDFGYCDLNYCEEAVEEETGDGLDEDSDRAIEGRTATSEYQTFFNPRTFG
SGEADCGLRPLFEKKSLEDKTERELLESYIDGRIVEGSDAEIGMSPWQVMLFRKSPQELLCGASLI
SDRWVLTAAHCLLYPPWDKNFTENDLLVRIGKHSRTRYERNIEKISMLEKIYIHPRYNWRENLDRD
IALMKLKKPVAFSDYIHPVCLPDRETAASLLQAGYKGRVTGWGNLKETWTANVGKGQPSVLQVVNL
PIVERPVCKDSTRIRITDNMFCAGYKPDEGKRGDACEGDSGGPFVMKSPFNNRWYQMGIVSWGEGC
DRDGKYGFYTHVFRLKKWIQKVIDQFGE

Compound

Assay

SAR Data

$K_i$ = 4.5 nM

APTT = 11 min

EMBL-EBI

# ChEMBL 14 (latest version)

**ChEMBL14**

Compounds: **1,213,239**

Assays: **644,734**

Targets: **9,003**

Publications: **46,133**

Activities: **10,129,256**

Data sources: **10**

Increase of >200,000 compounds
from literature since ChEMBL01

# ChEMBL interface

# Chemical search in ChEMBL

- A variety of sketchers (Marvin, JDraw, JME)

Compound Sketcher: Please select....

✓ Substructure Search
Similarity

100%   Fetch Compounds
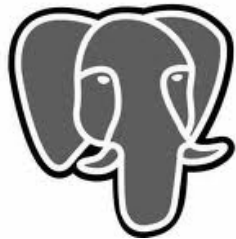
**List Search**

○ SMILES Search  ⊙ ChEMBL ID Search  ○ Keyword Search

Please enter a list of Compound IDs, keywords, or SMILES separated by newlines

Fetch Compounds

- ChEMBL chemical cartridge:

**accelrys®**
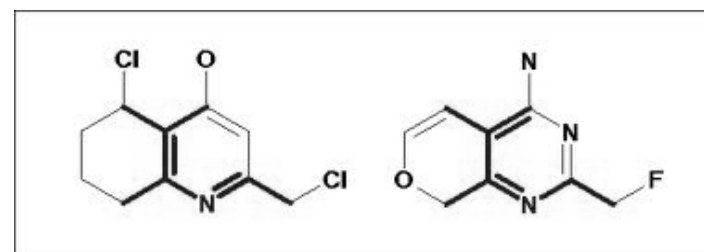Symyx

# LOADING CHEMBL

# RDKit overview

- An open source collection of chemoinformatics and machine-learning software written in C++ and Python

- There is available a PostgreSQL cartridge

(http://code.google.com/p/rdkit/wiki/DatabaseCartridgeReferenceDocumentation)

RDKit
Open-Source Cheminformatics
and Machine Learning

# RDKit PostgreSQL cartridge

- Molecular conversion (SMILES, SMARTS, CTAB)
- Substructure search
- Similarity Search

    Fingerprints: Morgan (ECFP-like), Atom-Pair, Torsion …
    Similarity Coefficients: Tanimoto, Dice

- Molecular properties calculation

# Bingo chemical cartridge

- Open Source chemical cartridge provided by Indigo
- Support PostgreSQL and Oracle schemas

Unfortunately it could not be compared with the others chemical cartridges (bug reported in the Google group forum)

GGA
GGA Software Services

# PostgreSQL ChEMBL version

• ChEMBL originally is available in Oracle (Based on this version a MySQL dump file is available)

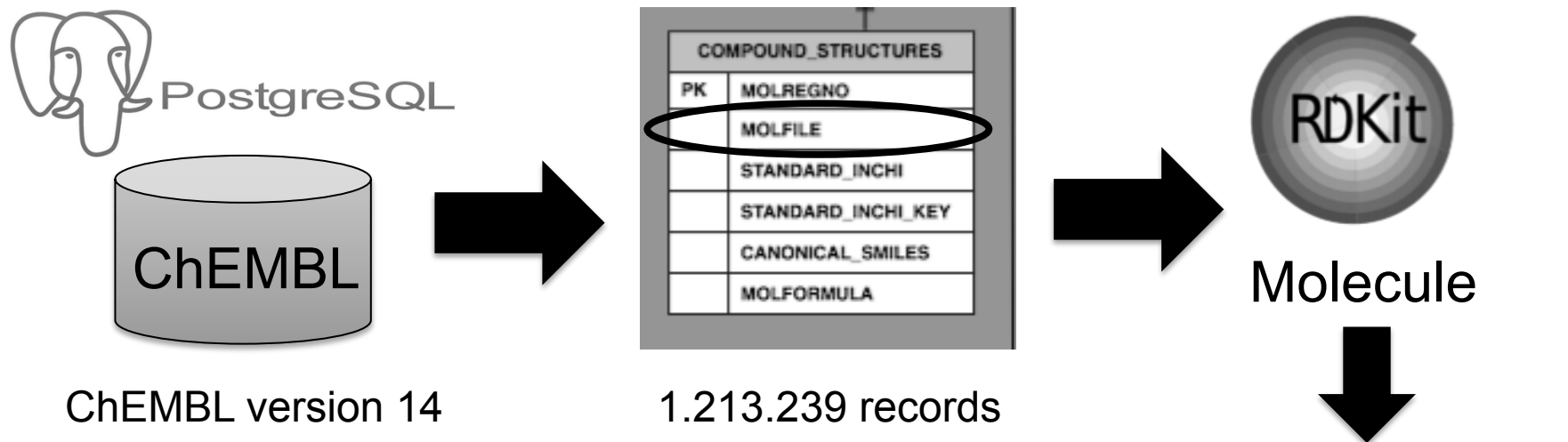**Task:** Migrate the Oracle Schema to a PostgreSQL schema

**Tool:**



Ora2Pg
Moves Oracle databases to PostgreSQL

Open Source Kit of Perl libraries

Highly configurable

http://chembl.blogspot.co.uk/2012/08/chembl-postgresql.html

EMBL-EBI

# ChEMBL + cartridges

PostgreSQL

**ChEMBL**

**ChEMBL version 14**

| COMPOUND_STRUCTURES | |
|---|---|
| PK | MOLREGNO |
| | MOLFILE |
| | STANDARD_INCHI |
| | STANDARD_INCHI_KEY |
| | CANONICAL_SMILES |
| | MOLFORMULA |

**1.213.239 records**

RDKit

**Molecule**

| | mols_rdkit Table |
|---|---|
| | fps_rdkit Table |

| Chemical Cartridge | # Molecules correctly built |
|---|---|
| RDKit | 1.146.045 |
| Symyx | 1.147.790 |
| Bingo | 1.138.682 |

is_valid_ctab() RDKit function

EMBL-EBI

# Building differences

## Not built by RDKit (1565)

Non pure organics (1506)

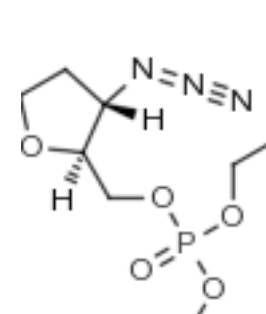Mainly boron clusters & organometallics



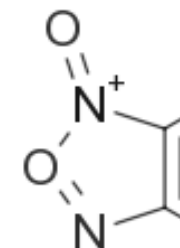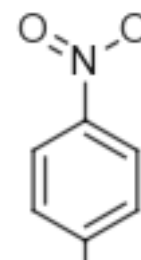EMBL-EBI

# High valency (164; overlap with inorganics)

Some are legitimate structures and probably should be handled:

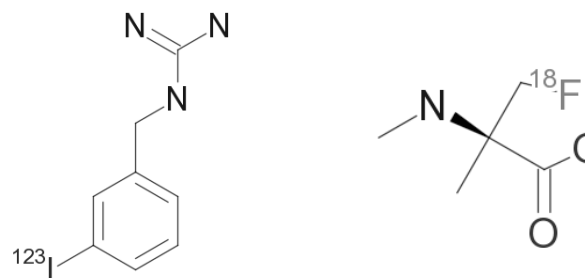Some need standardizing in ChEMBL to *e.g.* charge-separated form:
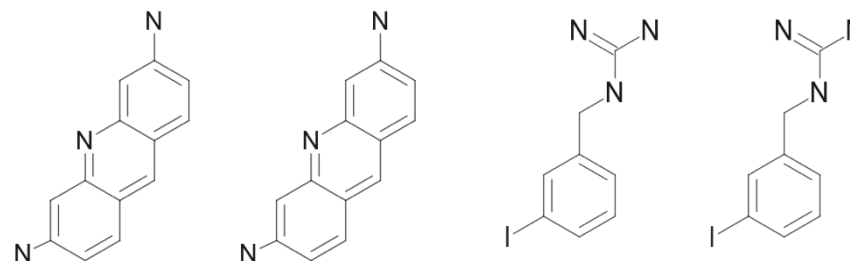
A few need fixing in ChEMBL:
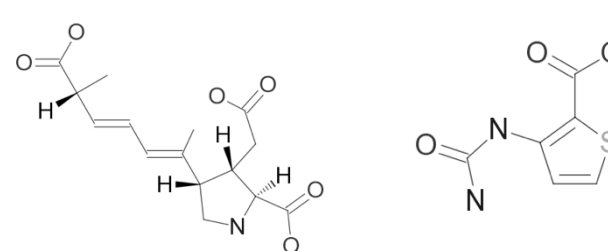
# Not built by Symyx (50)

Unusual isotopes (10); used in PET, as tracers *etc*:

Duplicated structures (35):

No obvious reason for failure (5):

**1,145,994** parent structures built by both systems
*i.e.* problems with only a very small minority of structures!

**COMPARISON**

# Cartridge comparison

- Chemical cartridges: RDKit and Symyx
- Type of queries: SMILES
- Objective: Detect discrepancies and the reasons behind them.
- RDKit configuration:

Version: 2012_06

OS: Linux Ubuntu 12.04 64 bits (Virtual Machine)
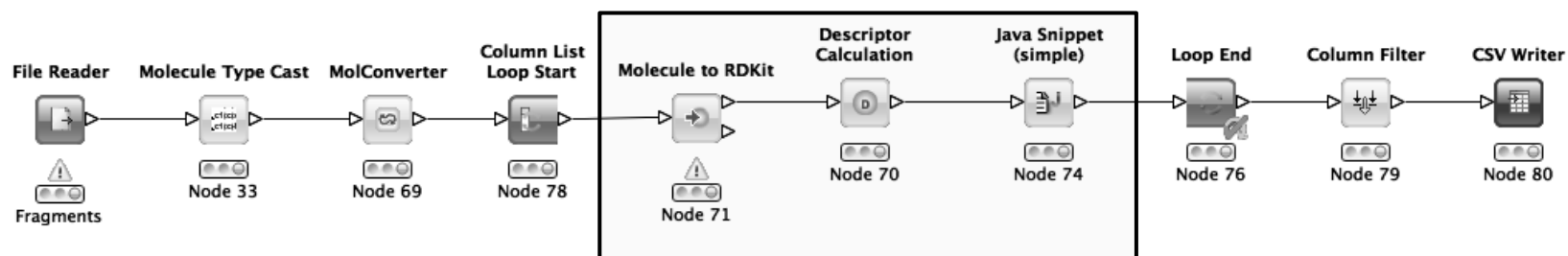
PostgreSQL version: 9.1.4

- Symyx configuration:
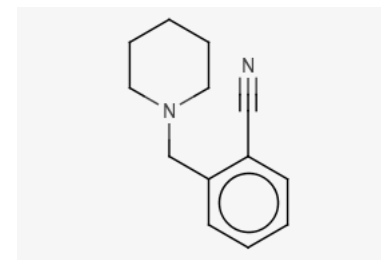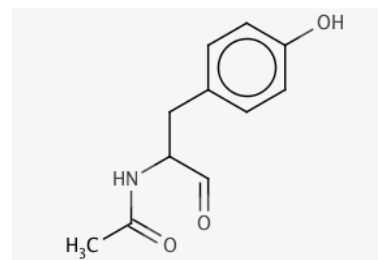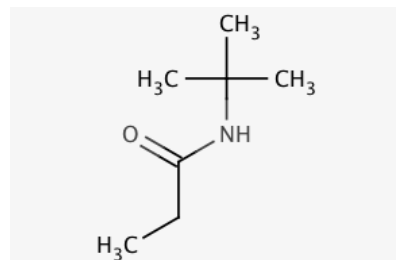
Version: 6.2
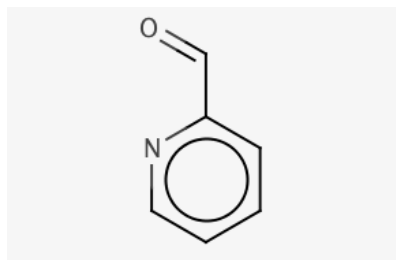
OS: Red Hat 5.8 64 bits (Server)

Oracle version: 11.1.0.7.0

# Fragment library

- From a set of ChEMBL molecules, RECAP* algorithm was applied.
- A filter based on the number of heavy atoms and molecular weight was done using the RDKit nodes in KNIME
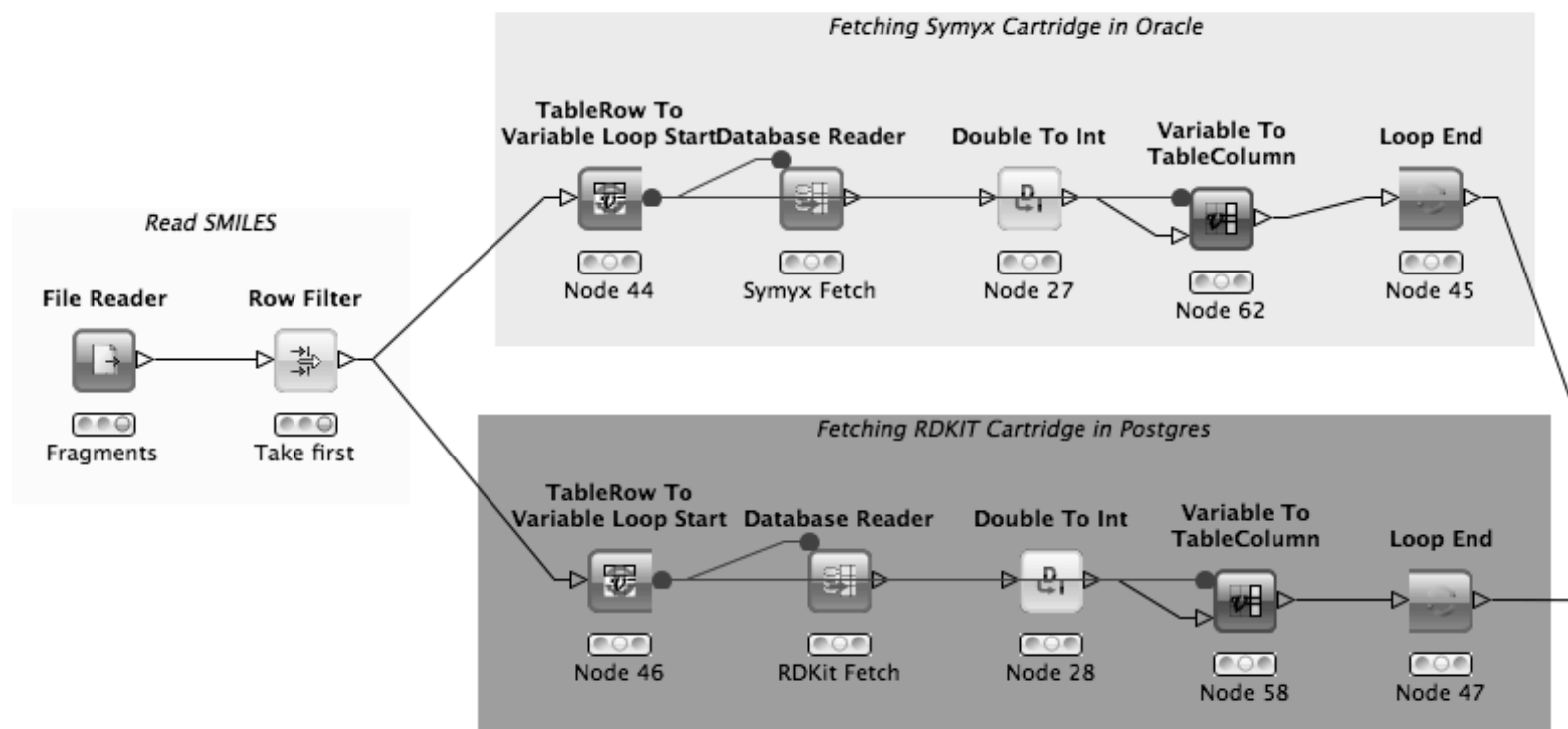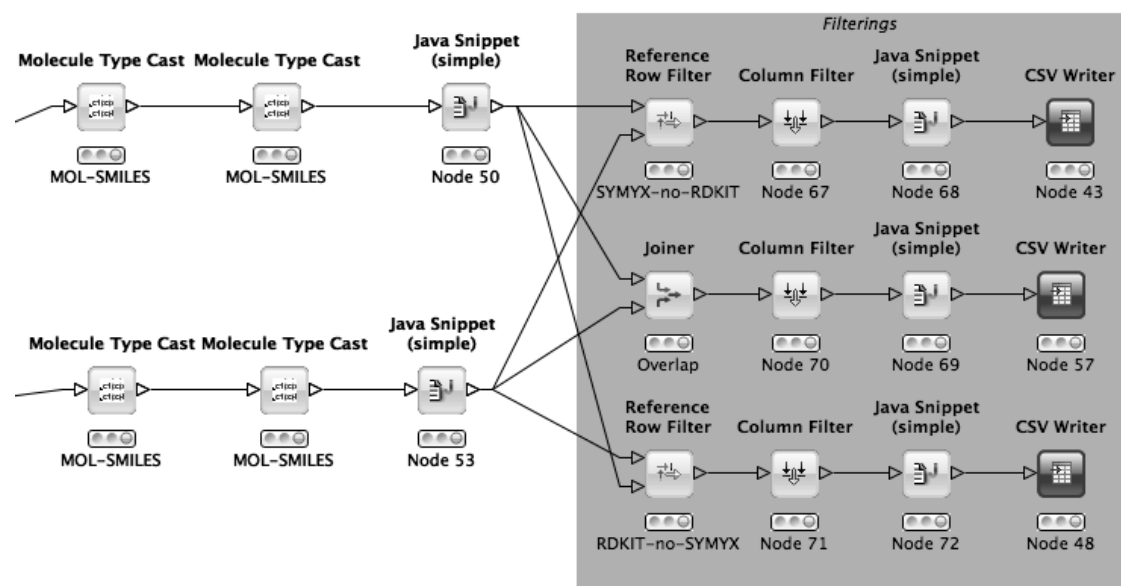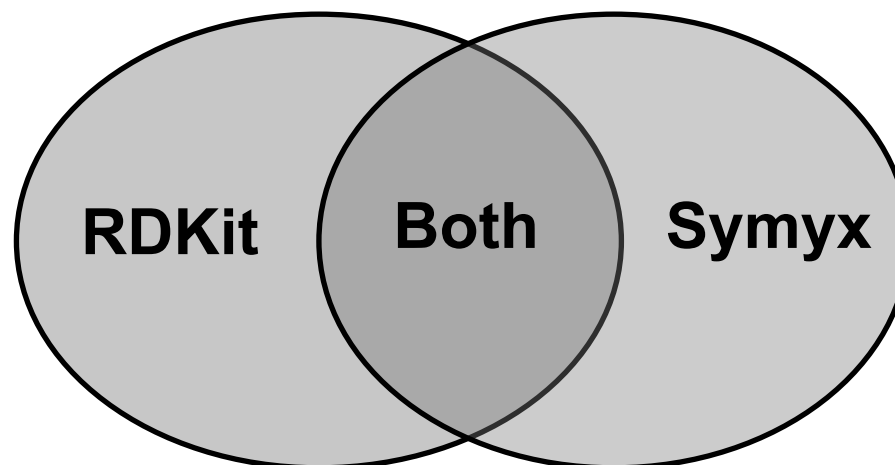


- 40 fragments were chosen

* Lewell XQ, Judd DB, Watson SP, Hann MM. RECAP - Retrosynthetic Combinatorial Analysis Procedure: A powerful new technique for identifying privileged molecular fragments with useful applications in combinatorial chemistry. *Journal of Chemical Information and Computer Sciences* 1998, 38:511-522.
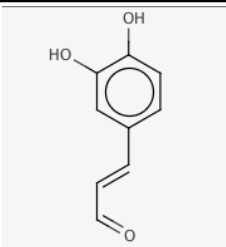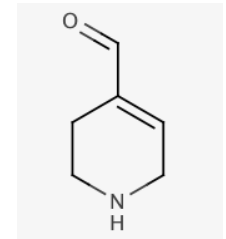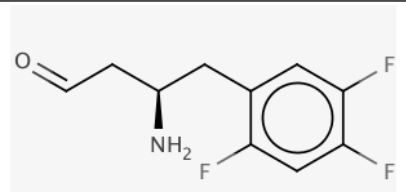
EMBL-EBI

# KNIME protocol

# KNIME results

# Some results …

| Fragment | Molecule | RDKit no Symyx | Symyx no RDKit | Join Symyx RDKit | Total RDKit | Total Symyx |
|---|---|---|---|---|---|---|
| 1 |  | 1595 | 5565 | 25404 | 26999 | 30969 |
| 5 |  | 52 | 268 | 18247 | 18299 | 18515 |
| 9 |  | 263 | 6 | 199 | 462 | 205 |
| 14 |  | 62546 | 2455 | 73510 | 136056 | 75965 |

# Aromaticity perception: benzene

RDKit SMARTS: `c1ccccc1`     Symyx molfile: 

- 1445 structures found by RDKit but not by Symyx
- Due to differing aromaticity models



- Sometimes highlights unusual choice of tautomer



EMBL-EBI

# Aromaticity perception: pyridone

RDKit SMARTS: `c1cccnc1=O`    Symyx molfile

- 6760 structures hit by RDKit but not by Symyx

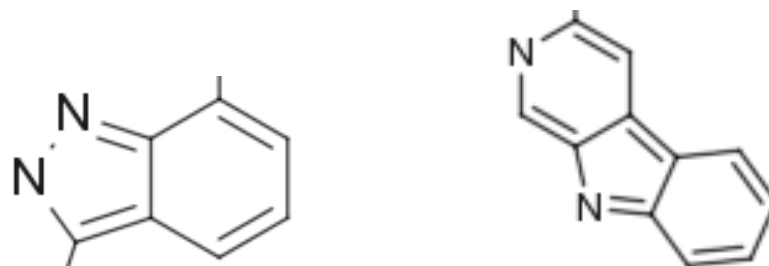  Bond is now aromatic and doesn't match double bond in query

- Can adjust Symyx query to catch these:

- Reduces discrepancy to 179...

# Aromaticity models

- Differing aromaticity models mean it can be difficult to code queries that give same hits in both systems

  Not a new problem!

  *e.g.* Symyx/MDL *vs*. Daylight

- MarvinSketch can be used to generate SMARTS queries appropriate for RDKit

  User would need to be aware of issues if using advanced query features

- Unlikely to be a real problem once transitioning users become accustomed to new conventions

# Stereochemistry

Chiral matching not currently handled by RDKit

Symyx *can* handle chiral queries

However, there are issues with chiral searching in ChEMBL

- Some inconsistences in setting of chiral flag

- Issue is being investigated

Best avoided for now !!!

# INTERFACE

# Interface

- Substructure, Similarity search against ChEMBL database
- Additionally, calculation of molecular properties

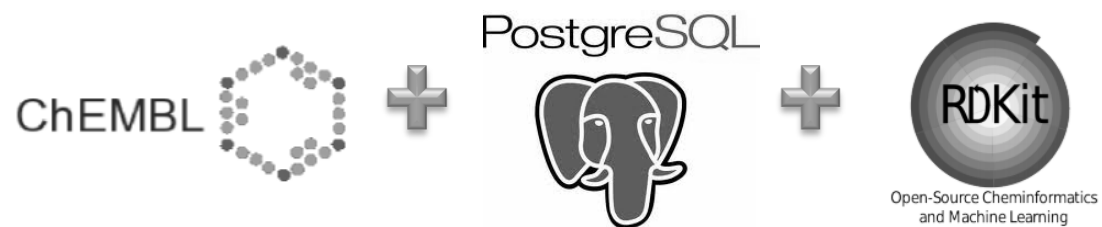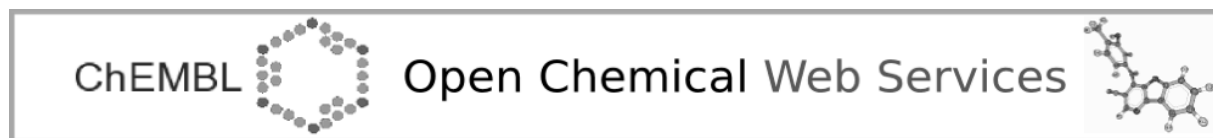- Substructure search using SMARTS, SMILES or MolFiles

## SUBSTRUCTURE AND EXACT SEARCH

**MENU**

- Home
- Tutorial
- Substructure and Exact Search
- Similarity Search
- Molecular Properties

**EXPLANATION**

In this section the user can choose between substructure and exact search, using as input format SMILES strings, SMARTS queries and MOL files stored in your computer
**NOTE:** The exact search only works for SMILES and MOL formats

**Select one option:**
- ○ Draw your structure
- ⦿ Input an string or a molfile stored in your computer

Please select one of the following formats for the input:  SMARTS ⬍

-- Enter the string of characters (for **SMARTS** queries):

[#6;X4]-1-[#6](=[#8])-[#7]-[#7]-[#6]-1=[#8]

( Run SMARTS )

EMBL-EBI

# SUBSTRUCTURE AND EXACT SEARCH

## MENU

- Home
- Tutorial
- Substructure and Exact Search ›
- Similarity Search ›
- Molecular Properties

### EXPLANATION

In this section the user can choose between substructure and exact search, using as input format SMILES strings, SMARTS queries and MOL files stored in your computer

**NOTE:** The exact search only works for SMILES and MOL formats

### Select one option:

- ⦿ Draw your structure
- ○ Input an string or a molfile stored in your computer

File  Edit  View  Insert  Atom  Bond  Structure  Calculations  Tools  Help



-- Select one kind of search: ⦿ SUBSTRUCTURE ○ EXACT

[ Search ]

EMBL-EBI

- Query results using substructure search



Generated with RDKit

- Similarity search using different Fingerprints and Similarity Coefficients



EMBL-EBI

- Query results using Similarity search

- Compound report (with links to the ChEMBL website)

# CHEMBL160451

ChEMBL Link: CHEMBL160451

| Canonical SMILES | CC1=CC(C)(C)Nc2cc3nc(O)cc(c3cc12)C(F)(F)F |
|---|---|
| Standard InChI | InChI=1S/C16H15F3N2O/c1-8-7-15(2,3)21-13-6-12-10(4-9(8)13)11(16(17,18)19)5-14(22)20-12/h4-7,21H,1-3H3,(H,20,22) |
| Standard InChI-Key | REOPBPDUXSVBMY-UHFFFAOYSA-N |



# Bioactivity Data

| Assay ID | Assay Type | Assay Relation | Value | Units | Target |
|---|---|---|---|---|---|
| 36279 | Ki | = | 115 | nM | CHEMBL1871 |
| 159377 | IC50 | = | 49 | nM | CHEMBL208 |
| 36107 | IC50 | = | 28 | nM | CHEMBL1871 |

( Back )

EMBL-EBI

- Molecular properties calculation

## MOLECULAR PROPERTIES

**QUERY DETAILS**

| Query | CC1=CC(C)(C)Nc2cc3oc(=O)cc(C(F)(F)F)c3cc21 |
|---|---|

**SUMMARY RESULTS**

| | |
|---|---|
| Molecular Weight | 309.287 |
| LogP | 4.2019 |
| Lipinski H-Bond Acceptors | 3 |
| Lipinski H-Bond Donors | 1 |
| Number of atoms | 36 |
| Number of heavy atoms | 22 |
| Number of rotatable bonds | 1 |
| Number of Heteroatoms | 6 |
| Number of Rings | 3 |
| Topological Polar Surface Area | 42.24 |

EMBL-EBI

- Similarity search RESTful web service (also for substructure)

## SIMILARITY SEARCH WEB SERVICE

**MENU**

- Home
- Tutorial
- Substructure and Exact Search ›
- Similarity Search ›
- Molecular Properties ›

**EXPLANATION**

In this section you can find the documentation to run a simple URI query to retrieve similarity searches, selecting different class of Fingerprints, and select between the Tanimoto and Dice coefficients. The input formats can be in SMILES strings or SMARTS queries

**NOTE:** All the categories are case-sensitive

### CATEGORIES

*Mandatory:*
- smiles : You can paste here an SMILES string
- smarts : You can paste here an SMARTS string

**(Warning: The user must be aware about the special URI characters)**

*Optional:*
- fingerprint : You can select between *'Morgan (ECFP-Like)'*, *'Morgan (FCFP-Like)'*, *'Torsion'*, or *'Atom-Pair'* fingerprints
- method : You can select between *'Tanimoto'* or *'Dice'* similarity coefficients

**Example 1 (SMILES):**

http://localhost/rest/api_chembl.php?action=similarity&smiles=C(=O)C1=CCnCC1&fingerprint=FCFP&method=Dice

**Example 2 (SMARTS):**
NOTE: The character "#" was replaced by the string "%23"

http://localhost/rest/api_chembl.php?action=similarity&smarts=[%236;X4]-1-[%236](=[%238])-[%237]-[%237]-[%236]-1=[%238]&fingerprint=ECFP&method=Tanimoto

**Check the PYTHON Client**

EMBL-EBI

## Check the PYTHON Client

```python
#! /usr/bin/env python

import urllib2
import urllib
import json


########################################################################
"""Functions"""
########################################################################
def translateURI(query):
    quoted_url = urllib.quote(query) # change the characters with trouble
    return quoted_url


########################################################################
"""Main"""
########################################################################

# Options for the query
smiles='C(=O)C1=CCnCC1'
smarts='[#6;X4]-1-[#6](=[#8])-[#7]-[#7]-[#6]-1=[#8]'
fingerprint = 'ECFP'
method = 'Tanimoto'

# Changing the query to an URI format
smiles=translateURI(smiles)
smarts=translateURI(smarts)


########################################################################
"""1. Example using SMILES"""
########################################################################
print "Results from SMILES ..."

# Storing the json file with the results
similarity_data = json.loads(urllib2.urlopen("http://10.7.248.227/rest/api_chembl.php?\
action=similarity&smiles=%s&fingerprint=%s&method=%s" % (smiles,fingerprint,method)).read())

# Printing the records
for record in similarity_data:
    print "ChemblID: %s" % record['ChEMBL_ID']
    print "Molregno: %s" % record['Molregno']
    print "Similarity: %s" % record['Similarity']
    print
```

EMBL-EBI

# Summary and future plans

• Provide a complete PostgreSQL version of ChEMBL per new release, available in the FTP site

    https://www.ebi.ac.uk/chembldb/index.php/downloads


• Integrate the Open Chemical Cartridge (RDKit) with open chemogenomics data (ChEMBL) in a Virtual Machine available for users, who can personalize the configuration and update the data.


• Provide a personal interface, with web applications and RESTful web services.

# Acknowledgements

ChEMBL

Training Opportunities

Mark Davies
Francis Atkinson
George Papadatos
John Overington

GGA
GGA Software Services

accelrys®
Symyx

RDKit
Open-Source Cheminformatics
and Machine Learning

KNIME

EMBL-EBI