

# Open-source from/in the enterprise: the RDKit

Gregory Landrum

NIBR Informatics

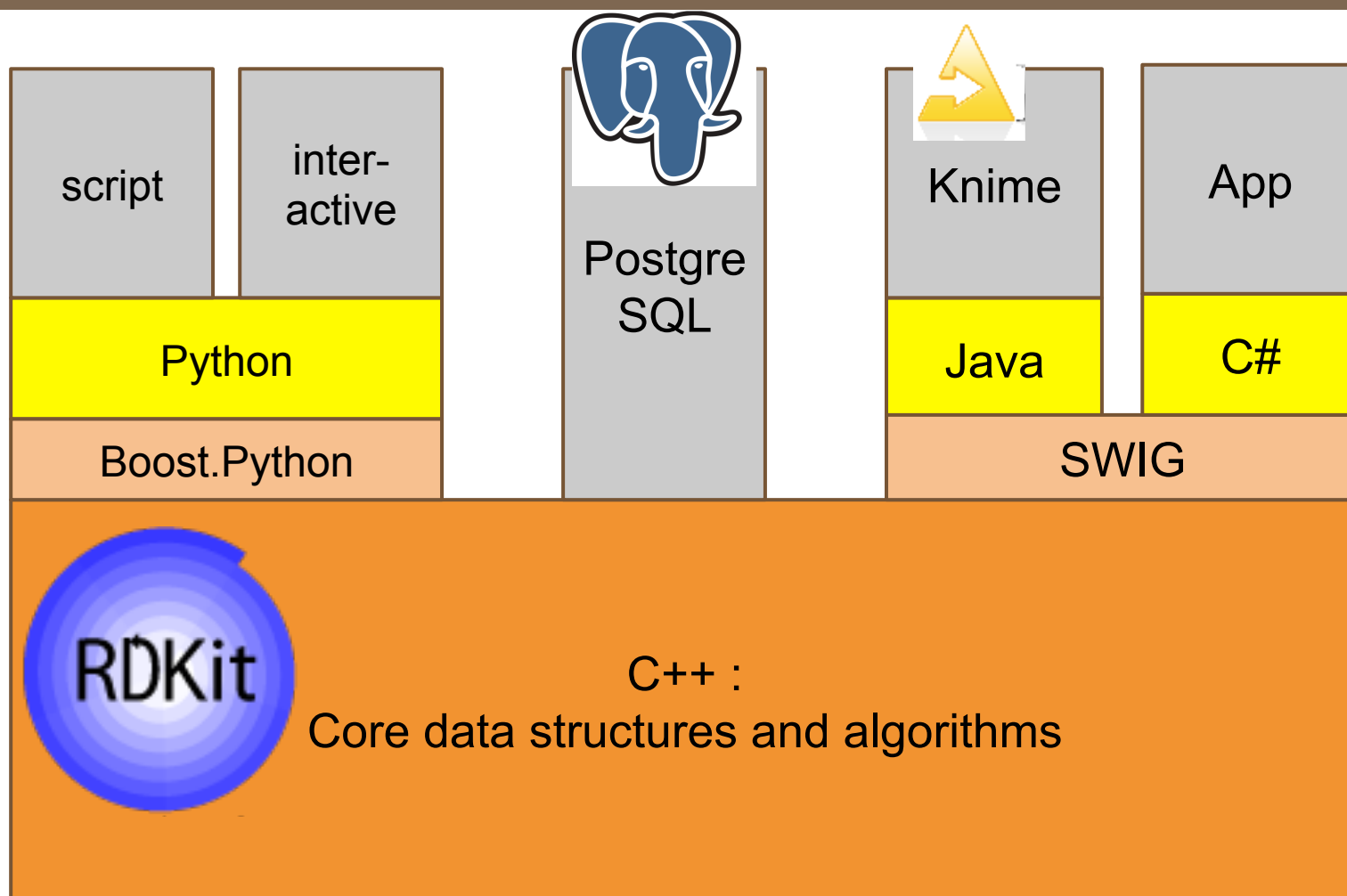
Novartis Institutes for BioMedical Research, Basel, Switzerland

# Outline

---

- RDKit integration with other open-source projects
  - Knime
  - PostgreSQL
  - IPython
  - Pandas
  - Lucene
- RDKit in NIBR, some case studies

# What is this all about?



*Exact same algorithms/implementations accessible from many different endpoints*

# Knime integration

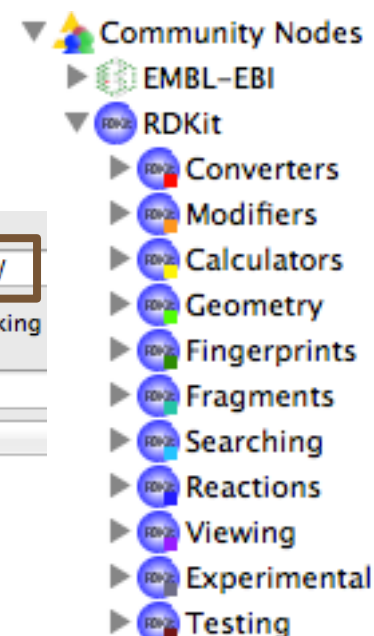


- Open-source RDKit-based nodes for Knime providing cheminformatics functionality
- Trusted nodes distributed from knime community site

Work with: Trusted Community Contributions – <http://tech.knime.org/update/community-contributions/trusted/2.9/> Find more software by working

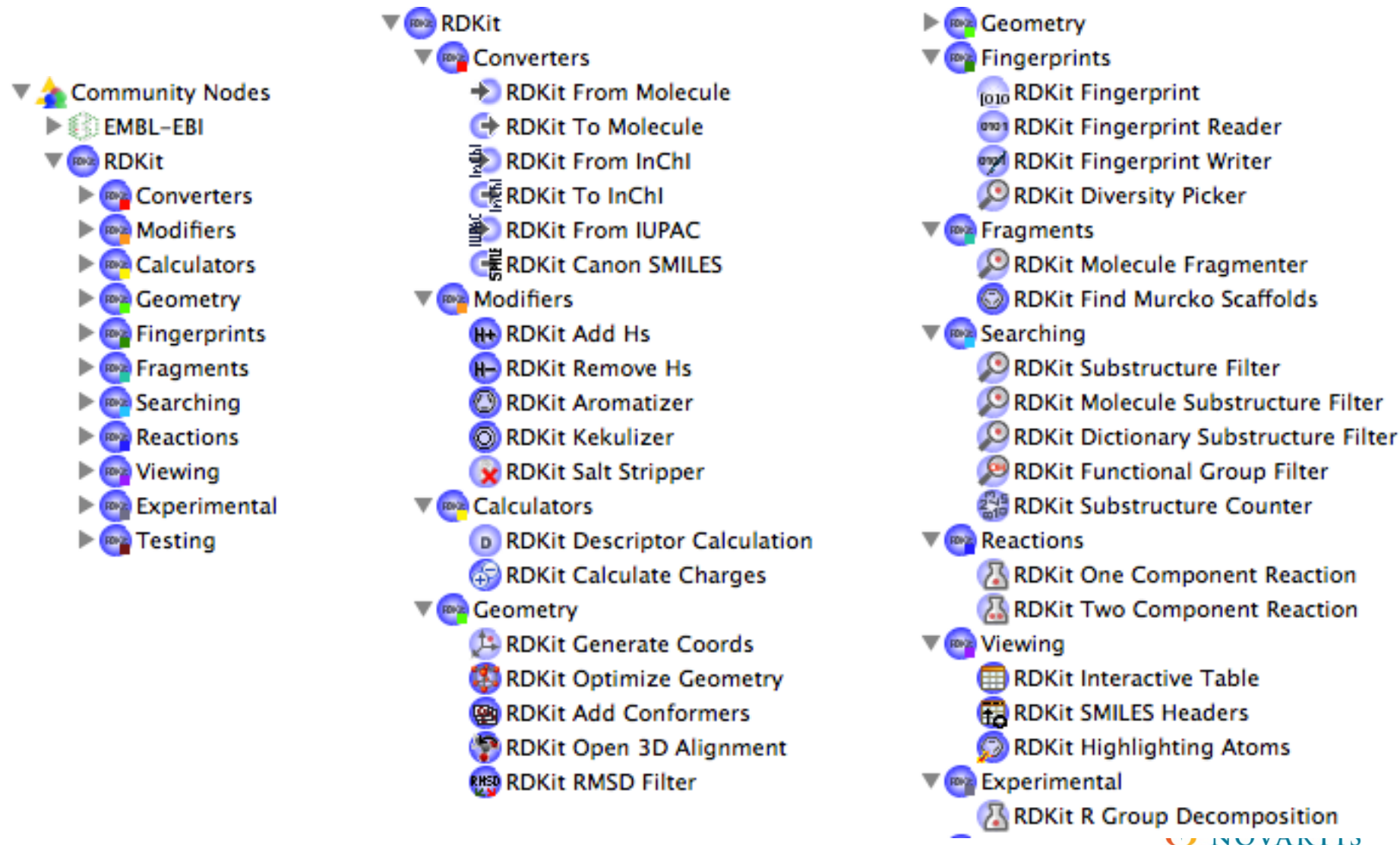
type filter text

Name	Version
<input type="checkbox"/> ▶ KNIME Community Contributions – Bioinformatics & NGS	
<input type="checkbox"/> ▼ KNIME Community Contributions – Cheminformatics	
<input type="checkbox"/> CheS–Mapper–Node extension for KNIME Workbench	1.0.0.201307180737
<input type="checkbox"/> CIR KNIME integration	1.0.0.201312051437
<input type="checkbox"/> EMBL–EBI Nodes for KNIME	1.0.3.201312121642
<input type="checkbox"/> KNIME CDK Integration (based on CDK 1.5.2)	1.4.4.201310241550
<input type="checkbox"/> RDKit KNIME integration ←	2.3.0.201401281545
<input type="checkbox"/> RDKit KNIME Wizards	2.3.0.201310111300
<input type="checkbox"/> Vernalis Knime Nodes	1.0.4.201401201138



- Work in progress: more nodes being added (new wizard makes it easy)

# What's there?



# RDKit Interactive Table



- KNIME interactive table with molecules as column headers

Diagram showing the workflow:

```

graph LR
    In[ ] --> Node34[Node 34: RDKit Substructure Counter]
    In --> Node35[Node 35: RDKit Interactive Table]
    Node34 --> Node35
  
```

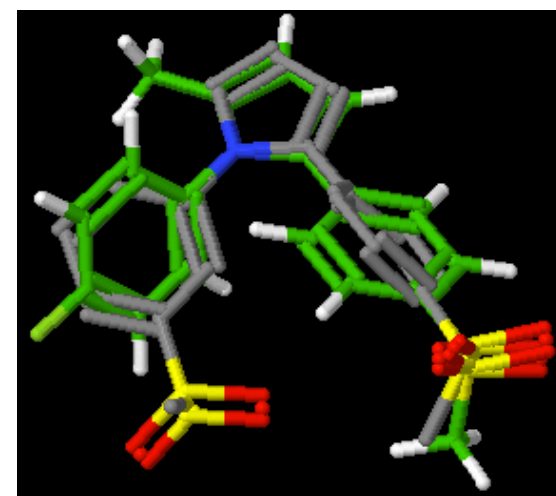
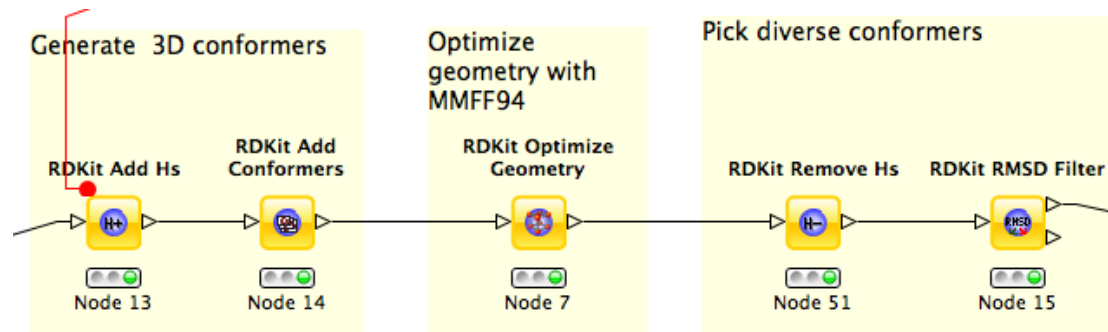
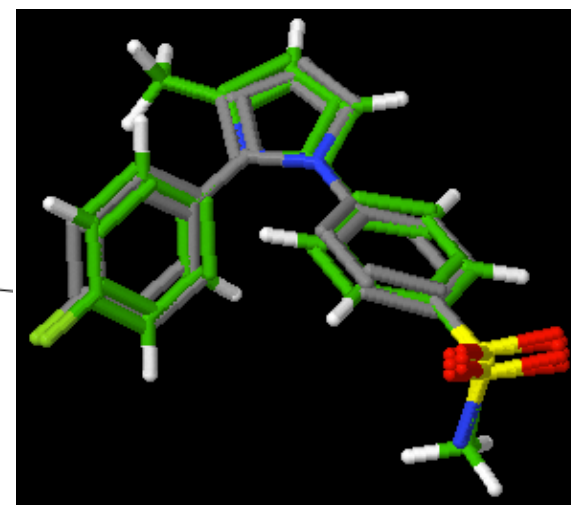
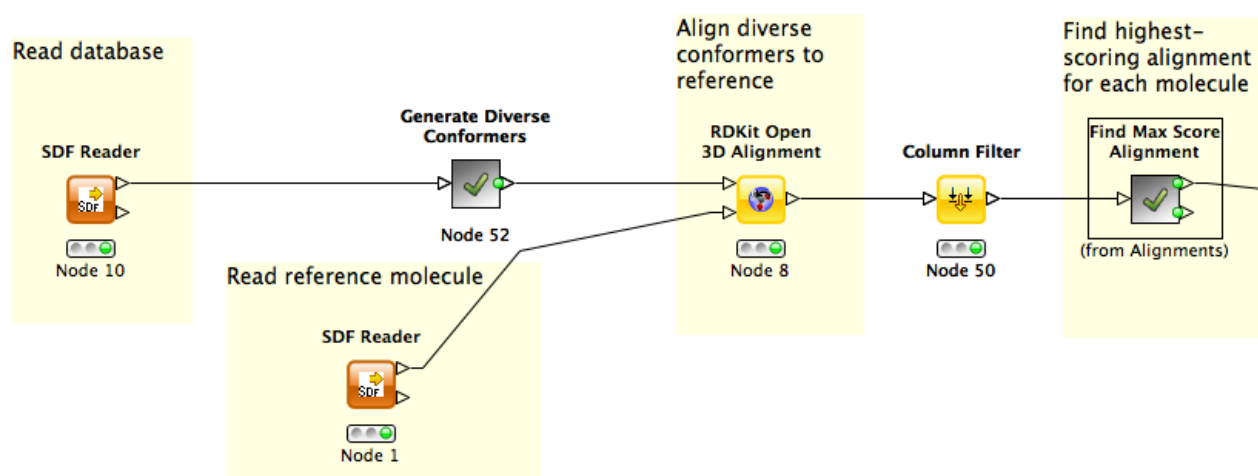
Below the diagram is a screenshot of the **RDKit Interactive Table View** window (10000 x 482). The window displays a table with the following columns:

Row ID	smiles	wehi-id	<chem>c1:c:c(</chem>	<chem>c1(:[c&amp;!H</chem>	<chem>c1(-,:n:,</chem>	<chem>[c&amp;!H0]2</chem>	#
Row294		WEHI-00947... 0	1	0	0	0	
Row1060		WEHI-00970... 0	1	0	0	0	



# Functionality for working with 3D molecules

- Example: flexible molecule-molecule alignment





## PostgreSQL integration

---

- PostgreSQL (<http://www.postgresql.org>): a robust, flexible, and extensible relational open-source database. Rich collection of extensions available
- RDKit “cartridge”:
  - Fast substructure and similarity search
  - Fingerprints (count-based and bit-vector):  
Morgan (ECFP-like), FeatMorgan (FCFP-like), RDKit (Daylight like), atom pair, topological torsion, MACCS
  - Standard molecule properties and descriptors
- Basis for myChEMBL (<http://chembl.blogspot.co.uk/2013/10/chembl-virtual-machine-aka-mychembl.html>) Ochoa, R., Davies, M., Papadatos, G., Atkinson, F., & Overington, J. P. (2014). myChEMBL: a virtual machine implementation of open data and cheminformatics tools. *Bioinformatics*, 30(2), 298–300.



# PostgreSQL integration

## Substructure search



```
chembl_17=# select molregno,m from rdk.mols where  
m@>'c1ccc2c(c1)C(=NN(C2=O)Cc3nc4cc(ccc4s3)C)CC(=O)O';
```

molregno	m
7502	O=C(O)Cc1nn(Cc2nc3cc(C(F)(F)F)ccc3s2)c(=O)c2cccc12
23364	O=C(O)Cc1nn(Cc2nc3cc(C(F)(F)F)cc(C(F)(F)F)c3s2)c(=O)c2cccc12
23439	O=C(O)Cc1nn(Cc2nc3cc(C(F)(F)F)cc(Cl)c3s2)c(=O)c2cccc12
23462	O=C(O)Cc1nn(Cc2nc3cc(C(F)(F)F)cc(F)c3s2)c(=O)c2cccc12
24192	Cc1cc2nc(Cn3nc(CC(=O)O)c4cccc4c3=O)sc2c(C)c1
24190	COc1cc2sc(Cn3nc(CC(=O)O)c4cccc4c3=O)nc2cc1C(F)(F)F
24194	Cc1ccc2sc(Cn3nc(CC(=O)O)c4cccc4c3=O)nc2c1
24237	O=C(O)Cc1nn(Cc2nc3cc(C(F)(F)F)c(O)cc3s2)c(=O)c2cccc12
24331	CC(c1nc2cc(C(F)(F)F)ccc2s1)n1nc(CC(=O)O)c2cccc2c1=O

(9 rows)

Time: 112.325 ms

# PostgreSQL integration

## Similarity search



```
chembl_17=# select * from get_mfp2_neighbors('O=C(O)Cc1nn(Cc2nc3cc(C(F)
(F)F)ccc3s2)c(=O)c2cccc12') limit 5;
```

molregno	m	similarity
7502	<chem>O=C(O)Cc1nn(Cc2nc3cc(C(F)(F)F)ccc3s2)c(=O)c2cccc12</chem>	1
24184	<chem>O=C(O)Cc1nn(Cc2nc3ccc(C(F)(F)F)cc3s2)c(=O)c2cccc12</chem>	0.859649122807018
24153	<chem>O=C(O)Cc1nn(CCc2nc3cc(C(F)(F)F)ccc3s2)c(=O)c2cccc12</chem>	0.830508474576271
24152	<chem>O=C(O)Cc1nn(Cc2nc3cccc3s2)c(=O)c2cc(C(F)(F)F)ccc12</chem>	0.813559322033898
24150	<chem>O=C(O)Cc1nn(Cc2nc3cccc3s2)c(=O)c2ccc(C(F)(F)F)cc12</chem>	0.813559322033898

(5 rows)

Time: 1222.426 ms

Notice that results come back in sorted order

# PostgreSQL integration

## Other functionality



```
chembl_17=# select mol_formula('O=C(O)Cc1nn(Cc2nc3cc(C(F)(F)F)ccc3s2)c(=O)c2cccc12');
mol_formula
```

```
-----
C19H12F3N3O3S
```

```
(1 row)
```

```
chembl_17=# select mol_logp('O=C(O)Cc1nn(Cc2nc3cc(C(F)(F)F)ccc3s2)c(=O)c2cccc12');
mol_logp
```

```
-----
3.7004
```

```
(1 row)
```

```
chembl_17=# select mol_inchi('O=C(O)Cc1nn(Cc2nc3cc(C(F)(F)F)ccc3s2)c(=O)c2cccc12');
mol_inchi
```

```
-----
-----
InChI=1S/C19H12F3N3O3S/
c20-19(21,22)10-5-6-15-14(7-10)23-16(29-15)9-25-18(28)12-4-2-1-3-11(12)13(24-25)8-17(26)27
/h1-7H,8-9H2,(H,26,27)
```

```
(1 row)
```

# PostgreSQL integration

## Other functionality

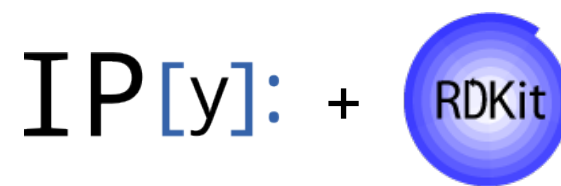


```
chembl_17=# select mol_to_ctab('CC'::mol);
               mol_to_ctab
```

```
-----+
      RDKit      2D      +
      +          +      +
      +          +      +
      2  1  0  0  0  0  0  0  0  0999 V2000      +
      0.0000  0.0000  0.0000 C  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0+
      1.2990  0.7500  0.0000 C  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0+
      1  2  1  0      +
M  END      +

(1 row)
```

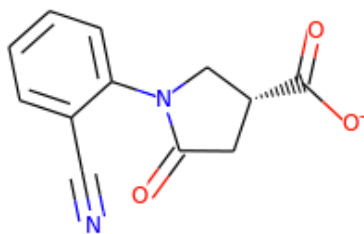
# IPython notebook integration



- IPython: a very powerful interactive shell for python  
<http://www.ipython.org>
- IPython notebook: IPython in the browser, with graphics
  - combines code and output in one place
  - great tool for reproducible research
  - [Example notebook with graphics.](#)
- RDKit integration:
  - Display molecules, substructure matches, reactions, graphics from PyMOL

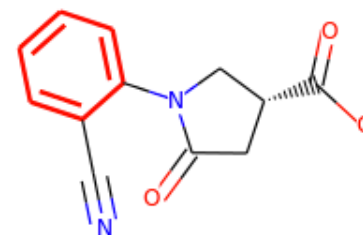
```
In [3]: Chem.SetHybridization(m)  
m
```

Out[3]:

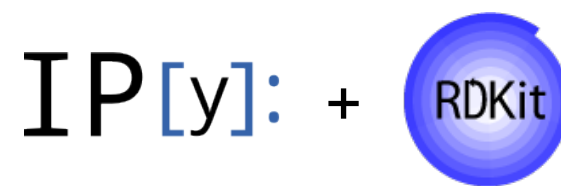


```
In [4]: tmp=m.GetSubstructMatch(Chem.MolFromSmarts('c1ccccc1'))  
m
```

Out[4]:

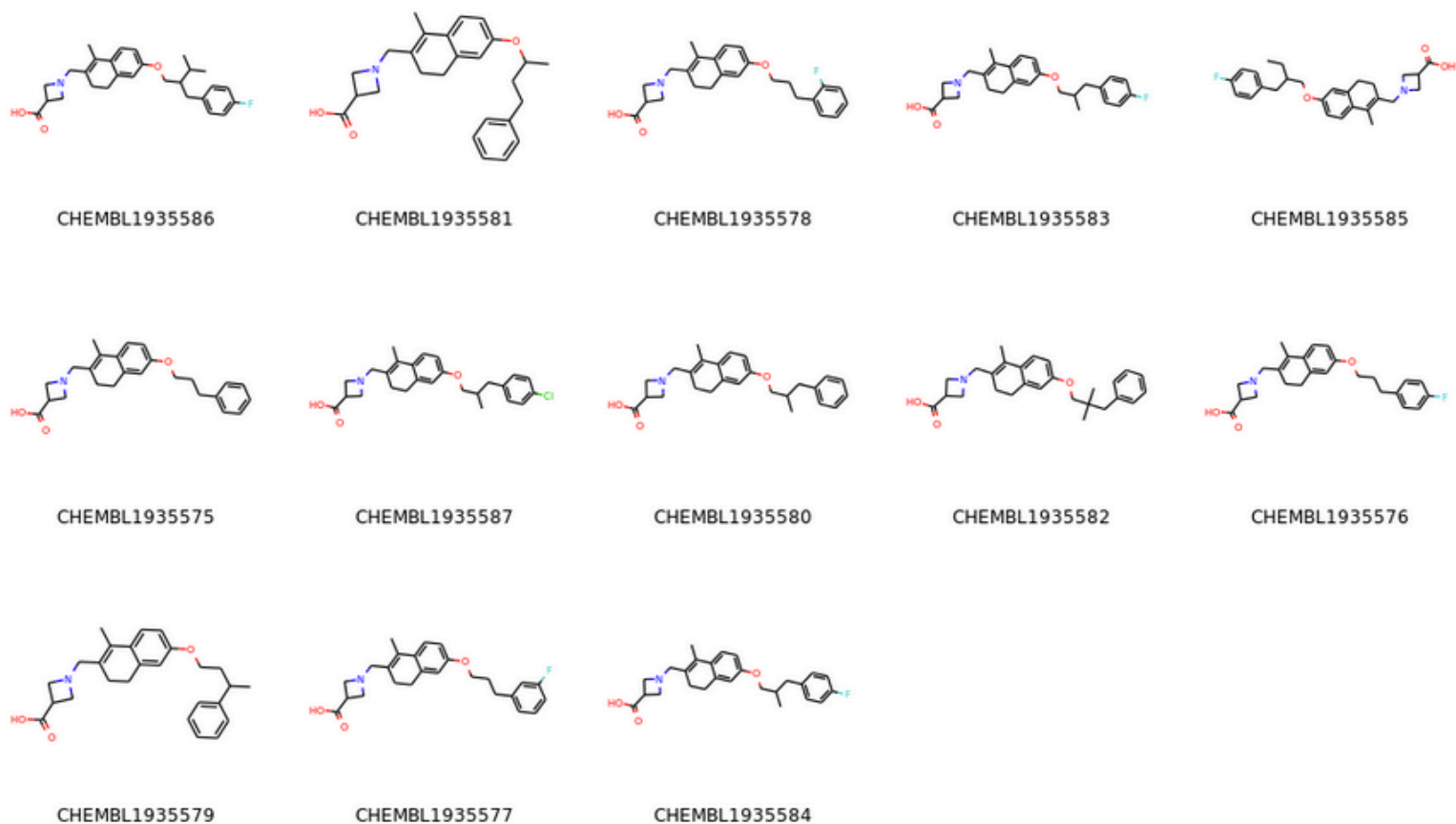


# IPython notebook integration: Molecule tables



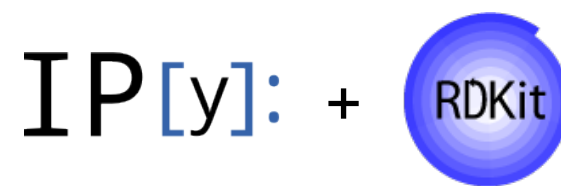
```
In [8]: ids,smis= zip(*cmpds)
mols = [Chem.MolFromSmiles(x) for x in smis]
Draw.MolsToGridImage(mols,molsPerRow=5,legends=ids)
```

Out[8]:



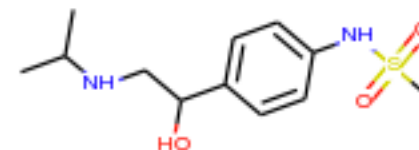
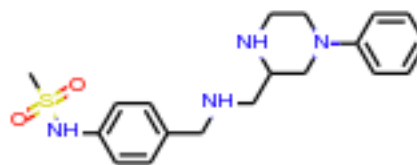
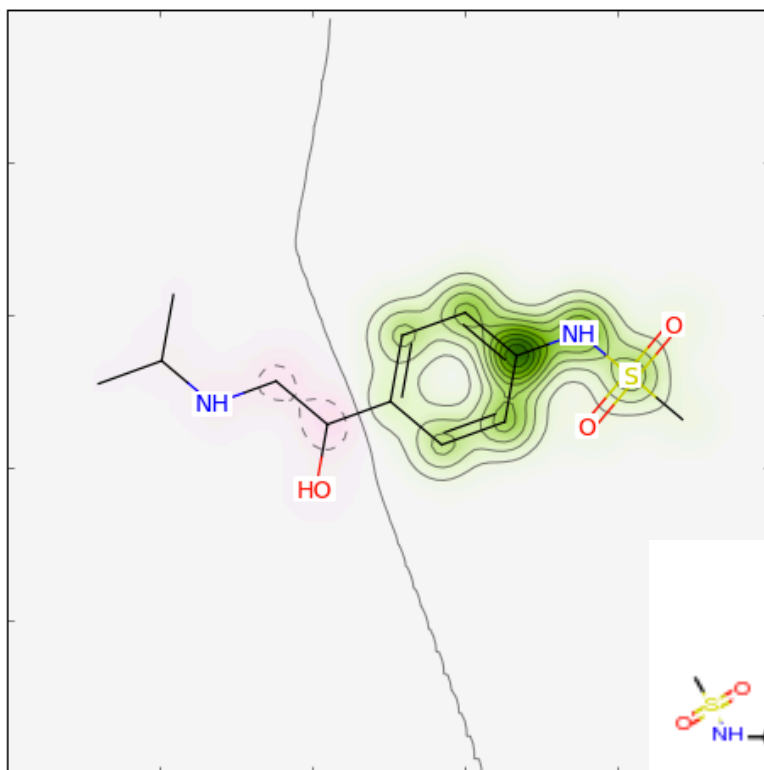
<http://rdkit.blogspot.ch/2014/02/more-on-datasets-ii.html>

# IPython notebook integration: Similarity Maps



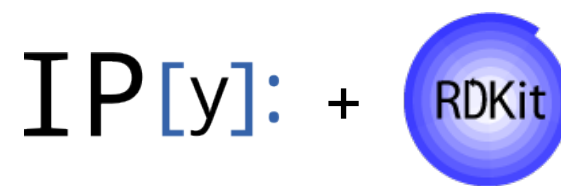
```
In [56]: SimilarityMaps.GetSimilarityMapForFingerprint(ms[0],ms[16],SimilarityMaps.GetTTFingerprint)
```

```
Out[56]: (<matplotlib.figure.Figure at 0x109786850>, 0.31663113006396593)
```





# IPython notebook integration: PyMol

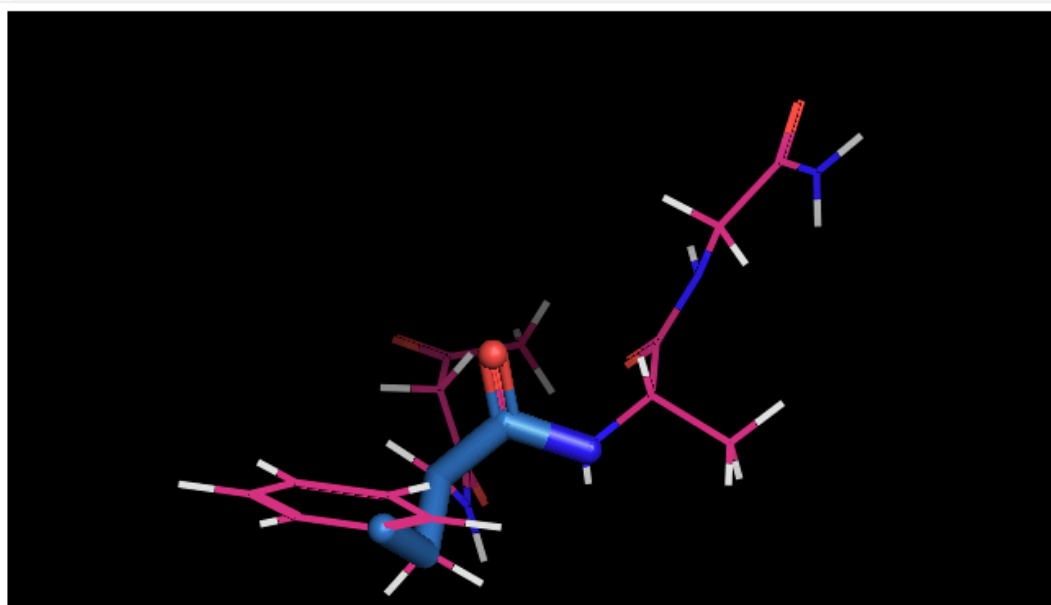


```
In [12]: nm = Chem.Mol(matchingMols[3])
GetFF=lambda x,confId=-1:AllChem.MMFFGetMoleculeForceField(x,AllChem.MMFFGetMoleculeProperties(x),confId=confId)
AllChem.ConstrainedEmbed(nm,core1,getForceField=GetFF)
rms = float(nm.GetProp('EmbedRMS'))
print 'RMS:',rms

v.ShowMol(core1,molB=Chem.MolToMolBlock(core1,kekulize=False),name='core')
v.SetDisplayStyle('core','sticks')
v.ShowMol(nm,'constrained embed',showOnly=False)
v.GetPNG(preDelay=2)
```

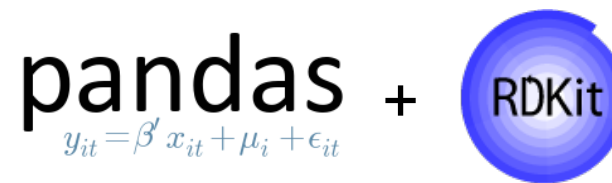
RMS: 0.0407469103071

Out[12]:



<http://rdkit.blogspot.ch/2013/12/using-allchemconstrainedembed.html>

# Pandas integration



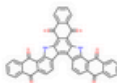
- Pandas: library for working with data tables in Python. Integrates well with matplotlib and ipython  
<http://pandas.pydata.org/>
- RDKit integration:
  - Load smiles tables or SD files into Pandas data tables
  - Adds molecule columns to existing tables with smiles/SD columns
  - Enables substructure filters on tables
  - Integration with IPython notebook to render molecules

# Pandas integration

## Molecules in tables

```
In [6]: PandasTools.AddMoleculeColumnToFrame(data, smilesCol='smiles', molCol='molecule', includeFingerprints=False)
data.head(2)
```

Out[6]:

	smiles	mutagenic	molecule
cas			
2475-33-4	<chem>O=C1c2ccccc2C(=O)c3c1ccc4c3[nH]c5c6C(=O)c7ccccc7C(=O)c6c8[nH]c9c%10C(=O)c%11ccccc%11C(=O)c%10ccc9c8c45</chem>	0	
820-75-7	<chem>NNC(=O)CNC(=O)\C=N\#N</chem>	1	None

## Substructure filters

```
In [11]: data.groupby(data['molecule'] >= nitroso).describe().unstack()
```

Out[11]:

	mutagenic							
	count	mean	std	min	25%	50%	75%	max
molecule								
False	5217	0.461760	0.498583	0	0	0	1	1
True	1233	0.838605	0.368044	0	1	1	1	1

# Lucene integration



- Still in the experimental stage
- Adds substructure search functionality with fingerprint screenout to Lucene
- Includes demo app for testing

Search

Enter Search Term(s):

Search Results

Execution of Search By Substructure: c1ncnc1

ZINC39279791

ZINC00141286

ZINC12504456

ZINC03953815

ZINC03873955

ZINC03873956

ZINC03873957

ZINC03873958

ZINC02524720

ZINC02046906

ZINC02046907

ZINC02047002

ZINC02169830

ZINC00000842

ZINC00057125

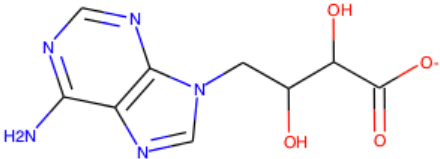
ZINC00598852

ZINC04514079

ZINC04514080

Result Details

Nc1ncnc2c1N=CN2CC(O)C(O)C([O-])=O

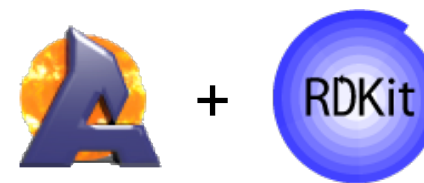


363 Hits in 273 ms.

## RDKit in NIBR

---

- Extensive use by CADD, informaticians, and IT
- Lots of convenience code/wrappers for accessing internal data sources and tools
- Combined with the Avalon toolkit (another NIBR-supported open-source project), provides the underpinning for many of our global chemistry-based applications



# The Avalon toolkit

---

- C/Java cheminformatics toolkit
- Primary author: Bernd Rohde (NIBR Informatics Basel)
- <http://sourceforge.net/projects/avalontoolkit/>
- Functionality:
  - Canonical SMILES
  - Avalon fingerprint (highly optimized substructure fingerprint)
  - Molecular standardization (STRUCHK)
  - 2D Coordinate generation
  - Tomcat webapp for 2D rendering
- The RDKit has (optional) Python bindings for much of the functionality

# RDKit in NIBR

## Case study 1: Clx Framework



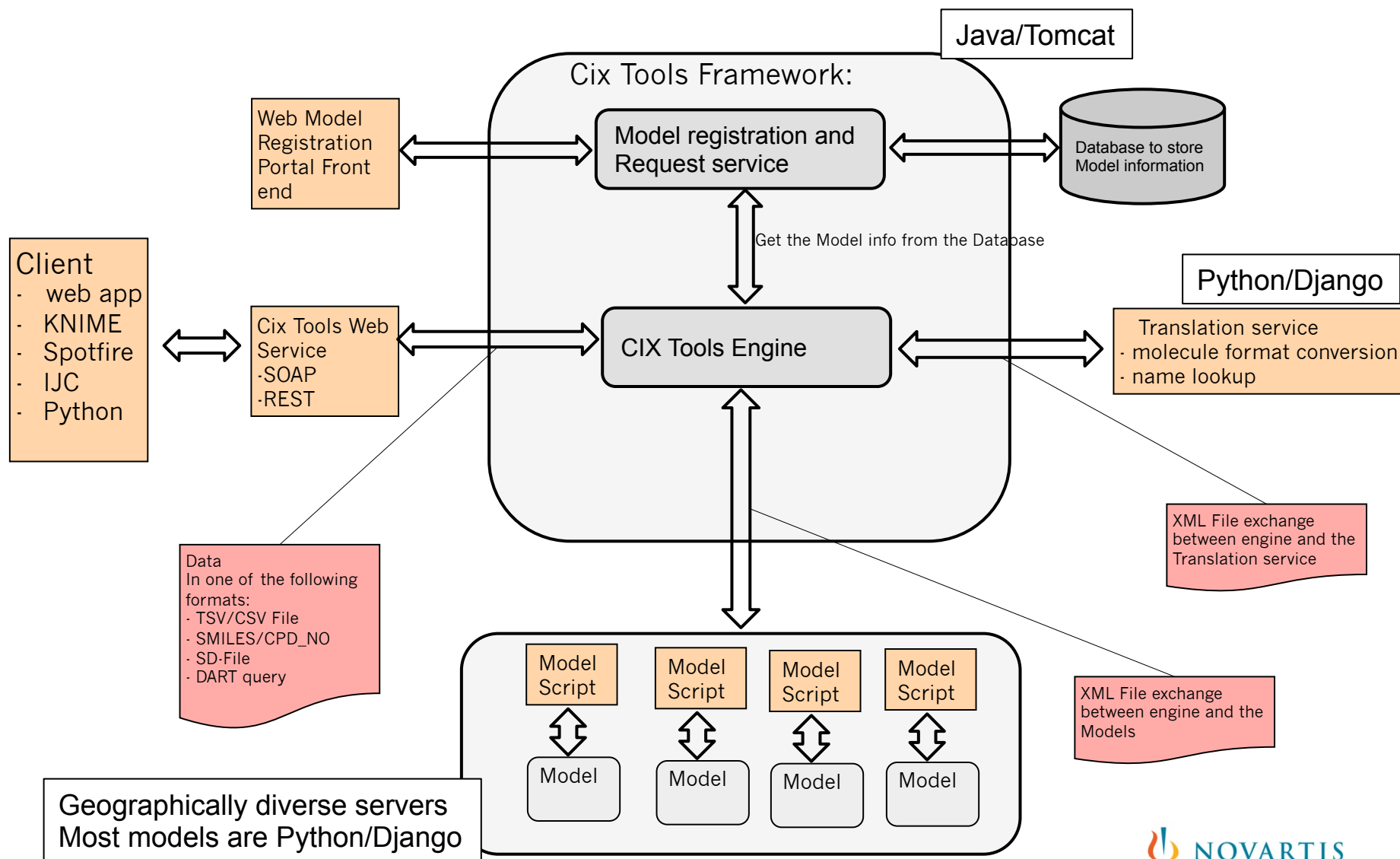
+



- “Service bus” for cheminformatics/CADD services
- Handles format conversions for input/output automatically  
i.e. callers can provide SMILES input to a service/model wants CTABs with 3D coordinates
- Supports versioning of models/services
- Tight integration with scientific tools (e.g. Tibco Spotfire, Knime, Instant JChem, etc.)
- Enables trivial addition of “chemical intelligence” to web apps
- Makes it easy to globally deploy models: once a new model/service (or new version of a model/service) is registered with the Framework, it is instantly globally accessible



# Cix Framework architecture



# RDKit in NIBR

## Case study 2: Small-Molecule Registration

---



+



- Internally developed web application for compound registration
- C#-based web services writing to Oracle
- RDKit + Avalon toolkit for structure standardization
- RDKit + InChI used for structure-key calculation
- Calls out to Clx Framework for standard computed properties
- Independent (but validated) Python implementation of standardization and structure-key calculation for standalone use

# RDKit in NIBR

## Case study 3: QSAR Toolkit

---



+



- Descriptor calculator providing access to all available internal descriptors
- Tools for pulling assay data from our data warehouse
- Standardized model-building
- Standardized reporting for evaluation and peer review
- Packaging for deployment via Clx Framework
- Model Watchdog:  
Pulls most recent data, generates predictions, creates report showing evolution of model accuracy over time

# RDKit in NIBR

## Case study 4: Similarity Server

---



+



- Central PostgreSQL database with easily available compounds
  - in-house available
  - available from reliable vendors
- Kept up-to-date
- Substructure search
- Similarity search with various fingerprints:
  - Avalon
  - Morgan2, Morgan3, FeatMorgan2
  - Atom Pairs, Topological Torsions
- Web services interface
- Available to chemists via one of their standard desktop tools
  
- Currently deploying a new version based on chemfp

# RDKit in NIBR

## Case study 5: Compound Series and Favorites

---



+



- Central system to store definitions of the chemical series that project teams are working on
- Captures annotations and relationships
- Searchable, including Markush search (using ChemAxon cartridge)
- Accessible via web services
- RDKit used for scaffold validation

# RDKit in NIBR

## Case study 5: Compound Series and Favorites



### Compound Series & Favorites BETA

Gregory Landrum

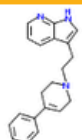
[Back](#) | [Home](#) / D4 Receptor (Public Data)

#### D4 Receptor (Public Data)



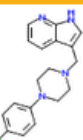
##### 7-Azaindoles

###### COMPOUND 2C



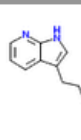
Target Selectivity **l**

###### L-745-870

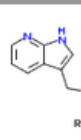


Start

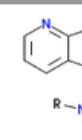
###### SCAFFOLD 4



###### SCAFFOLD 2



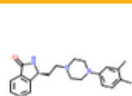
###### SCAFFOLD 3



Scaffold

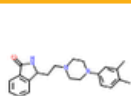
##### Isoindolinones

###### PD 172938



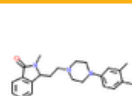
In vivo Active **s**

###### COMPOUND 10A

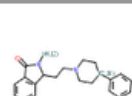


Target Selectivity **h**

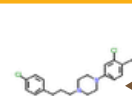
###### COMPOUND 14A



###### ISOINDOLINONE



###### COMPOUND 3

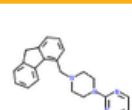


Start

Compound

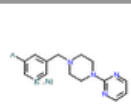
##### Biaryl methylamines

###### COMPOUND 8E

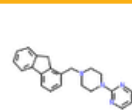


Target Selectivity **l**

###### SCAFFOLD 1

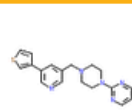


###### COMPOUND 8C



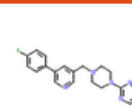
In vivo Active **s**  
Target Selectivity **h**

###### COMPOUND 5F



In vivo Active **s**  
Target Selectivity **h**

###### COMPOUND 3A



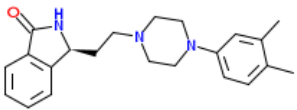
Start






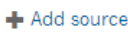
# RDKit in NIBR

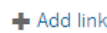
## Case study 5: Compound Series and Favorites

PD 172938 





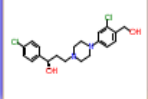
**Tags**    
In vivo Active **strong**

**Source info**    
Literature **1**

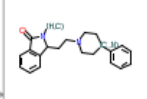
**Links** 

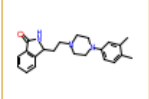
**Last modified** Anna Pelliccioli 21 Aug 2014 **Added** Anna Pelliccioli 21 Aug 2014 **Project** D4 Receptor (Public Data)  
**Series** Isoindolinones

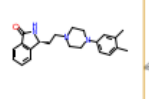
**Relationships**    
Derived from **1** Led to **0** **Path**

**COMPOUND 3**  



bioisosteric replacement

**ISOINDOLINONE**  


**COMPOUND 10A**  


**PD 172938**  


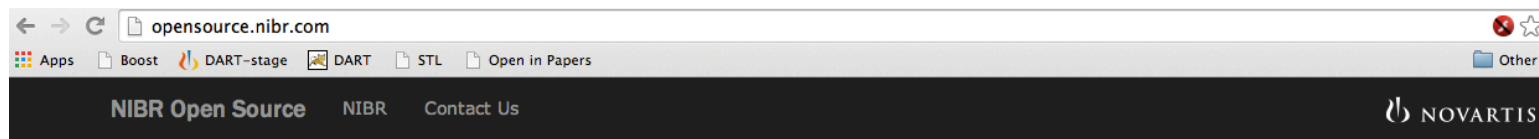
single enantiomer

**Notes** 



# NIBR Open Source

## Something new



## Open Source at NIBR

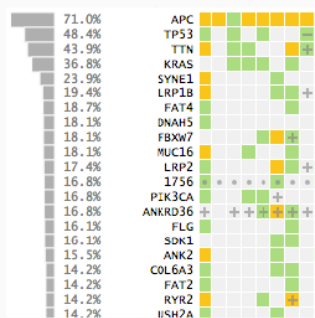
The Novartis Institutes for BioMedical Research (NIBR) is pioneering new informatics tools for drug discovery. We believe in the power of open-sourced, global collaboration for the greater good. Join us to help patients worldwide. [Read about the work we do.](#)

### Interested in working for NIBR Engineering?

At NIBR, you'll be at the forefront of technology – helping to shape it, develop it, and make it impactful. Partnering with scientists, our engineers create cutting-edge, state-of-the-art solutions that accelerate drug discovery and ultimately improve patients' lives.

[Learn More »](#)

### GridVar

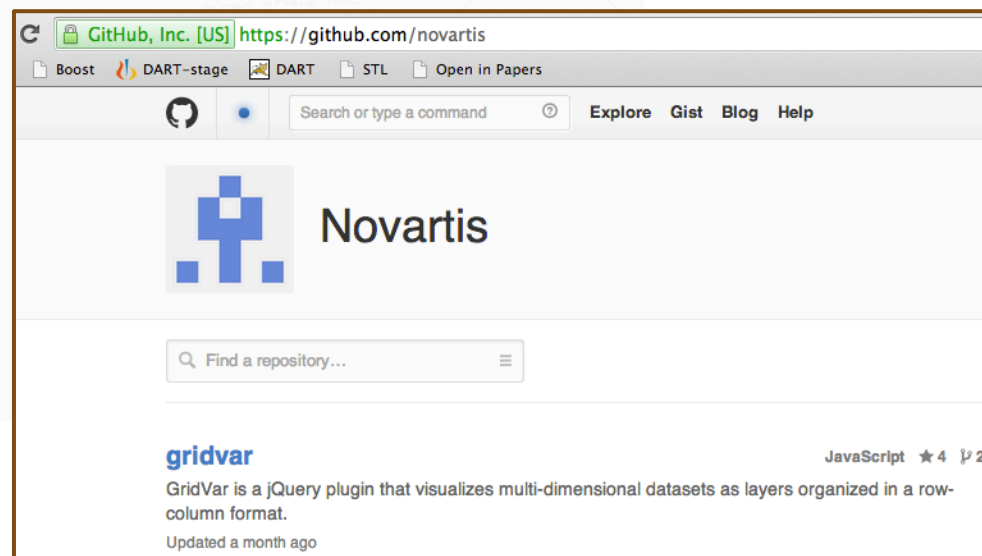


GridVar is a jQuery plugin that visualizes multi-dimensional datasets

### RDKit



A collection of cheminformatics and machine-learning software written in





# Acknowledgements

---

- General:
  - Remy Evard (NIBR/Informatics)
  - Richard Lewis (NIBR/GDC)
  - Tom Digby (NIBR/Legal)
  - Peter Gedeck (NIBR/GDC)
  - Nik Stiefl (NIBR/GDC)
- RDKit Community
  - Roger Sayle (NextMove): PDB Parser
  - Andrew Dalke (Dalke Scientific): FMCS
  - Paolo Tosco (University of Turin, now Cresset): MMFF94, Open3DAlign
  - Jameed Hussain (GSK, now CCG): Fraggie, mmpa
- Pandas, scikit-learn:
  - Sereina Riniker (NIBR/Informatics, now ETH)
  - Nikolas Fechner (NIBR/Informatics)
- Knime:
  - Manuel Schwarze (NIBR/Informatics)
  - Thorsten Meinl (knime.com)
  - Bernd Wiswedel (knime.com)
- SMR
  - Thomas Mueller (NIBR/Informatics)
  - Thomas Veith (NIBR/Informatics)
  - Dave Cotter (NIBR/Informatics)
- QSAR Toolkit:
  - Peter Gedeck (NIBR/GDC)
  - Nikolas Fechner (NIBR/Informatics)
- Clx Framework
  - Sandra Mueller (NIBR/Informatics)
  - Joerg Muehlbacher (NIBR/CPC)
  - Riccardo Vianello (NIBR/Informatics)
- Compound Series & Favorites
  - Anna Pelliccioli (NIBR/Informatics)
  - Manuel Schwarze (NIBR/Informatics)
  - Recca Chatterjee (NIBR/Informatics)
  - Mikhail Rybalkin (NIBR/Informatics)
  - Roman Bolshev (NIBR/Informatics)
- NIBR Open Source
  - Ken Robbins (NIBR/Informatics)
  - Dennis Jen (NIBR/Informatics)
  - Mark Schreiber (NIBR/Informatics)