Open for Innovation ®

# KNIME

# Guided Analytics in Action: Patent analysis from your browser[1]

Greg Landrum

KNIME.com AG

[1] and some other stuff

# Agenda

- A few more words on "What's new" and "What's cooking"
  - KNIME in the cloud
  - KNIME on Spark


- A Guided Analytics use case: patent analysis
- Quick advertisement for a workshop tomorrow

Open for Innovation   ®
KNIME

# KNIME in the cloud

- KNIME Analytics platform on AWS/Azure

- KNIME Server, Server + Big Data on AWS/Azure

- Launch Hadoop/Spark clusters on AWS/Azure

- Connectors for S3, Blob Store, Azure SQL DB, Redshift


- What else should we be thinking about/working on?

Open for Innovation    ®
KNIME

# KNIME on Spark

Demo

# The case study: working with patent data

- Start with the PDF for a med-chem patent

- Extract the structures from the patent

- Attempt to identify the key compound

Inspired by:
http://chembl.blogspot.ch/2014/11/finding-key-compounds-in-med-chemistry.html
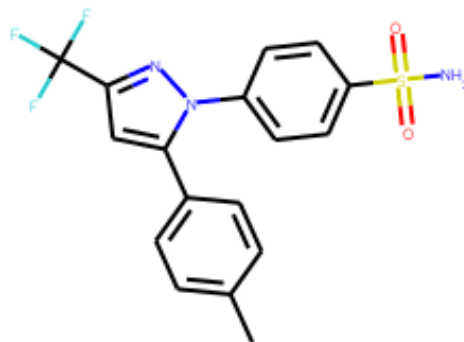
Open for Innovation    ®
KNIME

# The patent

The Front-page

## WO-1995015316-A1 / 1995-06-08

EN   SUBSTITUTED PYRAZOLYL BENZENESULFONAMIDES FOR THE TREATMENT OF INFLAMMATION

EN   ABSTRACT

A class of pyrazolyl benzenesulfonamide compounds is described for use in treating inflammation and inflammation-related disorders. Compounds of particular interest are defined by formula (II), whrein R2 is selected from hydrido, alkyl, haloalkyl, alkoxycarbonyl, cyano, cyanoa lkyl. carboxyl. aminoc arbonyl. alkyla mi nocar bonyl. cvcloalkvlaminocarb onvl. arvlamin oc arbonyl. carboxvalk vl aminocarbonvl.
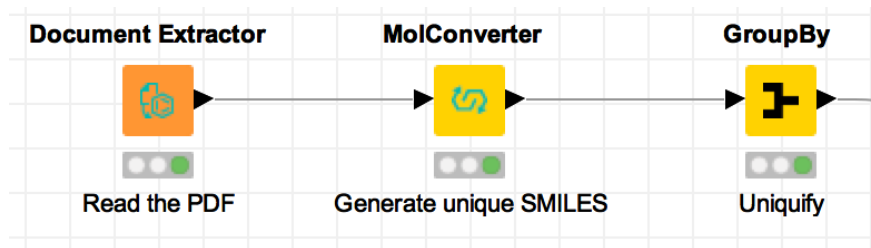
This is the patent "for Celebrex" :

One of the validation patents from:
Hattori, K., Wakabayashi, H. & Tamaki, K. Predicting Key Example Compounds in Competitors' Patent Applications Using Structural Information Alone. *J. Chem. Inf. Model.* **48,** 135–142 (2008).
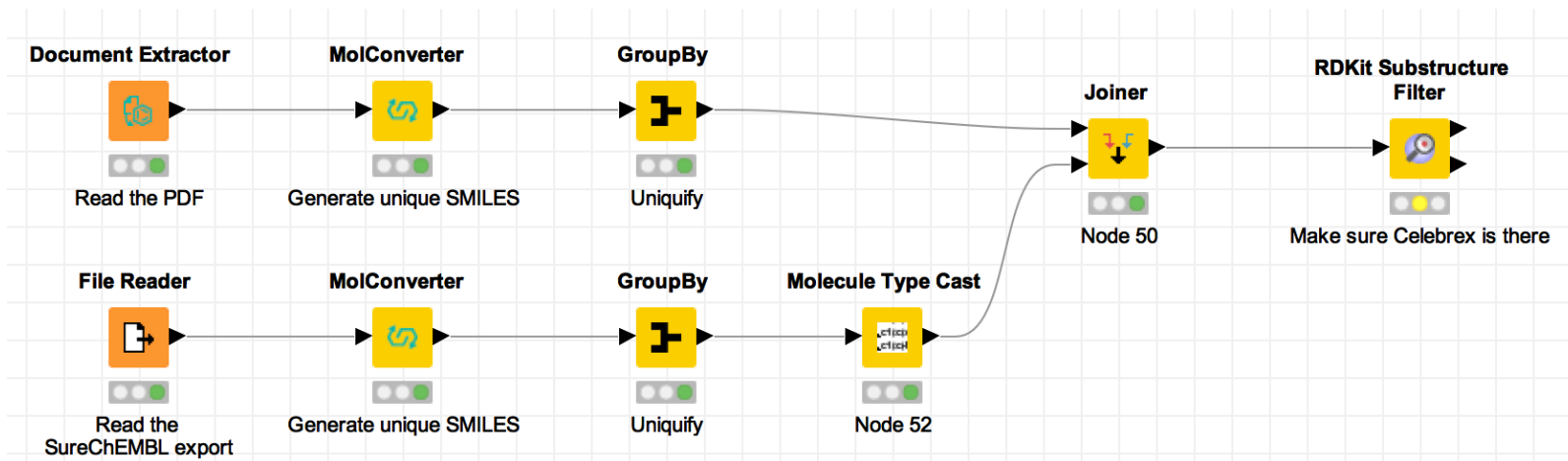
Open for Innovation ®
KNIME

# Extracting structures

Grab the PDF from SureChEMBL and the use the ChemAxon Document Extractor:

# Let's start with a bit of validation:

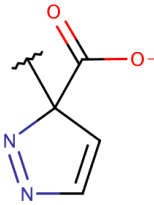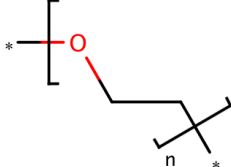Compare the Document Extractor to the structures downloadable from SureChEMBL:
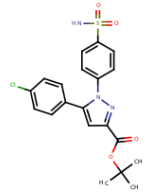


Results:

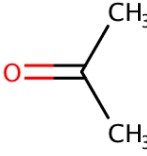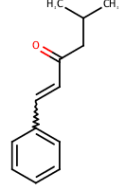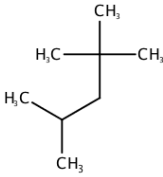- 451 structures from Document Extractor
- 548 structures from ChEMBL
- 192 are in common, and Celebrex is one of them.

# Cleaning up the data



The interesting structures are mixed in with a bunch of stuff we'd rather not see

# Cleaning up the data, but with which rules?

- Some things, like removing molecules with attachment points, are obvious. But what about the other properties?

- There are multiple different ways of filtering these down to the "interesting" ones.

- Instead of picking some hard-and-fast rule, let's do it interactively using a set of reasonable properties:
  - # of heavy atoms
  - # of rotatable bonds
  - # of rings

Open for Innovation ®
KNIME

# Finding the key compound in the patent

- Idea: the key compound is likely to have a lot of similar compounds

- Classic approach: find the compounds in the patent with the most near neighbors

- Alternate approach: use metrics from network analysis

Method from:
Hattori, K., Wakabayashi, H. & Tamaki, K. Predicting Key Example Compounds in Competitors' Patent Applications Using Structural Information Alone. *J. Chem. Inf. Model.* **48,** 135–142 (2008).
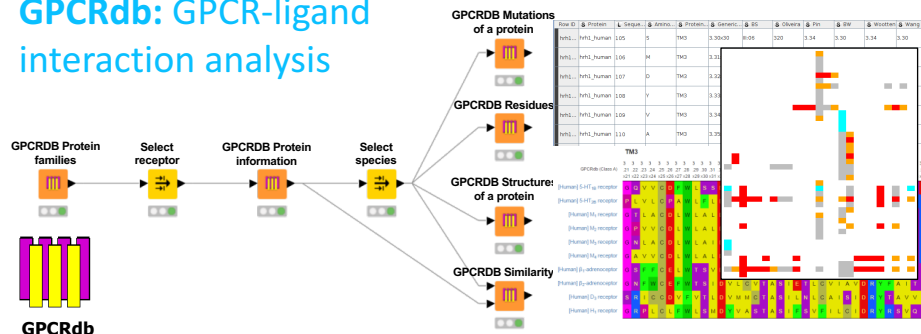
# Advertising!

Open for Innovation  ®

KNIME

# 3D-e-Chem KNIME Nodes for Integrated Structural Cheminformatics Analyses and Computer-Aided Drug Discovery

**Information:** https://tech.knime.org/3d-e-chem-nodes-for-knime

**Workshop:** https://github.com/3D-e-Chem/workflows (Clone or download - Download **ZIP**)

McGuire, Verhoeven, Vass, (...), De Graaf. *J Chem Info Model* 2017, 57: 115

The KNIME® trademark and logo and OPEN FOR INNOVATION® trademark are used by
KNIME.com AG under license from KNIME GmbH, and are registered in the United States.
KNIME® is also registered in Germany.

Open for Innovation    ®
KNIME