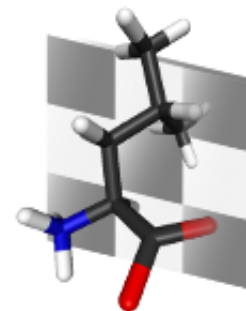


CheS-Mapper: New Developments



Martin Gütlein

mguetlein@fdm.uni-freiburg.de



JOHANNES GUTENBERG
UNIVERSITÄT MAINZ

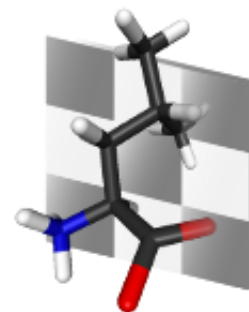
Andreas Karwath & Stefan Kramer



OpenTox Euro 2013, October 2

CheS-Mapper

- Chemical Space Mapping and Visualization in 3D

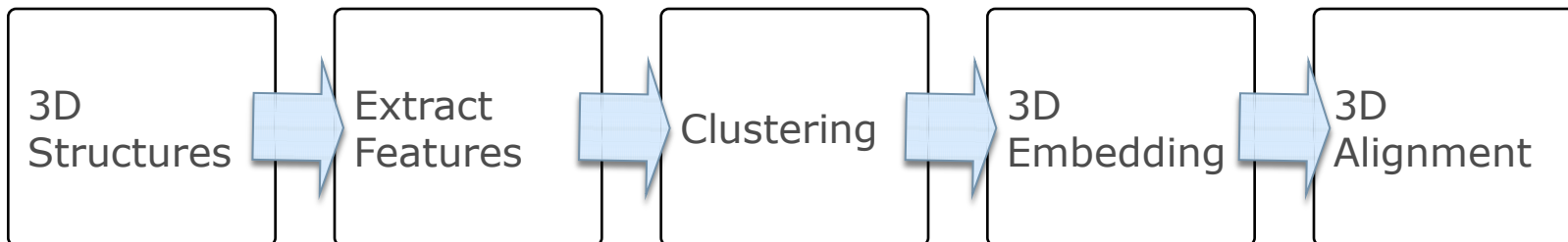


- 3D viewer for small molecule datasets
- Published in Journal of Cheminformatics, March 2012, >6000 accesses in 18 months
- Project homepage: <http://ches-mapper.org>
- Open-source java software
- Uses: Jmol, CDK, WEKA, OpenBabel, R
- Compatible to OpenTox dataset services

CheS-Mapper Workflow

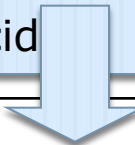


Chemical Space Mapping

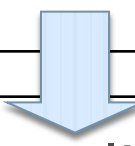
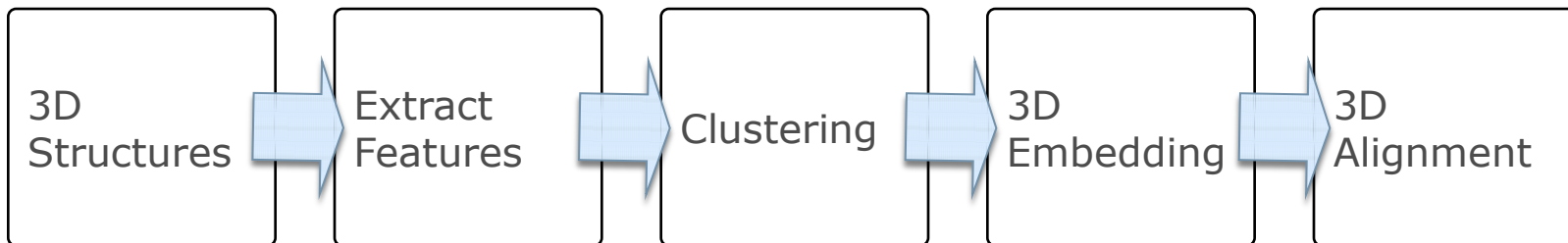


3D-Visualization

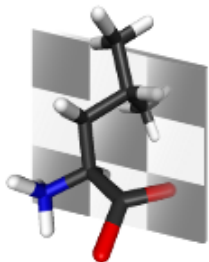
<chem>Cn1cnc2c1c(=O)n(c(=O)n2C)C</chem>	Caffeine	
<chem>CCCCC1CC2C(C3C=C(CCC3C(O2)(C)C)C)c(c1)O</chem>	THC	
<chem>O=C(Nc1c(ccc1C)C)CN(CC)CC</chem>	Lidocaine	
<chem>CN1CCCC1c1cccnc1</chem>	Nicotine	
<chem>CN1c2c(C(=NCC1=O)c1ccccc1)cc(Cl)cc2</chem>	Diazepam	
<chem>O1C(CO)C(O)C(O)C(O)C1O[C@@]1(OC(C(O)C1O)CO)CO</chem>	Sucrose	
<chem>OCC1OC(O)C(O)C(O)C1O</chem>	Glucose	
<chem>OC(=O)CC(O)(C(=O)O)CC(=O)O</chem>	Citric acid	
<chem>OS(=O)(=O)O</chem>	Sulfuric acid	
<chem>OP(=O)(O)O</chem>	Phosphoric acid	



Chemical Space Mapping



3D-Visualization



Load Dataset (step 1 of 6)

Select a dataset from your file system for clustering, embedding and visualization.

Select dataset file (Copy a http link into the textfield to load a dataset from the internet):

 Recently used datasets

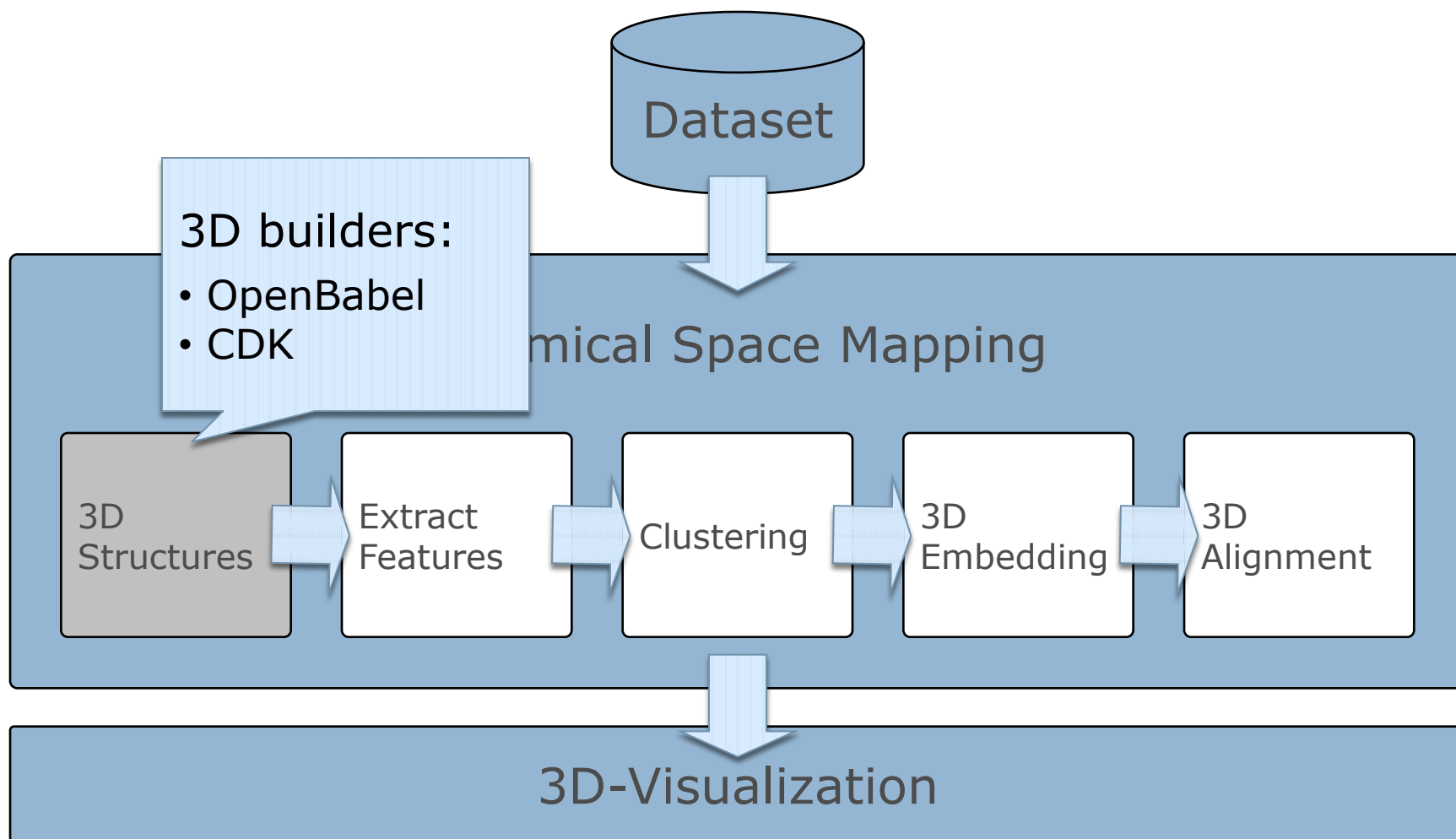
Dataset Properties:

File:	demo.smi
Num compounds:	10
Num properties per compound:	2
3D available:	true

Chemical space
mapping can be
configured with wizard



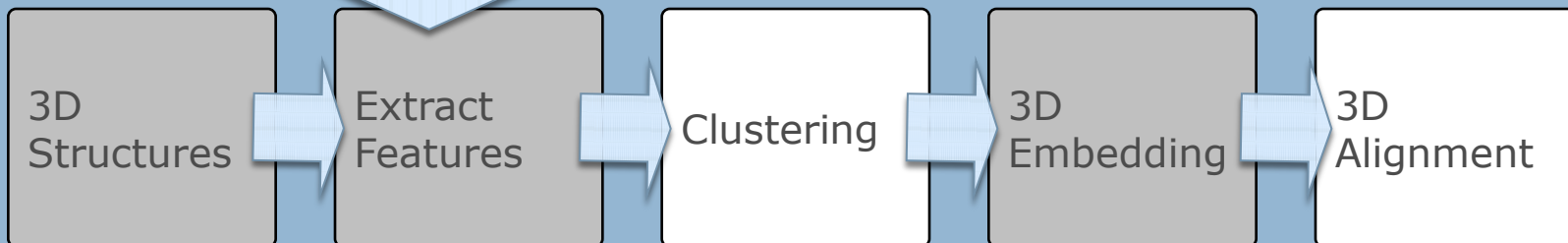
CheS-Mapper Workflow



CheS-Mapper Workflow

Feature Extraction:

- Included in the dataset
- PC-Descriptors (computed with CDK or OpenBabel)
- Structural Fragments
 - Subgraph mining
 - Smarts matching

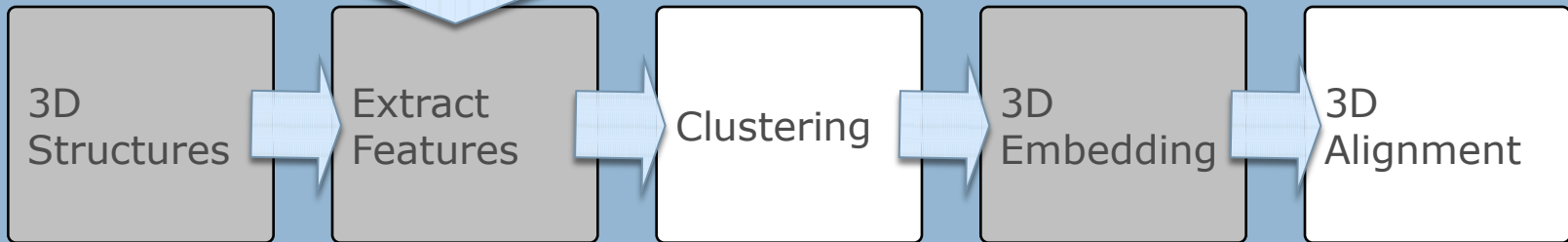


3D-Visualization

CheS-Mapper Workflow

Physico-chemical descriptors:

- MW (Molecular weight)
- logP (Octanol water coefficient)
- abonds (Number of aromatic bonds)
- HBD (Number of hydron bond donors)
- ...

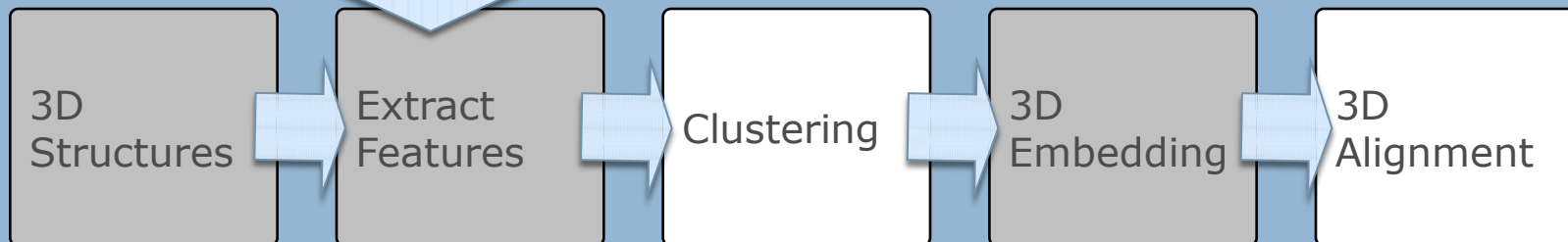


3D-Visualization

CheS-Mapper Workflow

MACCS list (166 SMARTS fragments):

- P
- [#8]~[#6]~[#8]
- *~[CH2]~[#8]
- ...

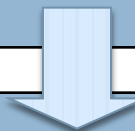
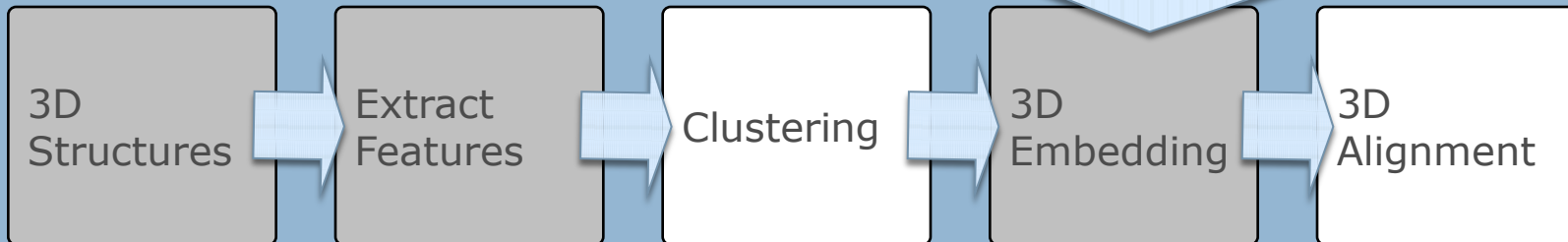


3D-Visualization

	MW	LogP	abonds	...	HBD
Caffeine					
THC					
Glucose					
Sulfuric acid					
...					



	X	Y	Z
Caffeine			
THC			
Glucose			
Sulfuric acid			
...			

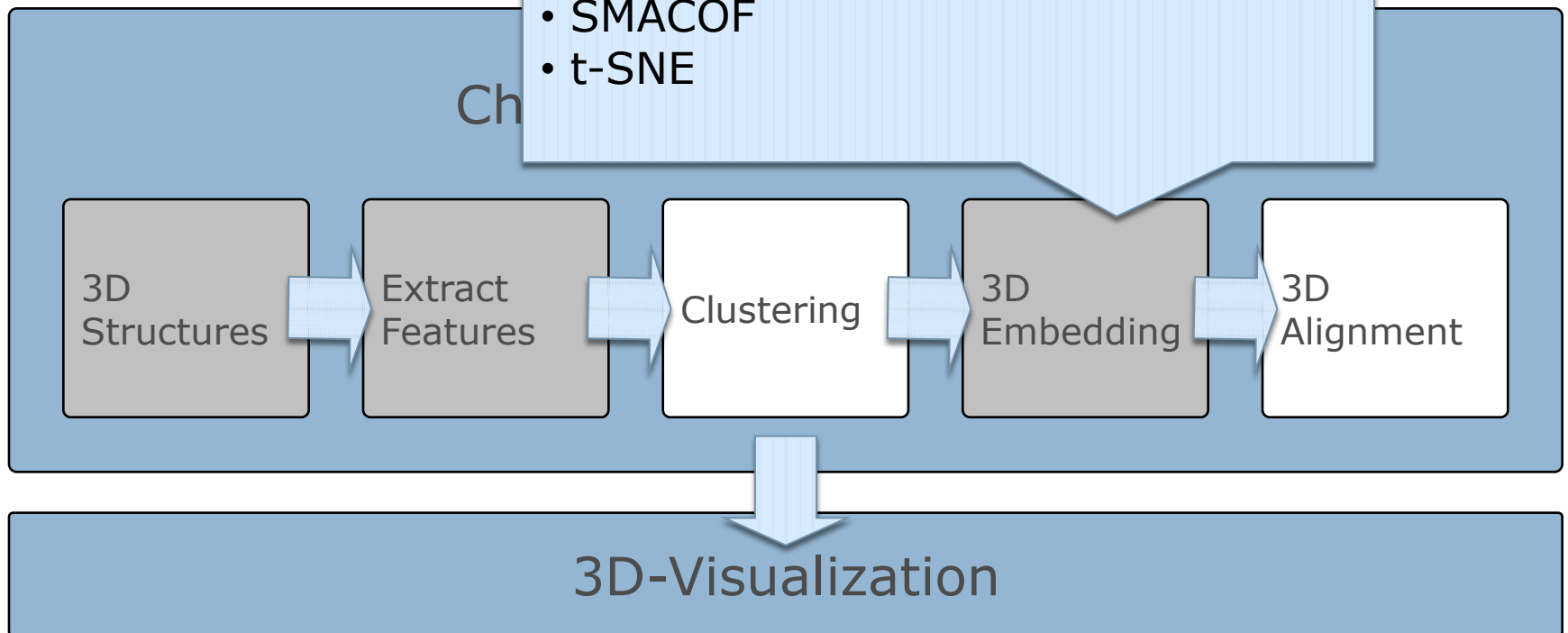


3D-Visualization

CheS-Mapper Workflow

3D embedding algorithms:

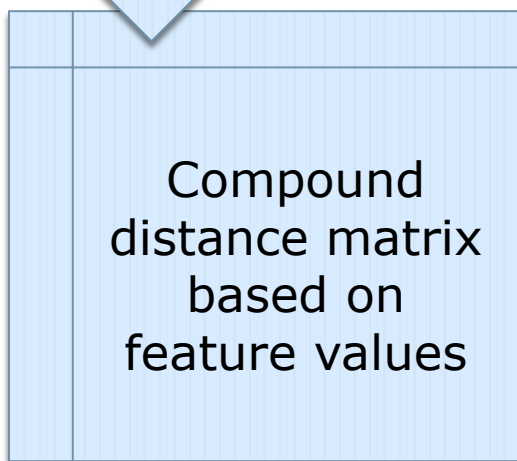
- PCA (fast)
- Sammon embedding (non-linear, configurable distance)
- SMACOF
- t-SNE



Quality of the 3D Embedding

	MW	LogP	abonds	...	HBD		X	Y	Z
Caffeine						Caffeine			
THC						THC			
Glucose						Glucose			
Sulfuric acid						Sulfuric acid			
...						...			

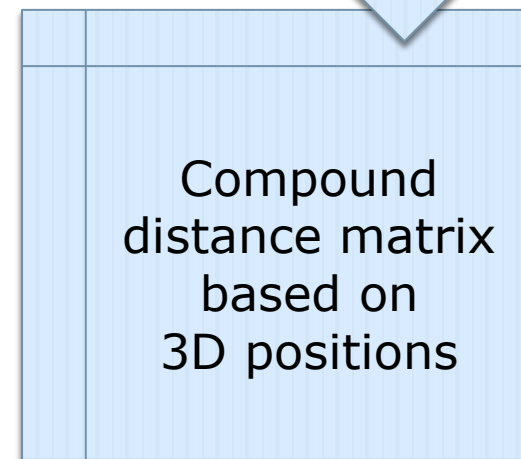
How well do the 3D positions reflect the feature values?



Global embedding quality: Correlation between distance matrixes (CCC, R^2)

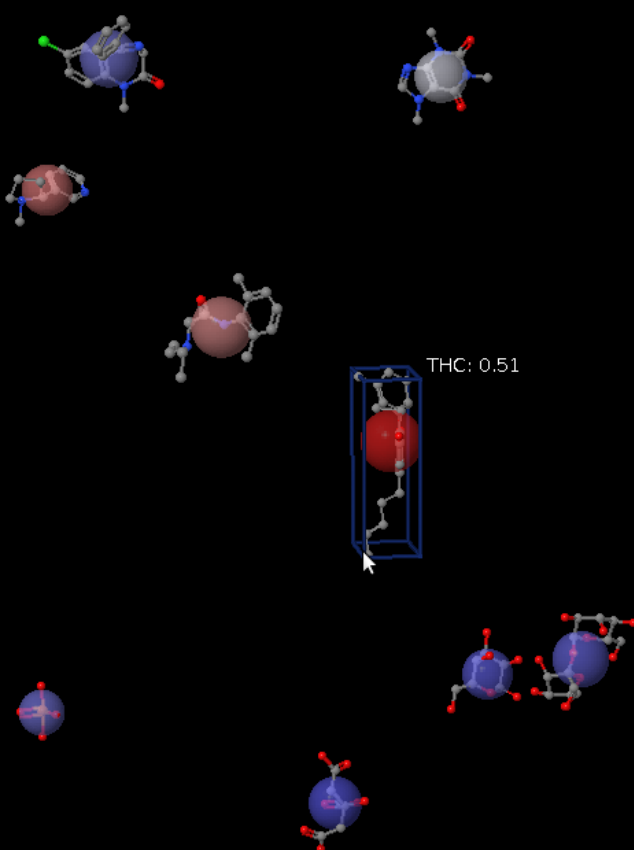
Embedding stress for each compound:

1 - correlation of the corresponding matrix rows



Citric acid	0.05
Glucose	0.07
Sucrose	0.08
Sulfuric acid	0.1
Phosphoric acid	0.11
Diazepam	0.12
Caffeine	0.26
Lidocaine	0.38
Nicotine	0.4
THC	0.51

■ Superimpose

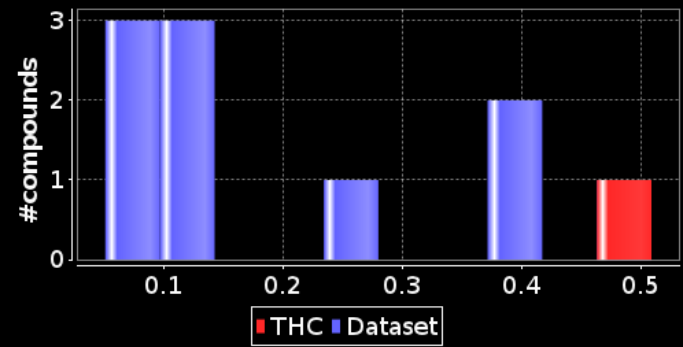


THC is the compound with the highest embedding stress

Dataset: demo.smi
Num compounds: 10
Cluster algorithm: No Dataset Clustering
3D Embedding: Sammon 3D Embedder (R)
3D Embedding Quality: good (CCC: 0.86, r²: 0.7)

Ground	THC
CCCCCc1cc...	
C=C(A)A	match
C=C	match
CC(C)(C)A	match
CH3AAACH2A	match
Onot%A%A	match
CH3CH2A	match
CH3ACH2A	match
ACH2CH2A ...	match
ACH2AACH2A	match
6M ring > 1	match
ACH2CH2A	match

Feature: Embedding stress
Values: 0.11 ±0.17
Description: 0 := perfectly embedded (no stress), 1: h
Usage: NOT used for clustering and/or embeddin
Missing values: 0



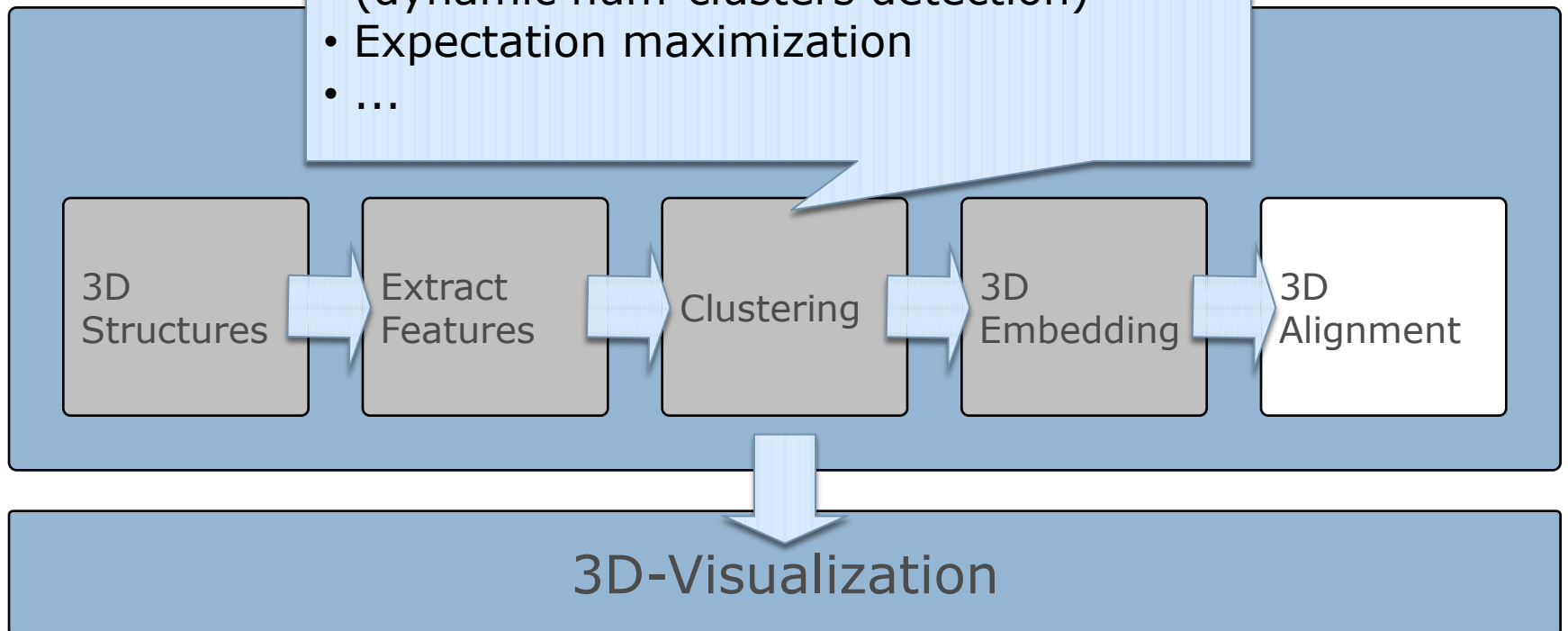
[-] [] [+] ○ Wireframe ● Balls & Sticks ● Dots

Feature: Embedding stress [v] Label

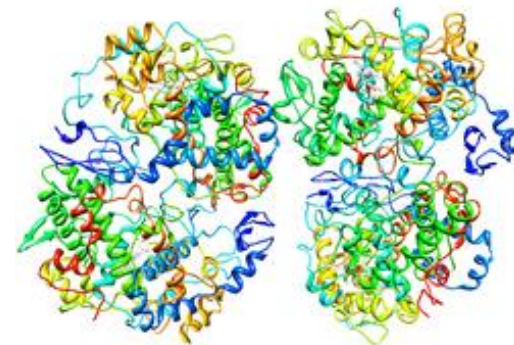
CheS-Mapper Workflow

Cluster Algorithms:

- k-Means (num clusters configurable)
- Hierarchical Clustering (dynamic num-clusters detection)
- Expectation maximization
- ...



COX-2 dataset



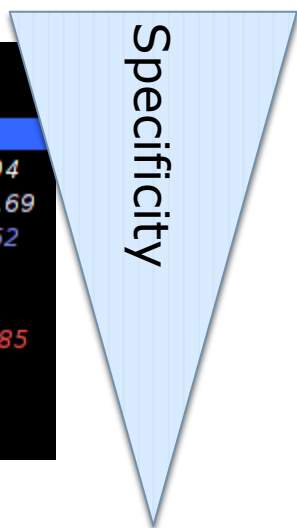
- 467 COX-2 Inhibitors
- Inhibition of the enzyme Cyclooxygenase-2 (COX-2) is investigated in cancer studies
- Compounds are structurally very similar (docking dataset)
- Endpoint: IC_{50} μ Mol (half maximal inhibitory concentration)

Spline-Fitting with a Genetic Algorithm: A Method for Developing Classification Structure–Activity Relationships

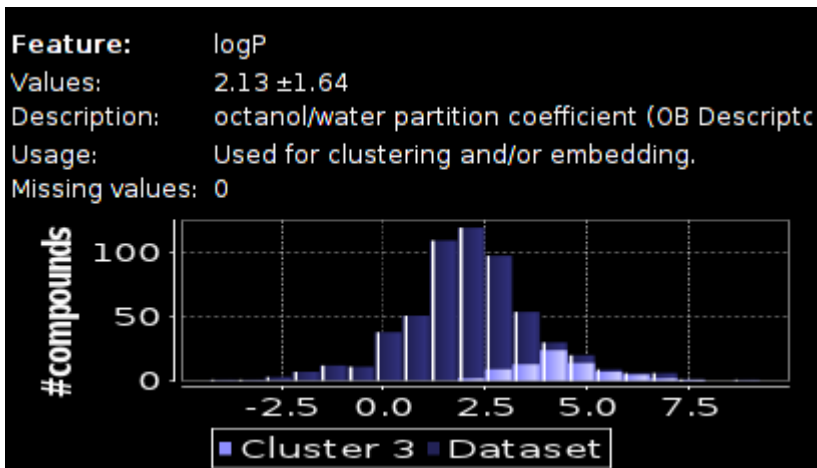
*Jeffrey J. Sutherland, Lee A. O'Brien, and, and Donald F. Weaver
Journal of Chemical Information and Computer Sciences 2003 43 (6),
1906-1915*

Sorting of features

Cluster	Cluster 3
Size	77
logP	4.35 ±1.05
MR	67.71 ±15.04
MW	246.39 ±76.69
TPSA	20.31 ±16.52
HBA1	1 ±0.72
HBA2	1 ±0.79
DSSTox_CID	2080 ±447.85
HBD	0 ±0.57
LC50_mmol	0 ±0.09
nF	0 ±0.23

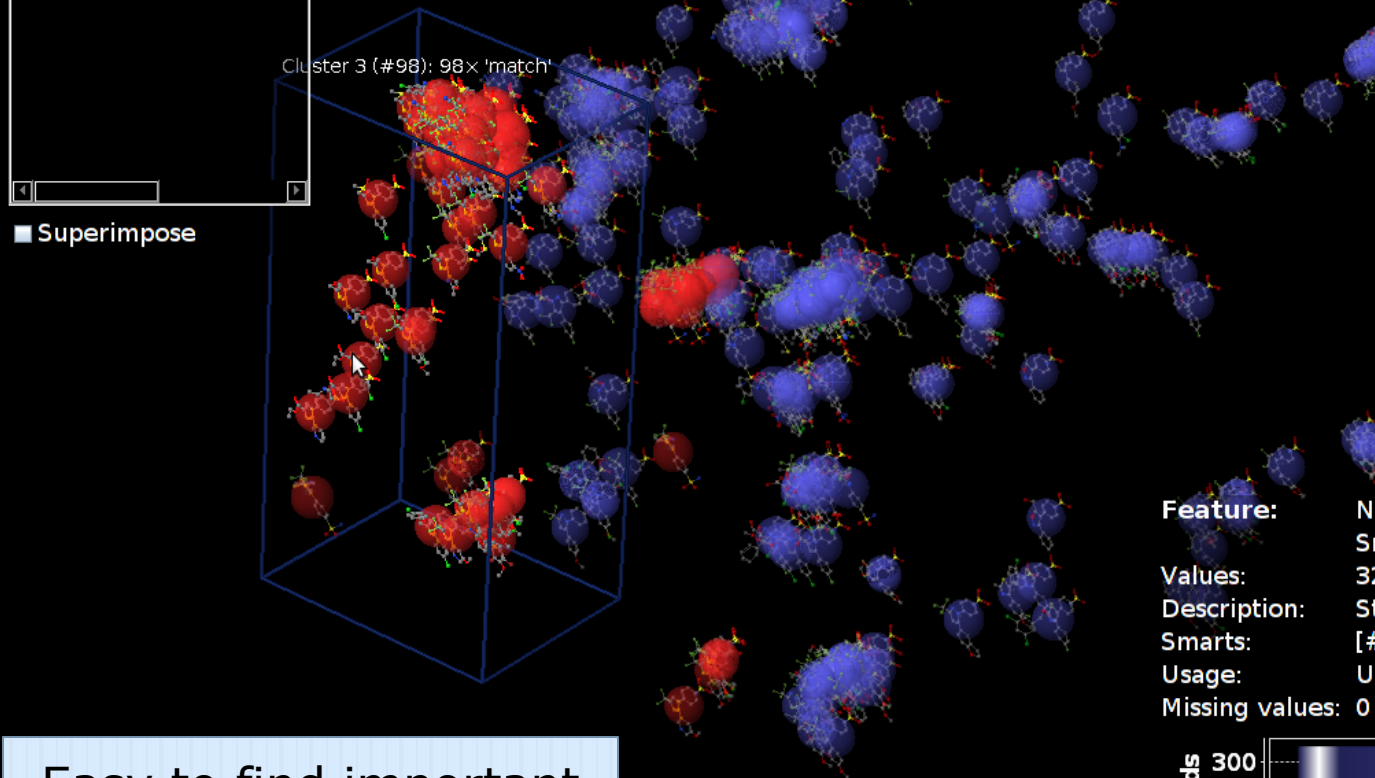


- Feature values are shown for the selected compound or cluster
- Features are sorted according to the p-value of
 - X² test for nominal features
 - ANOVA test for numerical features
- The most 'important' features are listed first: the feature values of this compound/cluster differ the most from the complete dataset



All clusters

Cluster 4 (#40)	30x 'match'
Cluster 6 (#39)	39x 'match'
Cluster 7 (#39)	39x 'match'
Cluster 5 (#51)	51x 'match'
Cluster 2 (#110)	79x 'match'
Cluster 1 (#90)	90x 'match'
Cluster 3 (#98)	98x 'match'



Dataset: cox2.sdf
Num compounds: 467
Cluster algorithm: Hierarchical - Dynamic Tree Cut
3D Embedding: Sammon 3D Embedder (R)
3D Embedding Quality: moderate (CCC: 0.64, r²: 0.12)

Cluster	Cluster 3
Size	98
MCS	O=S(=O)(c...
Family	71x 'B.1', 1...
NC(C)N	98x 'match'
NAN	98x 'match'
NAAN	98x 'match'
CN(C)C	98x 'match'
QN	96x 'no-ma...
CSN	98x 'no-ma...
NS	98x 'no-ma...
QQH	98x 'no-ma...
NAO	96x 'no-ma...
NH2	96x 'no-ma...
N > 1	98x 'match'

Specificity

Feature: NC(C)N
Smarts list: MACCS (OpenBabel MACCS)
Values: 328x *no-match*, 139x *match*
Description: Structural Fragment, matched with Open
Smarts: [#7]~[#6](~[#6])~[#7]
Usage: Used for clustering and/or embedding.
Missing values: 0



Easy to find important features for clusters

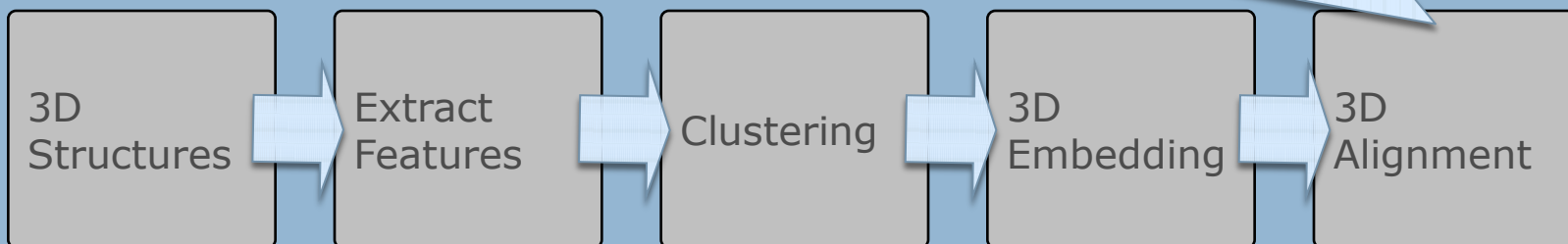
Wireframe
 Balls & Sticks
 Dots

Feature: Label

CheS-Mapper Workflow

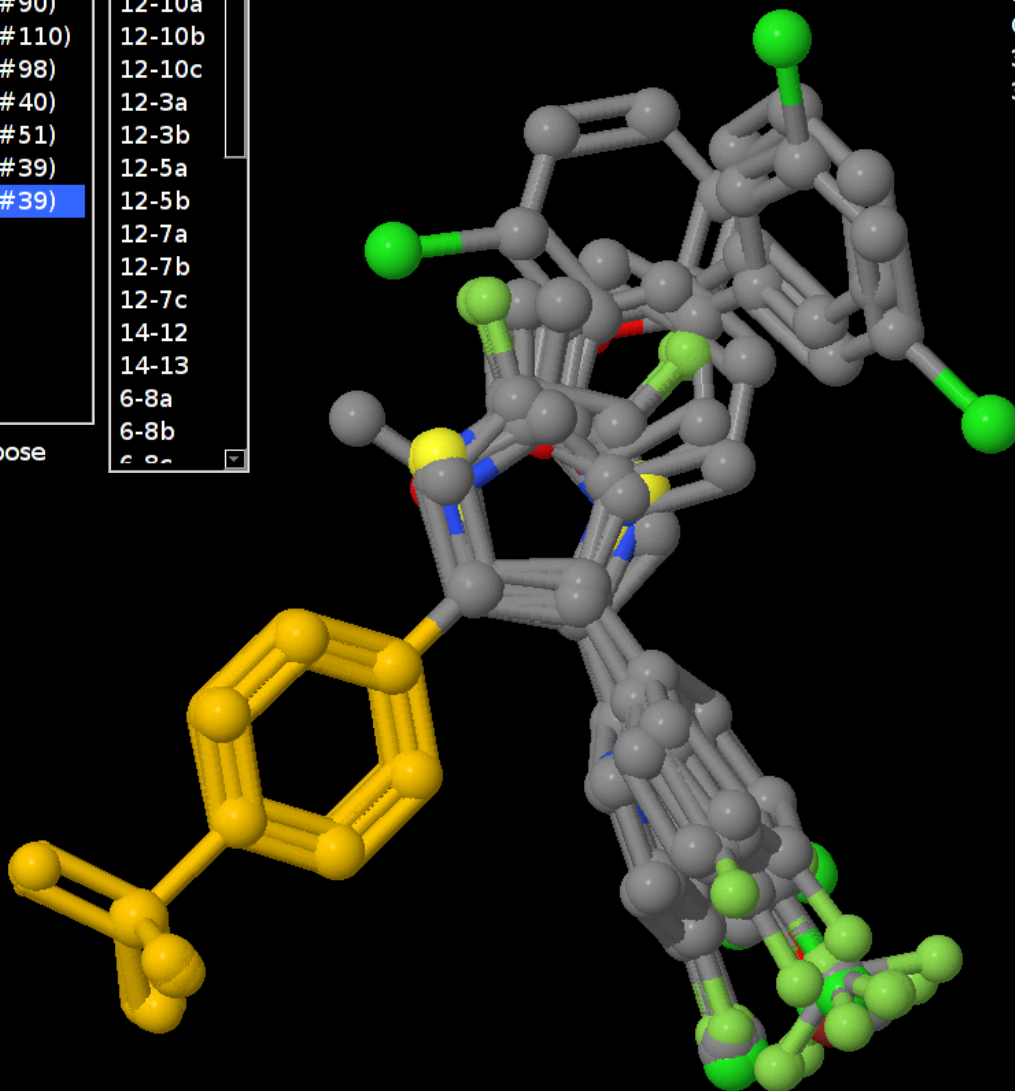
3D Alignment:

- Detect common subgraph of each cluster:
 - Compute MCS (maximum common subgraph)
 - Use largest fragment
 - Manually specify subgraph
- Align compounds according to common subgraph:
 - Obfit (OpenBabel)
 - Kabsch Alignment (CDK)



3D-Visualization

All clusters	10-2
Cluster 1 (#90)	12-10a
Cluster 2 (#110)	12-10b
Cluster 3 (#98)	12-10c
Cluster 4 (#40)	12-3a
Cluster 5 (#51)	12-3b
Cluster 6 (#39)	12-5a
Cluster 7 (#39)	12-5b
	12-7a
	12-7b
	12-7c
	14-12
	14-13
	6-8a
	6-8b
	6-8c

 Superimpose


Dataset: cox2.sdf
Num compounds: 467
Cluster algorithm: Hierarchical - Dynamic Tree Cut
3D Embedding: Sammon 3D Embedder (R)
3D Embedding Quality: moderate (CCC: 0.64, r²: 0.12)

Cluster	Cluster 7
Size	39
MCS	O=S(=O)(O)
AN(A)A	39x 'no-match'
A!N\$A	39x 'no-match'
C-N	38x 'no-match'
QCH3	36x 'no-match'
A\$A!N	38x 'no-match'
Nnot%A%A	38x 'no-match'
NH2	39x 'match'
CSN	39x 'match'
NS	39x 'match'
QQH	39x 'match'
NAO	39x 'match'
NH	39x 'match'
QN	39x 'match'
QAAAA@1	27x 'no-match'
Heterocycle	25x 'no-match'
CH3	23x 'no-match'
N Heterocycle	26x 'no-match'
C%N	26x 'no-match'
Family	10x 'D.1', ...
CN(C)C	39x 'no-match'
Heterocyclic atom > 1 (&....	28x 'no-match'
NAAAN	39x 'no-match'

Cluster compounds
aligned according to MCS

 Wireframe Balls & Sticks Dots

Feature: MCS

All clusters

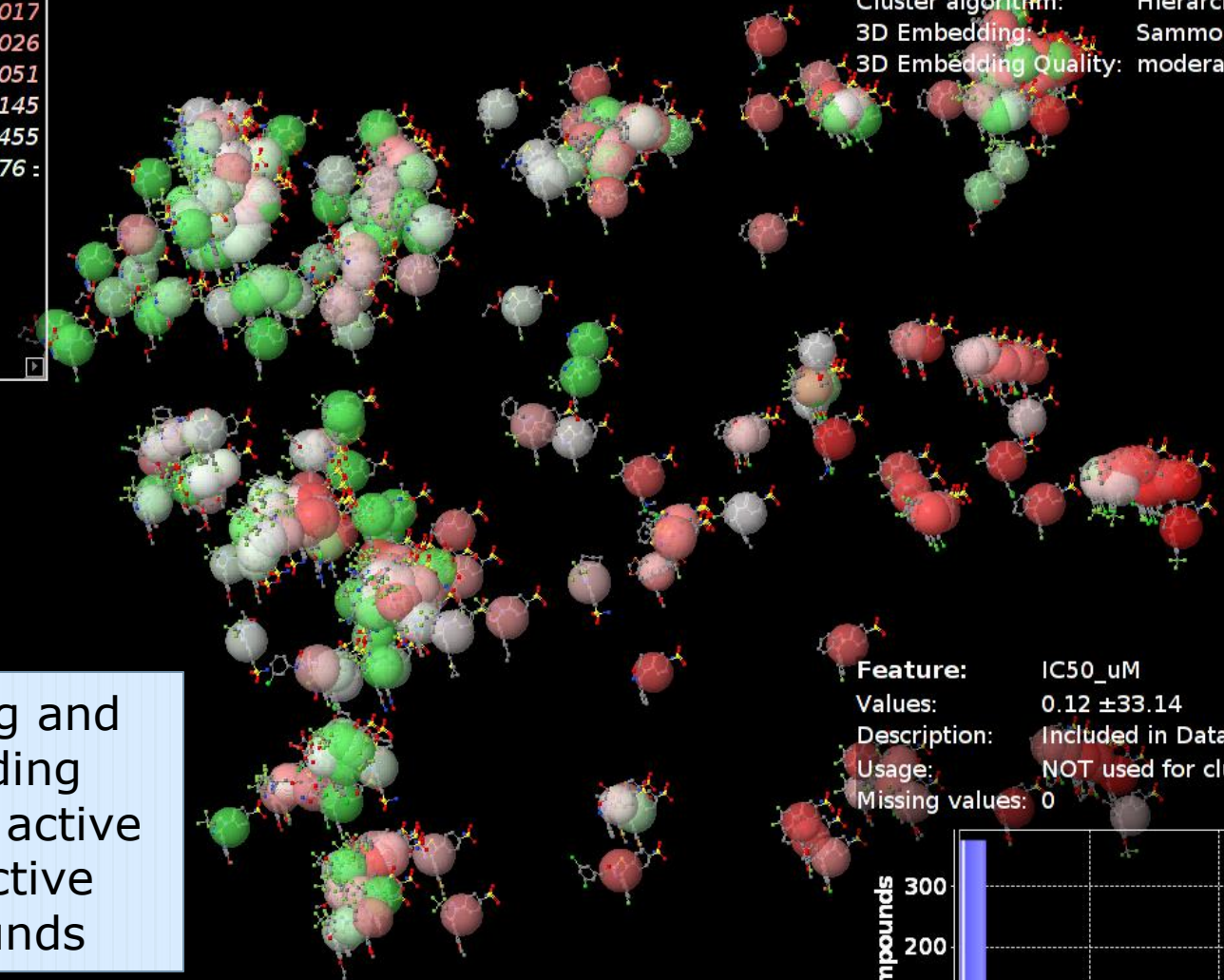
Cluster 7 (#39)	0.007
Cluster 6 (#39)	0.017
Cluster 5 (#51)	0.026
Cluster 4 (#40)	0.051
Cluster 2 (#110)	0.145
Cluster 1 (#90)	0.455
Cluster 3 (#98)	0.76 :

Dataset: cox2.sdf
Num compounds: 467
Cluster algorithm: Hierarchical - Dynamic Tree Cut
3D Embedding: Sammon 3D Embedder (R)
3D Embedding Quality: moderate (CCC: 0.64, r²: 0.12)

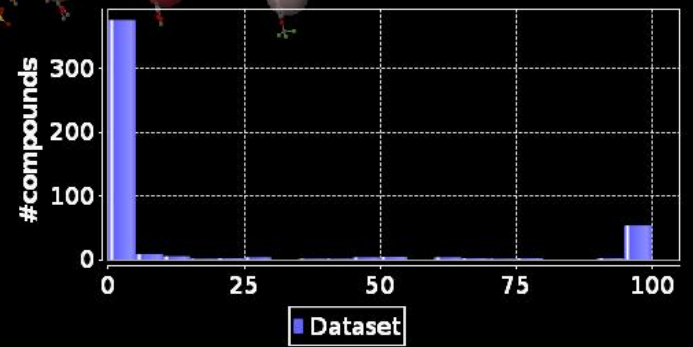
Superimpose

Clustering and embedding separates active and inactive compounds

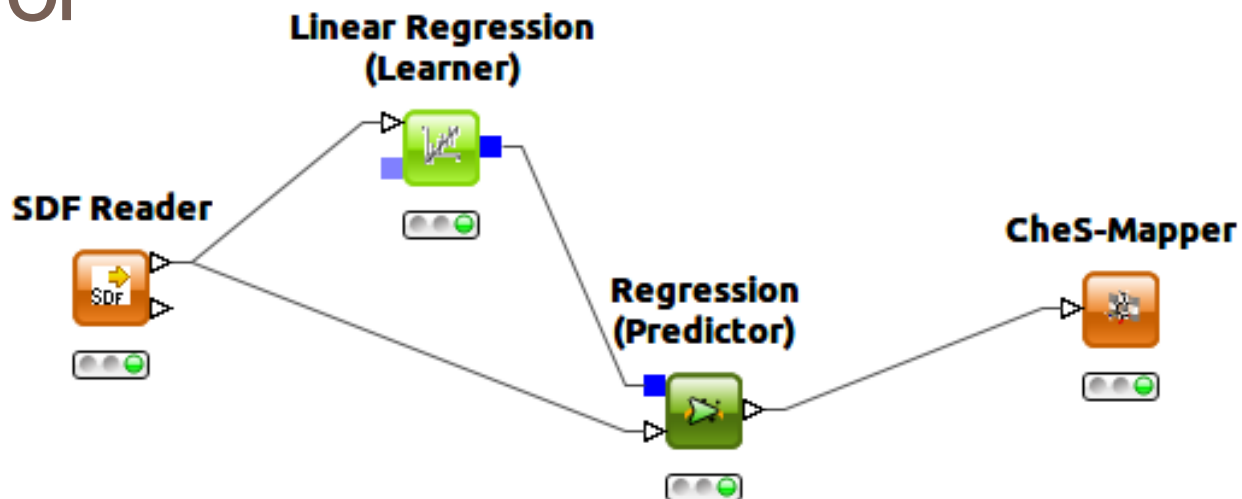
[-] [] [+] Wireframe Balls & Sticks Dots
Feature: IC50_uM Label



Feature: IC50_uM
Values: 0.12 ±33.14
Description: Included in Dataset
Usage: NOT used for clustering and/or embedding
Missing values: 0



The CheS-Mapper extension for KNIME



- **KNIME:**
 - graphical workbench for data access, investigation and predictive analysis
 - various extensions to process chemical data
- **CheS-Mapper integration as visualization node**

EPAFHM dataset



- EPAFHM:
US Environmental Protection Agency (EPA)
Fathead Minnow Acute Toxicity Database File
- 617 industrial organic chemicals
- Endpoint: LC50 mMol ((lethal) concentration that kills 50%)

Predicting modes of action from chemical structure: Acute toxicity in the fathed minnow (*Pimephales promelas*).

Russom, C.L., S.P. Bradbury, S.J. Broderius, D.E. Hammermeister, and R.A. Drummond (1997)

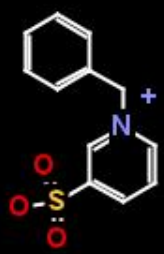
Environmental Toxicology and Chemistry, 16(5): 948-967.

- 0.74 2,2-Dimethyl-1-propyl:
- 0.74 Benzamide
- 0.8 Cyclohexanone
- 0.82 3-Methyl-3-pentanol
- 0.84 2-Cyanopyridine
- 0.84 2-Methoxyethylamine
- 0.85 Cyclohexanol
- 0.87 3-Acetamidophenol
- 0.88 tert-Butyl methyl ethe
- 0.88 3-Amino-5,6-dimethyl-
- 0.89 Isopropyl ether
- 0.94 5-Chloro-2-pyridinol
- 0.98 2-Picoline
- 0.99 1-Benzylpyridinium 3-s

Superimpose

Cluster: Single cl
 Num compounds: 579
 3D Alignment: No Clust
 LC50_mmol_log: [-6.38; 2]

Dataset: knime_inj
 Num compounds: 579
 Cluster algorithm: No Datas
 3D Embedding: Sammon
 3D Embedding Quality: good (r^



Compound:
 0.99 1-Benzylpyridinium 3-sulfonate

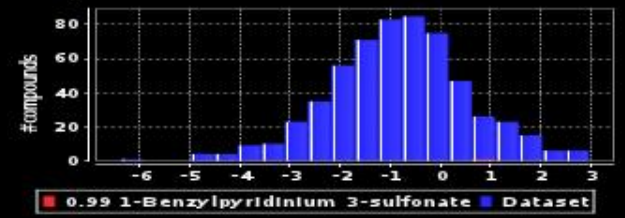
Endpoint:	LC50
Species:	fathead
ChemClass_FHM:	Pyridines
CLOGP:	-99.99
MLOGP:	null
LC50_mg:	2410
LC50_mmol:	9.67
ActivityOutcome_EPAFHM:	active
ActivityScore_EPAFHM:	21
LC50_Ratio:	1.39
LC50_Note:	null
MOA:	MOA not
MOA_Confidence:	null
MOA_MixtureTest:	null
ExcessToxicityIndex:	0
FishAcuteToxSyndrome:	null
FishBehaviorTest:	Conflictir
Note_EPAFHM:	null
LC50_mmol_log:	0.99
LC50_mmol_log (predicti...:	-1.53

Predicted

Actual

Highlight two features at once to detect high prediction errors

Feature: LC50_mmol_log
 Description: Included in Dataset
 Usage: NOT used for clustering and/or embed
 Missing values: 0



Feature: LC50_mmol_log Label

More features

- Export clusters/compounds/features
- Export high-res images
- Access ChEMBL database
- Save and share embedding settings
- Data tables to browse through raw compound/feature/cluster data
- Command line interface
- Configurable highlight and view settings
 - Adjust highlight color gradient
 - Enable log highlighting
 - Switch between sphere and atom-color highlighting
 - ...

Quotes about CheS-Mapper

... a very nice piece of software.

Basil Hartzoulakis, PhD, Xention Ltd, Cambridge

You help to make me look like a genius to my bosses ... this is the only open source tool I know of that is usable by a regular bench chemist.

*Kerry W. Fowler, Ph.D., Senior Scientist
Kineta Inc, Seattle*

CheS-Mapper has come in handy.

*Kaushik Hatti, Vittal Mallya Scientific Research
Foundation, Banagalore*

I am very impressed by the ches-mapper software ... the tool is very effective to spot trends of a chemical group

Hiroshi Nara, Organic chemist, Japan

thanks for the gorgeous ches-mapper.

Santi Villalba, University College Dublin

Many compliments because we have found that it is highly useful and well working.

Prof. Paola Gramatica, University of Insubria