



Tomorrow's Drugs Today™

trust value
productivity innovation
commitment speed
synergistic partnerships
innovation

Mark Brewer, 19th June 2007
4th Joint Sheffield Conference on
Chemoinformatics

A spectral clustering approach for the analysis of screening data

- Techniques for generating screening data
- Analysing screening data
- A spectral clustering approach
- Example application to COX-2 inhibitors
- Provision of results and applications
- Some possibilities for future work and summary

High throughput screening

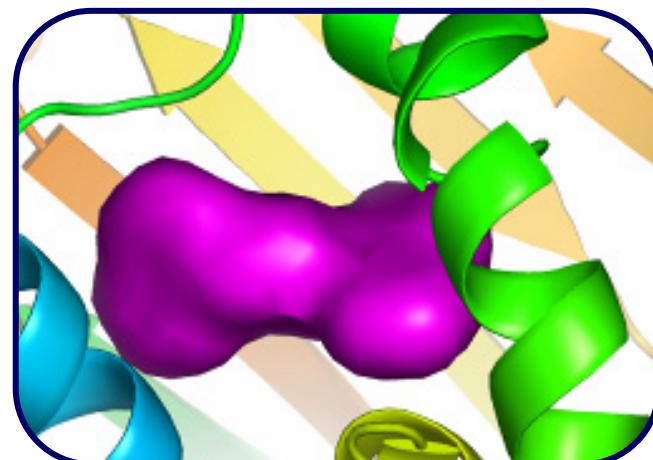
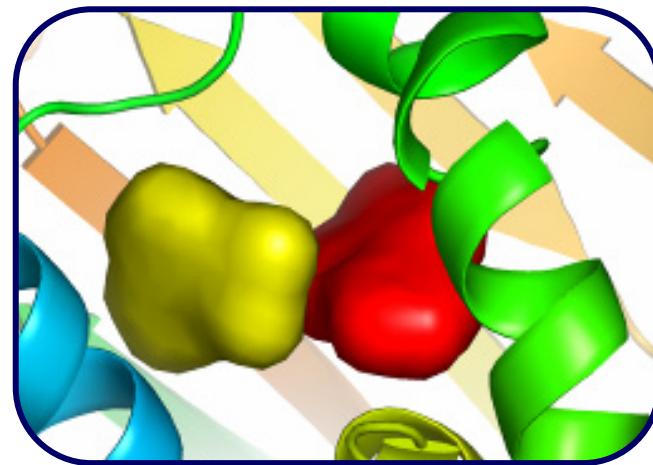
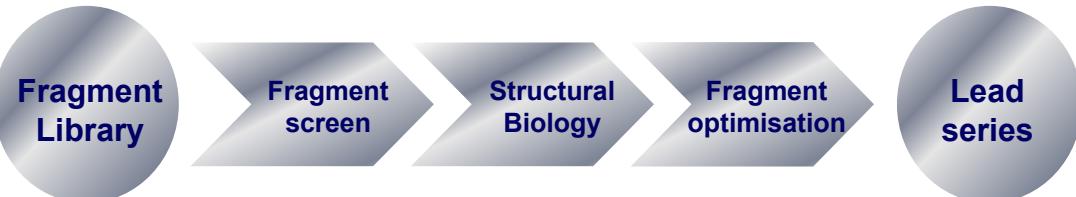
- High throughput screening factory that can generate 100,000 data points a day per machine
- Developed with partners:



Fragment screening

- Confocal Fluorescence Spectroscopy
- NMR Spectroscopy
- 30000 fragment library

Fragment based drug discovery

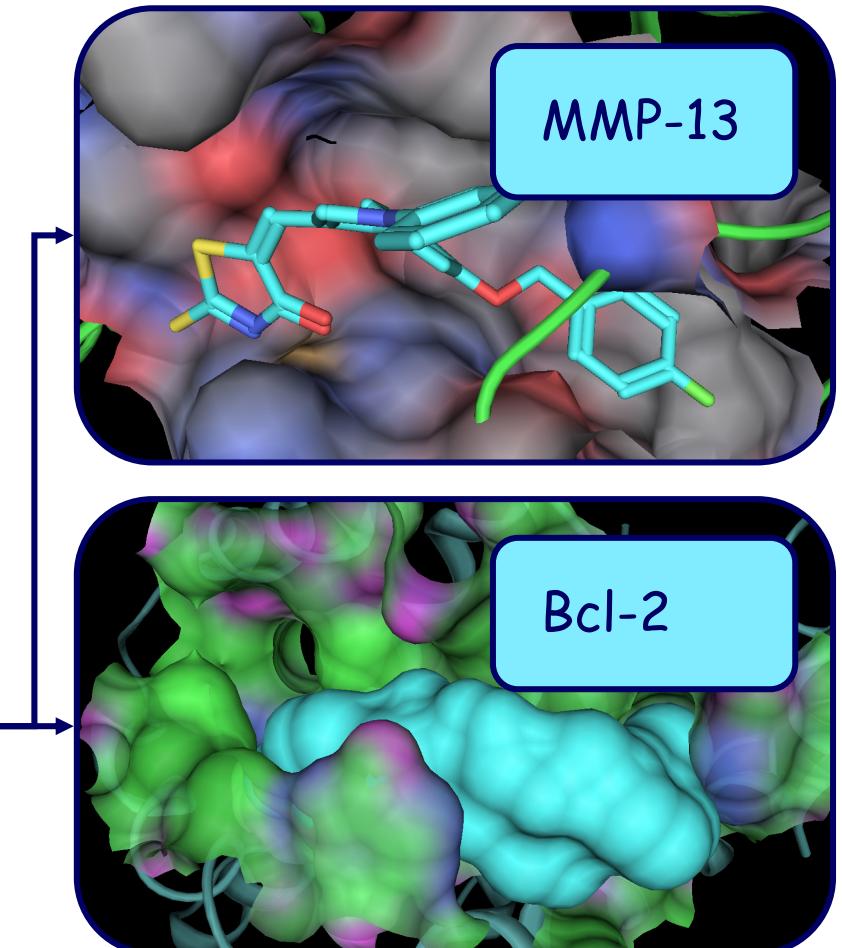


Virtual screening approaches

- Substructure searches
- Similarity searches
- Pharmacophore-based approaches
- Molecular docking
- Physiochemical properties
- *In silico* ADMET
-

Screening database

High throughput docking employs a computing grid distributed across 400 PCs



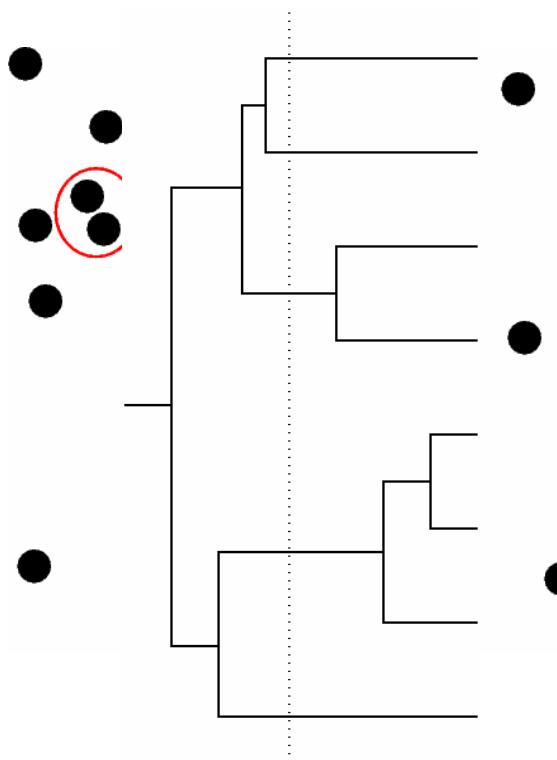
Clustering for the analysis of screening data

- Both *in vitro* and *in silico* screening yield lists of compounds
- Medicinal chemists will inspect *in vitro* hits to identify compounds to take into hit-to-lead programs
 - Assess number and representation of different scaffolds
 - Molecules with similar structure constitute hit series
 - Realised by software packages which cluster molecules
- Virtual screening hit lists contain list of compounds for purchase and/or synthesis
 - Diversity can be an important caveat so clustering can help here also
- Lots of other applications of clustering e.g. profiling of compound collections

Clustering molecular datasets

- Many types of clustering method available: hierarchical, non-hierarchical, nearest neighbour, relocation, etc.
- There is a common general procedure
 1. Assign descriptors
 2. Use descriptors to calculate similarity
 3. Process similarities with a clustering algorithm

Hierarchical clustering:
e.g. Wards Method



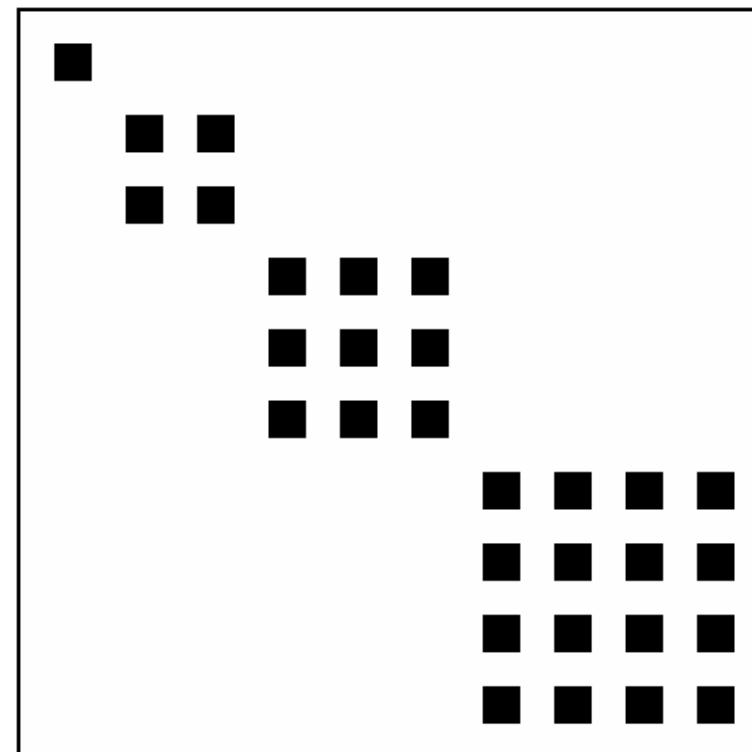
Clustering molecular datasets

- How many clusters?
- How to interpret and distribute clustering output?
- Ties in proximity?
- How to determine the validity of clusters?

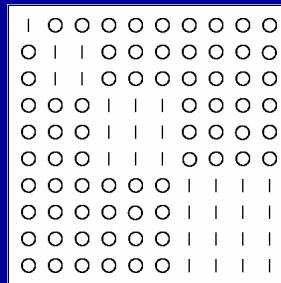
.... the best measure of the validity (and hence choice of clustering algorithm) has been the usefulness of the resulting classification ...[†]

A thought experiment

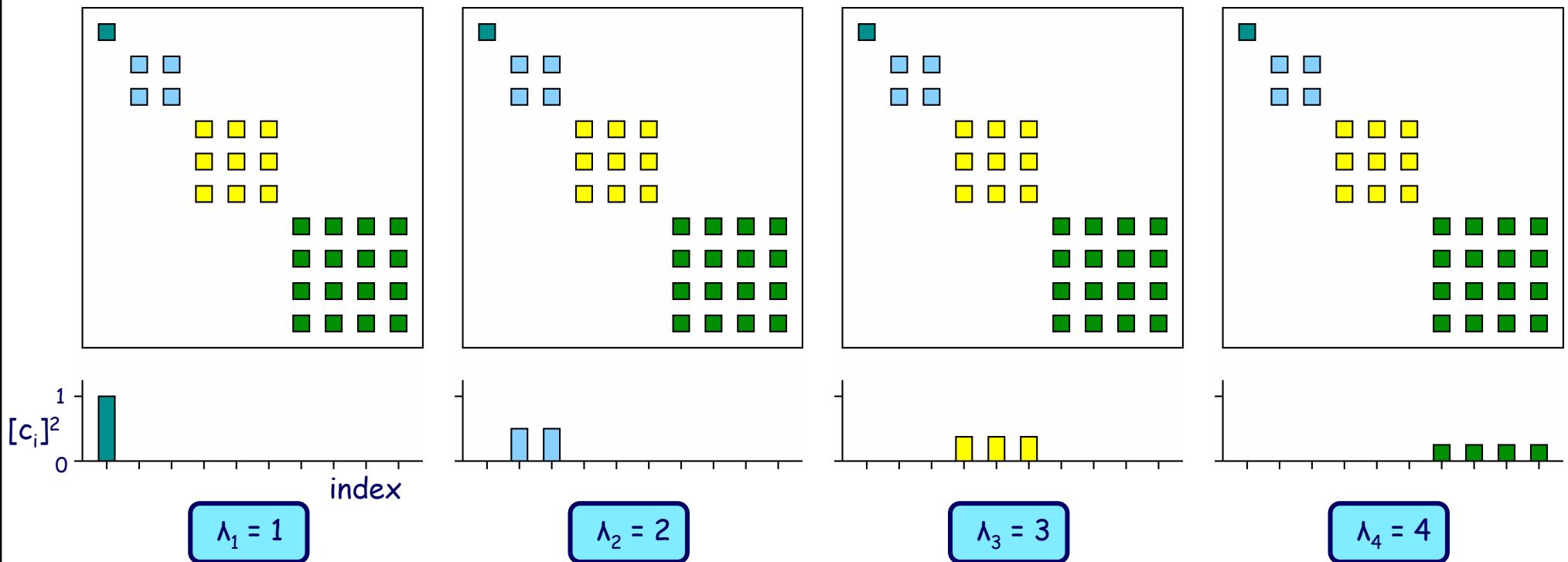
A hypothetical similarity matrix, where ■ denotes a pair of 'similar' molecules



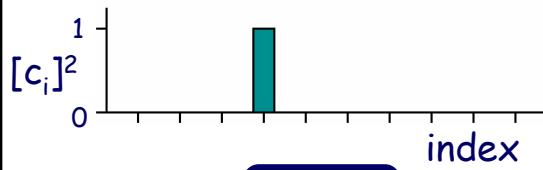
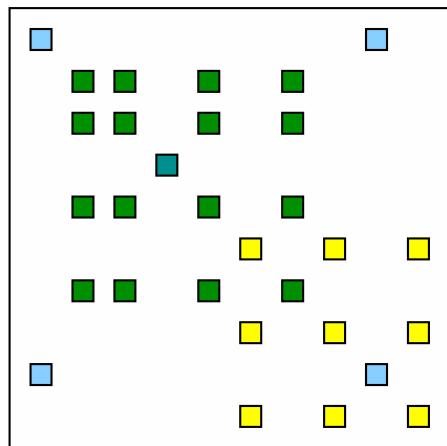
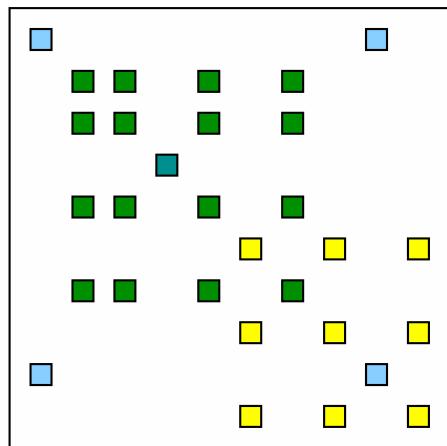
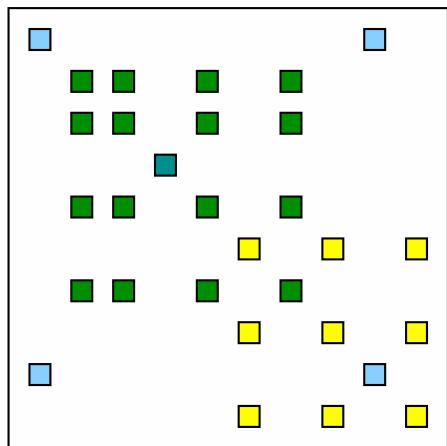
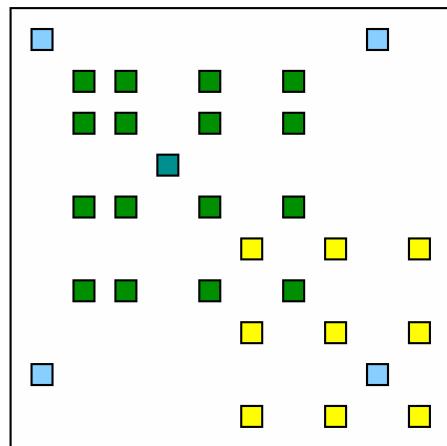
Ordered ideal matrix



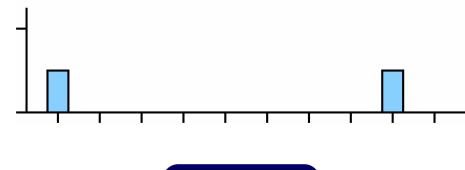
$$= S = C \Lambda C^T$$



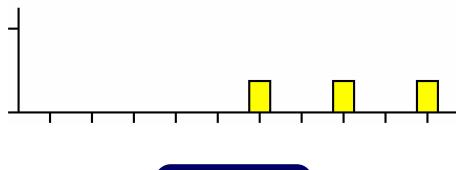
Permuted ideal matrix



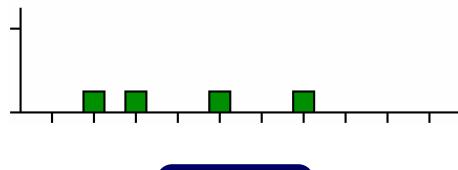
$$\lambda_1 = 1$$



$$\lambda_2 = 2$$



$$\lambda_3 = 3$$

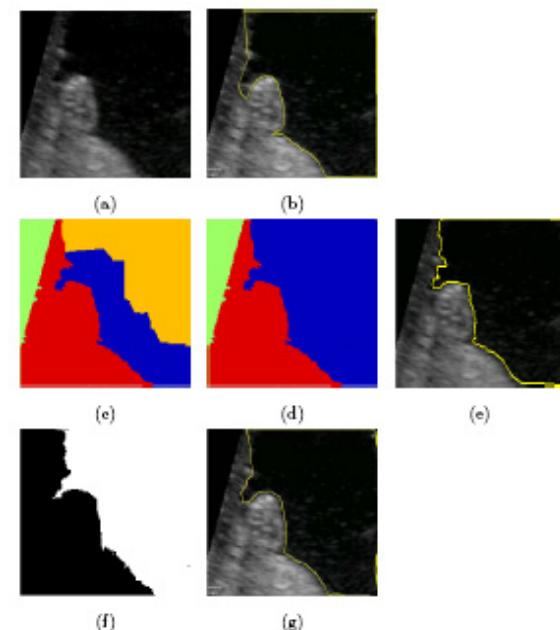


$$\lambda_4 = 4$$

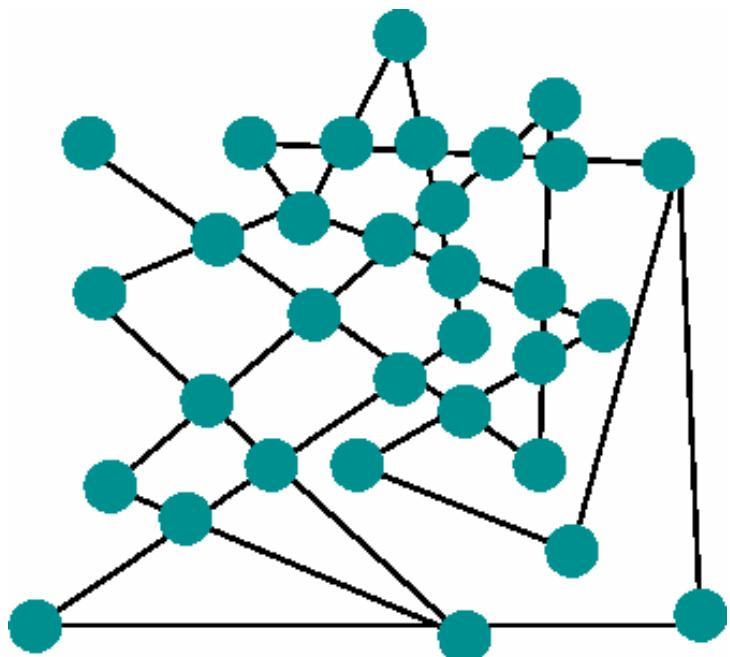
Existing applications of spectral clustering

- Eigenvalues and eigenvectors form the basis of spectral clustering algorithms
- Spectral clustering algorithms are used routinely to process digital images
 - Segmentation: To obtain a compact representation of what is useful in the image
- Spectral clustering has also been used in other cheminformatics applications
 - Guha and Wild have used singular value decomposition to cluster molecular datasets

Foetal images processed using spectral clustering



Further background



If we recognise the similarity matrix (S) as a weighted adjacency matrix and diagonalise:

$$S = C \Lambda C^T$$

Results from theory of graph spectra show that

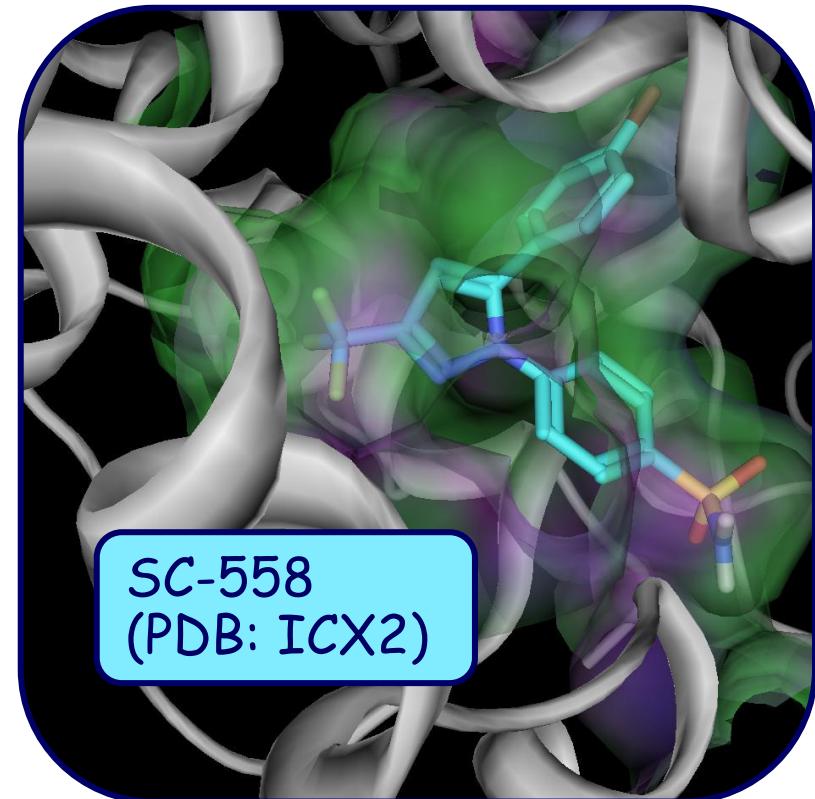
- 1) Eigenvalues (λ) relate to cluster cohesiveness
- 2) Eigenvectors (C) relate to cluster contribution and are normalised $CC^T = 1$
(analogous to fuzzy clustering membership functions)

Sorting by λ and C provide a means to inspect clusters

N.B. λ does not depend on the order of molecules

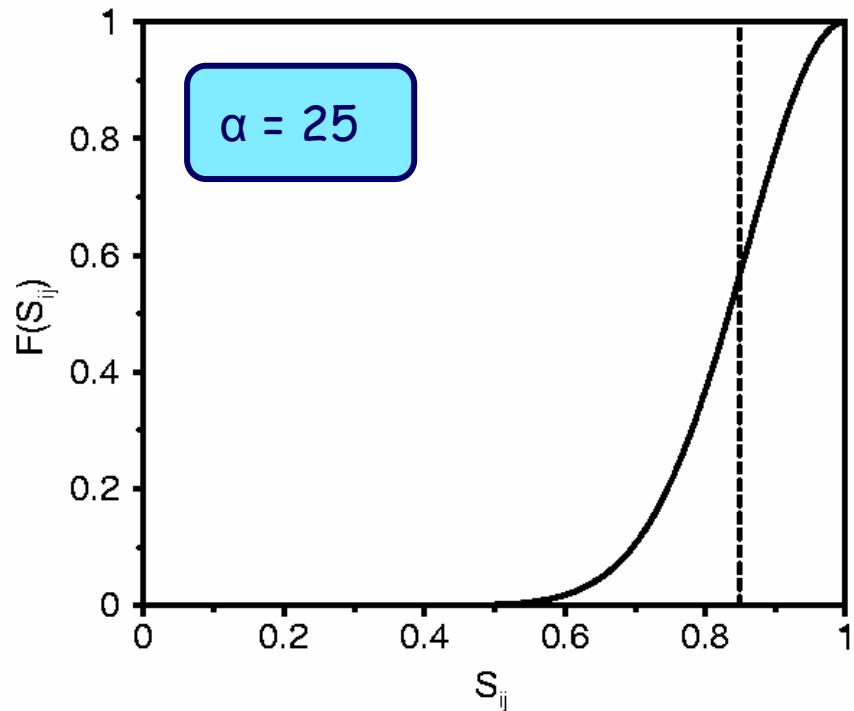
Example application: COX-2 inhibitors

- COX-2 is an important target for inflammation
- COX-2 dataset downloaded from the cheminformatics website
- Standard Tripos UNITY 2D molecular fingerprints for 125 unique chemical structures
- Programs written to compute Tanimoto similarities and diagonalise filtered similarity matrices
- Example clusters presented using the five largest eigenvalues and $\alpha = 25$



Filtering function

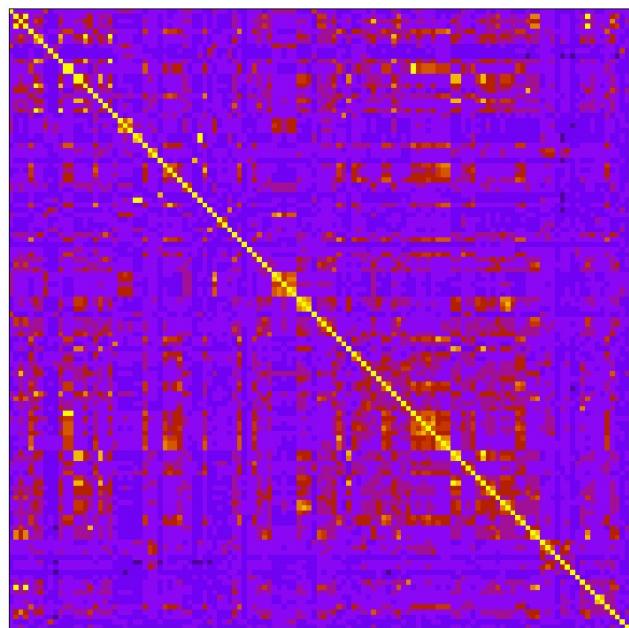
- Filtering function was found to be necessary to obtain useful clusters
- Removes background similarity but maintains chemically relevant levels of similarity
- Several filtering strategies investigated
- Gaussian filtering function selected
 - Single parameter with correct limiting behaviour



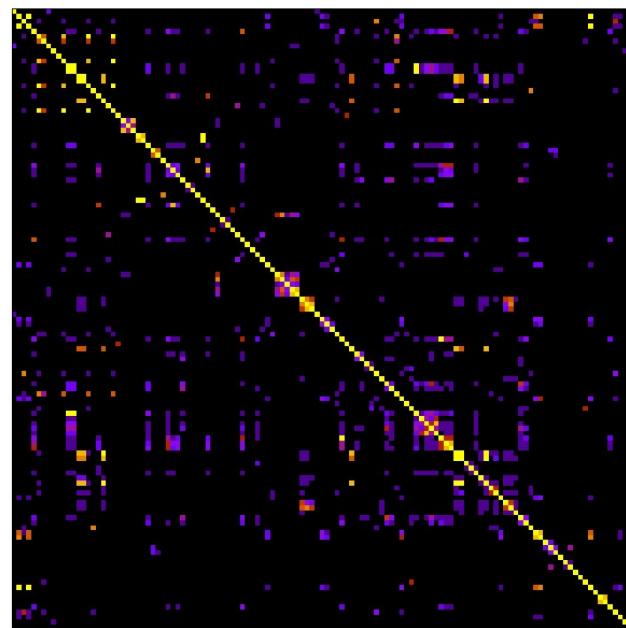
$$F(S_{ij}) = \exp[-\alpha(S_{ij} - 1)^2]$$

Effect of filtering function

Similarity

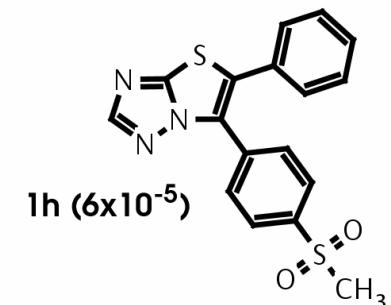
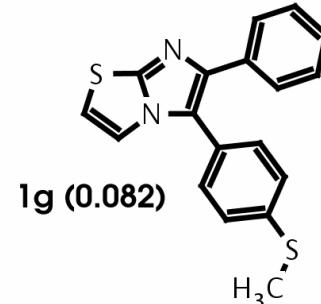
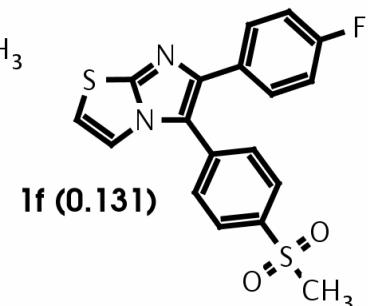
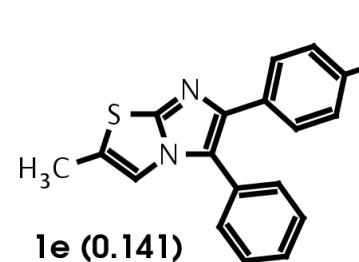
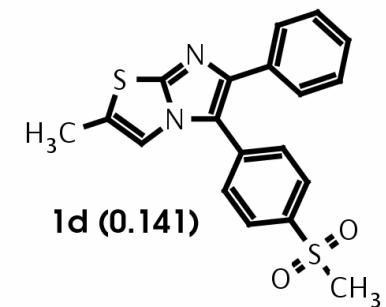
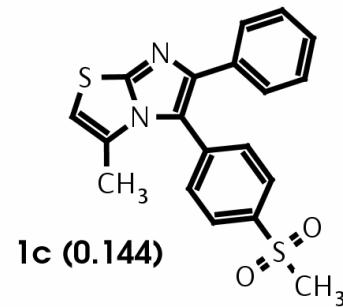
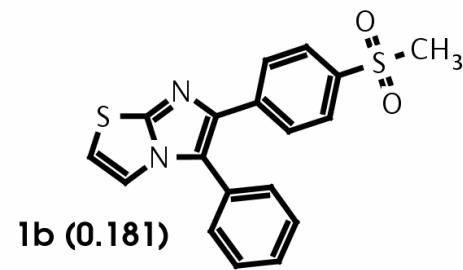
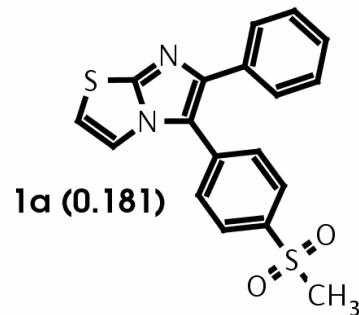


Input similarity matrix

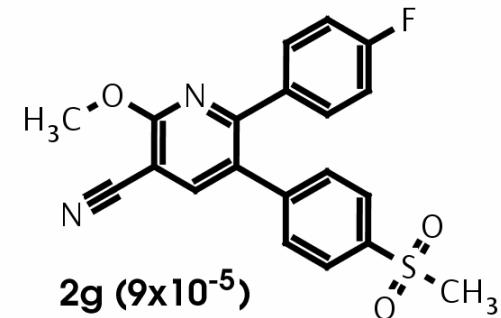
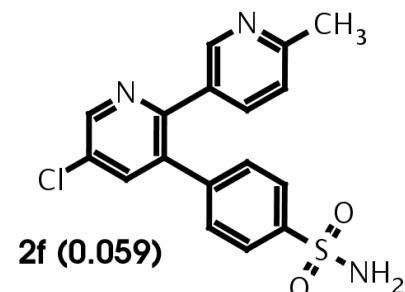
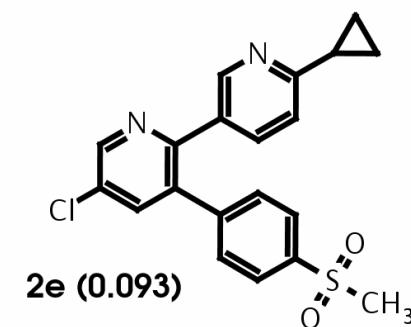
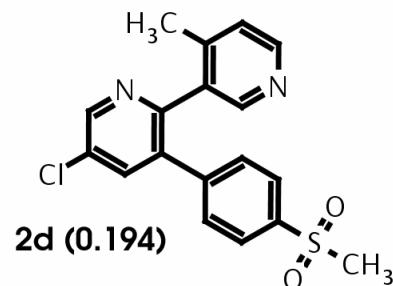
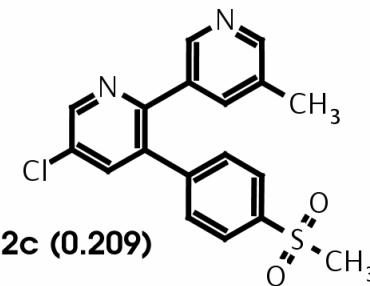
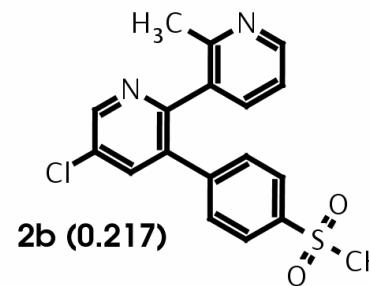
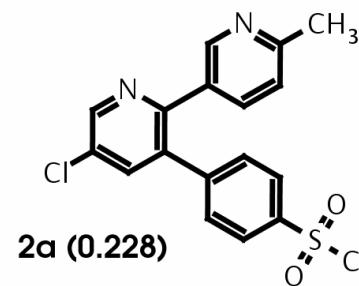


Gaussian filtering ($\alpha = 25$)

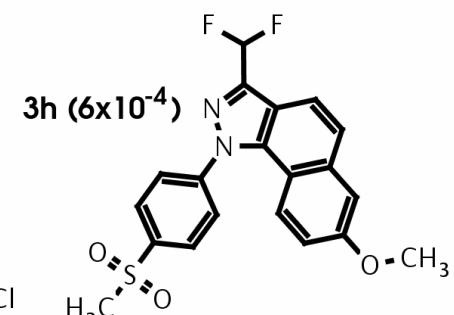
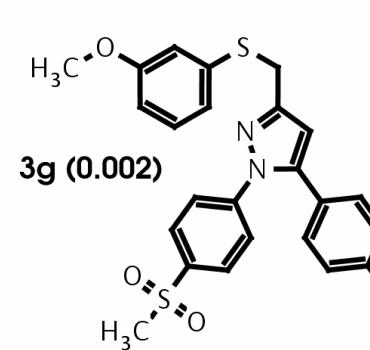
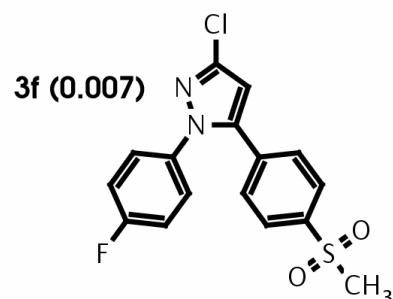
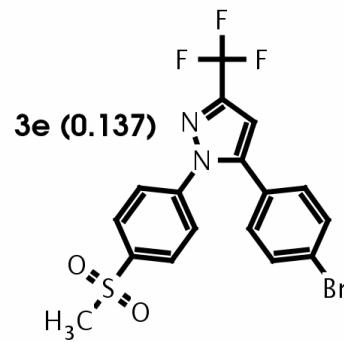
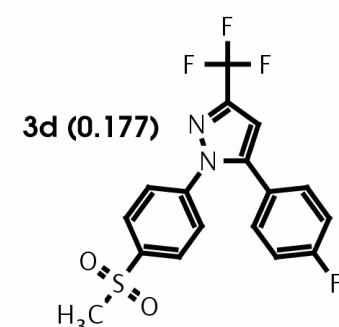
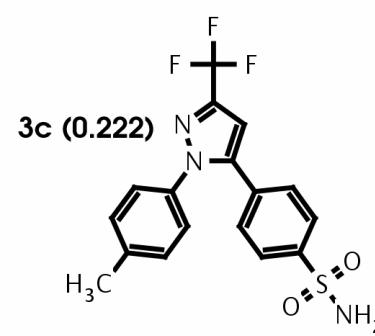
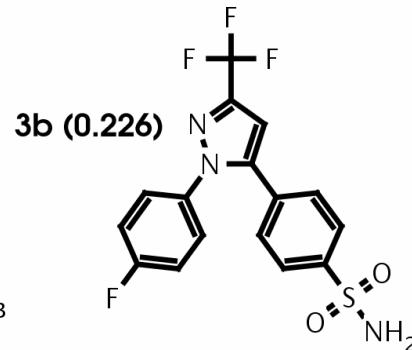
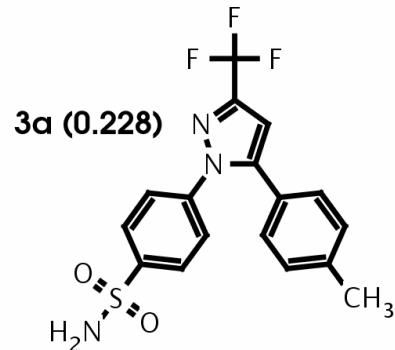
Cluster 1: Imidazothiazoles



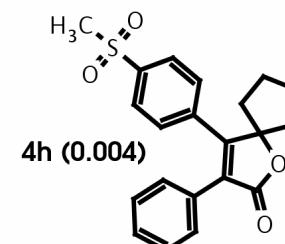
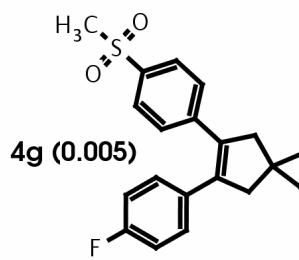
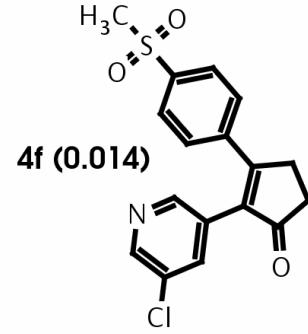
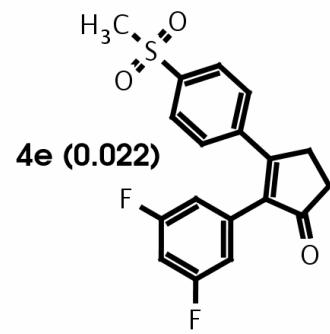
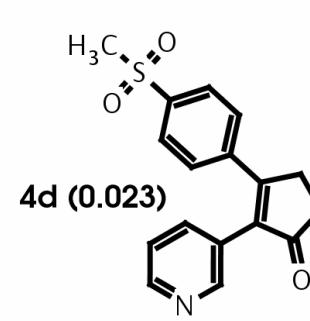
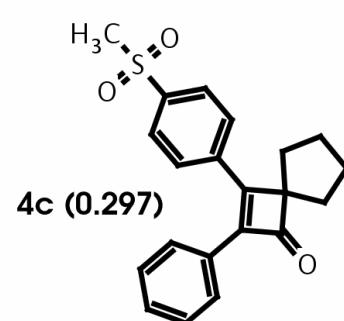
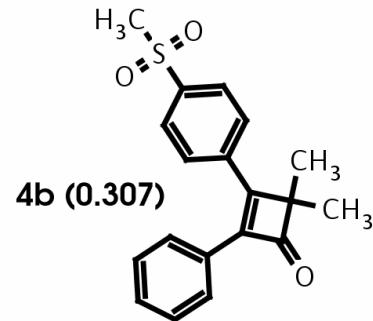
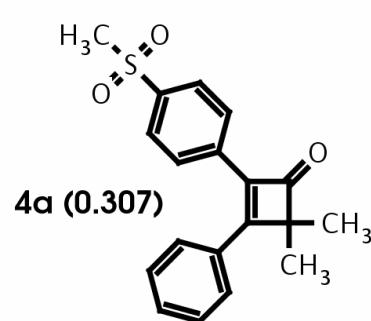
Cluster 2: Chloropyridines



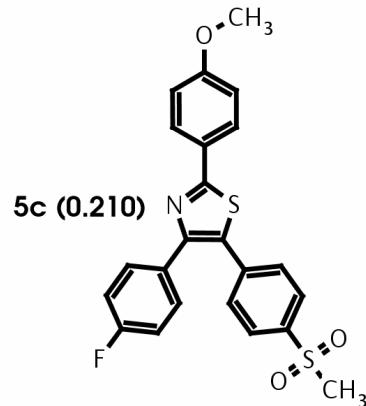
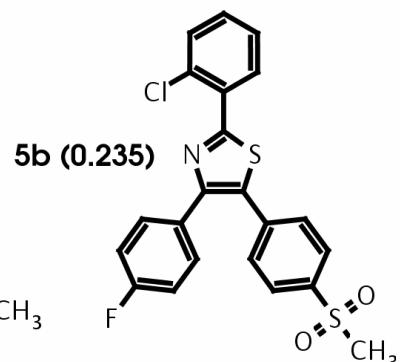
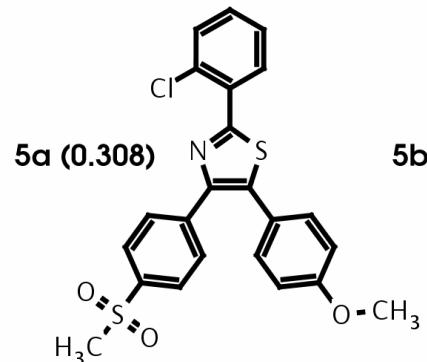
Cluster 3: Pyrazoles



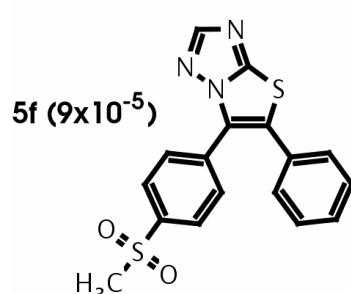
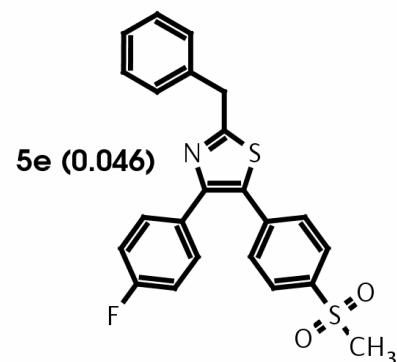
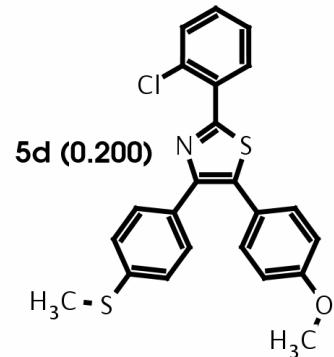
Cluster 4: Stilbenes



Cluster 5: Thiazoles



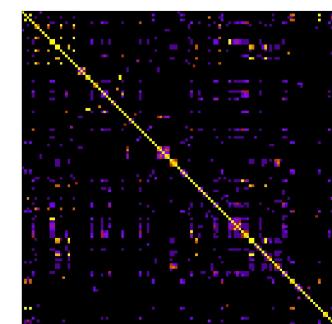
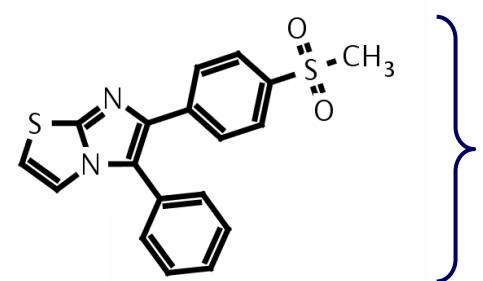
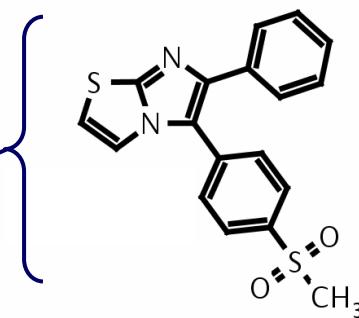
This molecule also made a minor contribution to cluster 1



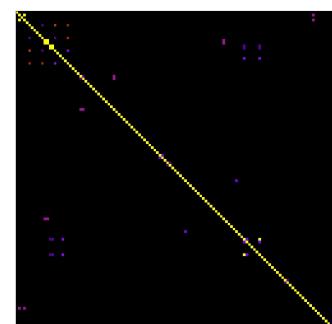
Effect of Gaussian width parameter (α)

	α					
	12.5	25	50	100	200	400
λ_1	6.367	5.586	4.752	3.806	2.917	2.272
1a	0.149	0.181	0.208	0.243	0.286	0.348
1b	0.149	0.181	0.208	0.243	0.286	0.348
1c	0.132	0.144	0.142	0.133	0.115	0.075
1d	0.131	0.141	0.139	0.134	0.119	0.097
1e	0.131	0.141	0.139	0.134	0.119	0.097
1f	0.127	0.131	0.118	0.098	0.073	0.034
1g	0.097	0.082	0.046	0.016	0.002	4×10^{-5}
1h	0.004	6×10^{-5}	3×10^{-8}	0	0	0

Normalised contributions
from molecules in cluster 1

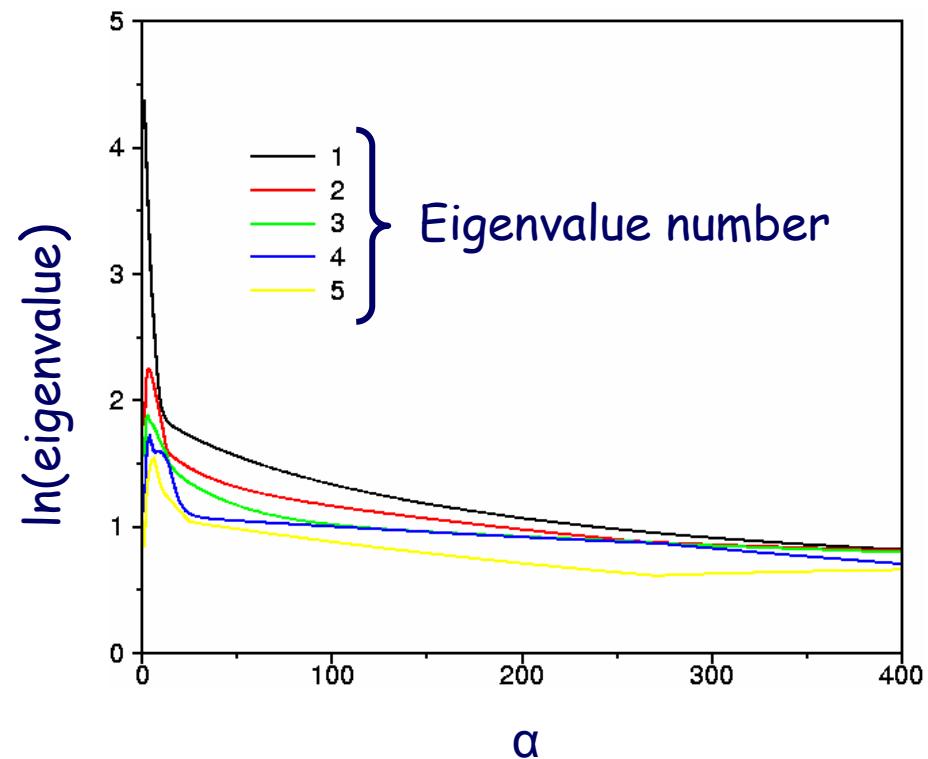


$\alpha = 12.5$

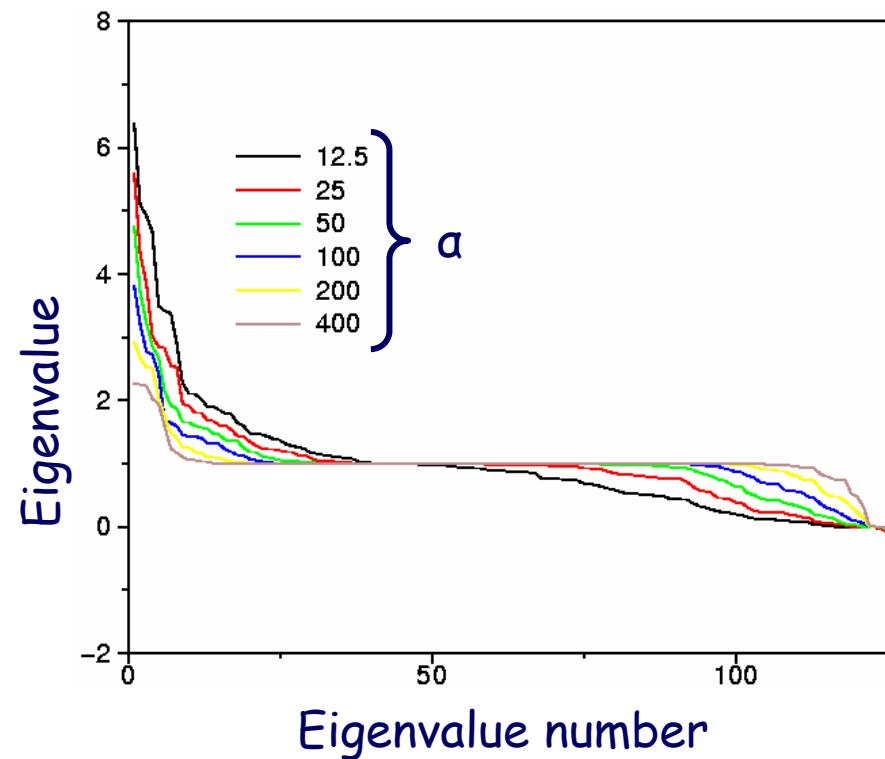


$\alpha = 400$

Effect of Gaussian width parameter (α)

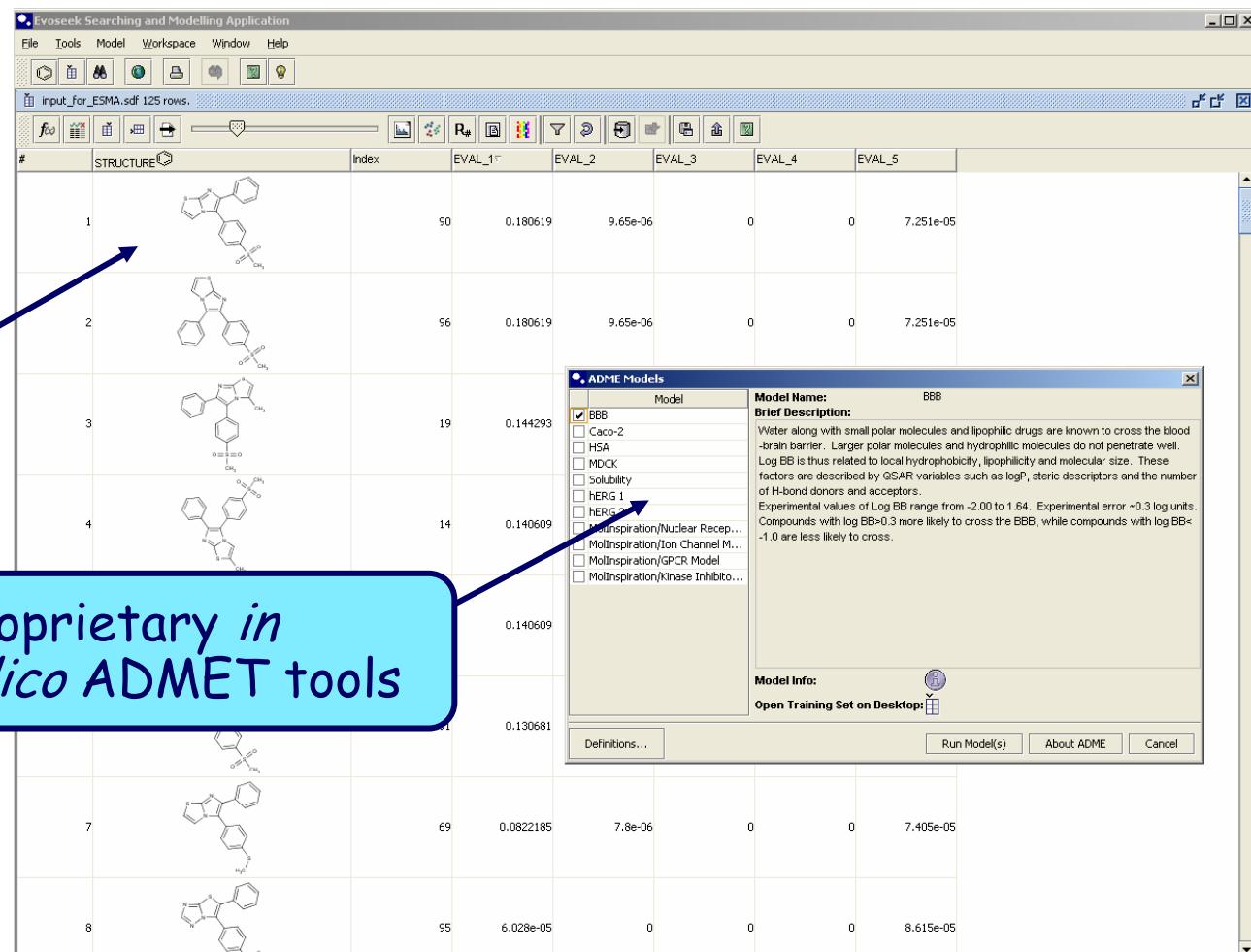


The five largest eigenvalues
as a function of α



Magnitude of eigenvalues for
different values of α

Provision of results



Applications

- Analysis of screening data
 - Prioritisation of hit series
 - Visualising different scaffolds within a dataset
 - Identifying chemically related scaffolds within and between different clusters
- Diversity assessment
 - Virtual screening hit lists
 - Virtual libraries

Future and on-going work

- Investigate alternative molecular descriptors and/or similarity methods and/or datasets and/or filtering schemes
- Investigate alternative spectral clustering methods
- Implementations which scale more favourably with the number of molecules
 - Matrix diagonalisation scales as N^3 where N is the number of molecules
- Develop the intuitive understanding of spectral clustering methods
 - A global rather than local approach to clustering

Summary

- Demonstrated the benefit of spectral clustering
 - Eigenvalues and the normalised cluster contributions (eigenvectors) provide a natural way to analyse a dataset
 - Single tunable parameter
 - Overlapping and non-overlapping clusters possible
- The clusters of molecules produced have proven effective for medicinal chemists
- Further details of this work to appear soon:
Mark Brewer, “*Development of a spectral clustering method for the analysis of molecular datasets*”, JCIM (accepted)

Acknowledgements

I would like to thank current and former colleagues at Evotec who have assisted with this work