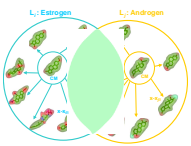
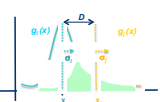


Gaussian ensemble screening (GES): A new approach to polypharmacology and virtual screening


Violeta I. Pérez-Nueno, Vishwesh Venkatraman, Lazaros Mavridis, David W. Ritchie
 Orpailleur Team, INRIA Nancy - Grand Est





LORIA (Laboratoire Lorrain de Recherche en Informatique et ses Applications),
 INRIA Nancy – Grand Est, 615 rue du Jardin Botanique, 54506 Vandœuvre-lès-Nancy, France

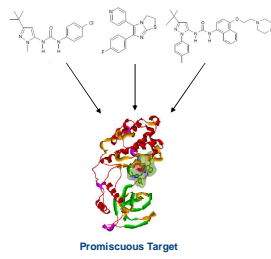
1



Polypharmacology

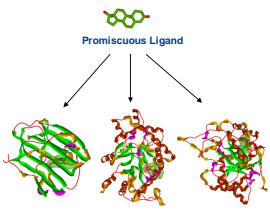
Polypharmacology (Drug selectivity)

Multiple drugs bind to a given target (promiscuous targets)




Promiscuous Target

A given drug binds to more than one target (promiscuous ligands)



Promiscuous Ligand

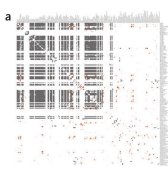
2



Previous work

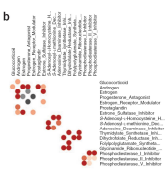
Relate receptors to each other quantitatively based on the similarity in the:

a



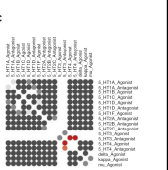
Ligand space
(chemical fingerprints)

b



Sequence space


c



Binding pocket space
(pharmacophoric descriptors)

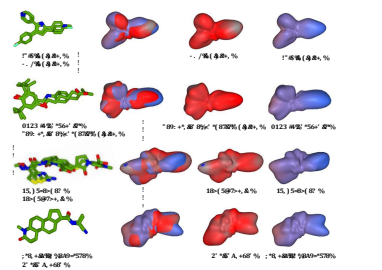
*Keiser et al. *Nature Biotechnol.* 2007, 25, 197-206. Similarity Ensemble Approach (SEA) relates proteins based on the set-wise chemical similarity among their ligands.
 *Vidal & Mestres. *Mol. Inf.* 2010, 29, 543. PHRAG, FPD, SHED molecular descriptors.
 *Weskamp et al. *Proteins* 2009, 76, 317-330. Similarity amongst binding pockets extracted by LIGSITE algorithm.
 *Milletti, F.; Vulpetti, A. *J. Chem. Inf. Model.*, 2010, 50, 1418-143. Binding pocket comparison using four-point pharmacophoric descriptors based on GRID.

3



Our approach

Gaussian Ensemble Screening (GES): 3D spherical harmonic (SH) shape-based approach which compares molecular surfaces and predicts quantitatively the relationships between drug classes very fast and efficiently.



4

Methodology

1. Calculating SH consensus shapes and center molecules
2. Ligand set representations
3. Gaussian ligand set comparisons
4. Finding the best clustering threshold
5. Gaussian p-values
6. MDDR polypharmacology interaction matrix
7. Examples of strongly related targets

5

1. Calculating spherical harmonic shapes

Surface shapes are represented as radial distance expansions of the molecular surface with respect to the center of the molecule.

- Real SHs: $y_{lm}(\theta, \phi)$
- Coefficients: a_{lm}
- Encode radial distances from origin as SH series...
- Solve coefficients by numerical integration...

$$r(\theta, \phi) = \sum_{l=0}^{15} \sum_{m=-l}^l a_{lm} y_{lm}(\theta, \phi)$$

Ritchie, D.W. and Kemp, G.J.L. *J. Comp. Chem.* 1999, 20, 383-395.

6

2. Calculating SH consensus shapes and center molecules

$$\bar{r}(\theta, \phi) = \frac{1}{N} \sum_{k=1}^N \sum_{l=0}^L \sum_{m=-l}^l a_{lm}^{(k)} y_{lm}(\theta, \phi)$$

"Consensus" shape

Pérez-Nueno et al. *J. Chem. Inf. Model.* 2008, 48, 2146-2165.

7

3. Ligand set representations

The idea is to represent a cluster of molecules as a Gaussian distribution with respect to a selected centre molecule (CM).

- Calculate SH molecular surfaces of each ligand in each ligand set and superpose them.
- Calculate the center molecule (CM) of the ligand set and the normalised SH distance (1-Similarity Score) between that of the CM and each cluster member.
- Assuming that these distances follow a Gaussian distribution, each cluster may be represented as a probability density function $g_i(x)$

L_i : Estrogen

$$g_i(x) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \cdot e^{-\frac{(x-x_i)^2}{2\sigma_i^2}}$$

σ : SD of the member distances

An illustration of a Gaussian ligand set cluster.

8

4. Gaussian ligand set comparisons

By considering the SD of the member distances as the Gaussian width of a distribution, we calculate a "distance" (D) between two clusters, i and j , and normalizing the distance term we can write it as a Hodgkin-like similarity score S_{ij} between two distributions.

$$S_{ij} = \frac{2 \int g_i(x) g_j(x) dx}{\int g_i(x)^2 dx + \int g_j(x)^2 dx}$$

$$S_{ij} = \frac{2 \sqrt{\frac{a \cdot b}{a+b}} \cdot e^{-\frac{ab}{a+b} \left(\frac{x_i - x_j}{D}\right)^2}}{a^2 + b^2}$$

$a = 1/2\sigma_i^2$
 $b = 1/2\sigma_j^2$
 x_i, x_j : distance between the CMs of clusters i and j

L_i : Estrogen
 L_j : Thrombin

$S_{ij} = 1.33! \cdot 10^{-41}$

Illustration of the very small Gaussian overlap between the estrogen and thrombin ligand sets.

9

4. Gaussian ligand set comparisons

The similarity between drug classes can be calculated rapidly and reliably by calculating the Gaussian overlap between pairs of such clusters.

Thus, it is straight-forward to calculate all-against-all cluster comparisons. It is worth noting that our cluster similarity score depends only on the similarity of pairs of centre molecules and the SDs of their respective clusters. It does not depend on the number of members of each cluster.

L_i : Estrogen
 L_j : Androgen

$S_{ij} = 0.57$

Illustration of the large Gaussian overlap between the estrogen and androgen ligand sets.

10

4. Gaussian ligand set comparisons

1. MDDR ANNOTATION FAMILIES SPACE

THERAPEUTIC ANNOTATION

We applied the approach to 270 specific therapeutic annotations in MDDR. Ligands which share an annotation define a set of functionally related molecules which we call a "ligand set". MDDR annotations are quite general and were primarily derived from the patent literature. A given annotation may thus contain a diverse set of compounds with a wide range of affinities.

11

4. Gaussian ligand set comparisons

2. MDDR ANNOTATION SHAPE CLUSTERS

ANNOTATION CLUSTER

In order to eliminate outliers, we used the CAST clustering algorithm to cluster the members of each annotation using their PARAFIT Tanimoto similarity scores. We then calculated the consensus shape and the center molecule for each cluster, and we eliminated any cluster members beyond 1.5 standard deviations (SDs) from the corresponding CM.

12

5. Finding the best clustering threshold

2. MDDR ANNOTATION SHAPE CLUSTERS

We clustered each annotation according to Parafit Shape Tanimoto using different similarity thresholds: 0.6, 0.65, 0.675, 0.7, 0.8, 0.85. Each ligand set was randomly split into two almost equally sub-clusters, and all-vs-all clustering was performed with the aim of split and reassemble the split clusters correctly.

13

5. Finding the best clustering threshold

3. SPLIT ANNOTATION SHAPE CLUSTERS + GAUSSIAN SCORING

Here are shown the C1 of different annotations split in two groups to obtain the distribution of scores for the true cases, where annotations are related to each other (L₁C1_A vs L₁C1_B, L₂C1_A vs L₂C1_B ...) and the false cases, where the annotations are not related (L₁C1_A vs L₂C1_A, L₁C1_A vs L₃C1_A ...).

If we can split and reassemble clusters of molecules that we know they are related, then we can identify interesting relationships between clusters of molecules that we don't know they are related

14

5. Finding the best clustering threshold

4. GAUSSIAN SCORES ALL VS ALL RANK & ROC

	L ₁ C1 _A	L ₁ C1 _B	L ₁ C2 _A	L ₁ C2 _B	L ₁ C3 _A	L ₁ C3 _B	L ₂ C1 _A	L ₂ C1 _B	L ₃ C1 _A	L ₃ C1 _B	...
TRUE											
L ₁ C1 _A		S _{C1A,C1B}	S _{C1A,C2A}	S _{C1A,C2B}	S _{C1A,C3A}	S _{C1A,C3B}					
L ₁ C1 _B	S _{C1B,C1A}		S _{C1B,C2A}	S _{C1B,C2B}	S _{C1B,C3A}	S _{C1B,C3B}					
REST											
L ₁ C2 _A	S _{C2A,C1A}	S _{C2A,C1B}			S _{C2A,C3A}	S _{C2A,C3B}					
L ₁ C2 _B											
L ₁ C3 _A											
L ₁ C3 _B											
L ₂ C1 _A											
L ₂ C1 _B											
L ₃ C1 _A											
L ₃ C1 _B											

SCORE	LIGAND SET 1	LIGAND SET 2	
S _{L1C1A,L1C1B} 0.999	L ₁ C1 _A	L ₁ C1 _B	T 1
S _{L1C2A,L1C2B} 0.998	L ₁ C2 _A	L ₁ C2 _B	T 1
S _{L2C1A,L3C1A} 0.854	L ₂ C1 _A	L ₃ C1 _A	F 0
S _{L1C1B,L4C1B} 0.153	L ₁ C1 _B	L ₄ C1 _B	F 0

We produce a matrix of Gaussian Overlap Scores for **true** target classes (members of the same annotation cluster) and the **rest** (members supposed not to be related).

15

5. Finding the best clustering threshold

ROC curves obtained for the range of similarity thresholds of 0.6, 0.65, 0.675, 0.7, 0.8, 0.85. Using a PARAFIT shape Tanimoto value of 0.65 gave the best early performance AUC(5%,10%). Hence, 0.65 was chosen as the appropriate for shape-based clustering.

16

6. Gaussian p-values

In order to transform a list of cluster similarity scores into a more meaningful list of probabilities, a statistical model was developed.

Each Gaussian similarity score was transformed into a probability value, or "p-value", from the observed distribution of scores. For a Gaussian distribution, it can be shown that the probability of finding at random from the distribution some value X greater than a given value x is given by:

$$p(X > x) = \int_x^{\infty} f(t) dt = \text{erfc}(x)$$

$$p(S > s) = \text{erfc}\left(\frac{S}{\sqrt{2}\sigma}\right)$$

where $f(t)$ is the standard normalized Gaussian probability density function and $\text{erfc}(x)$ is the complementary error function. For a normalized distribution of scores, we obtain $p(S > s)$

In our real data we can see that all the false pairs appear at the beginning of the list. We are looking at the really top end of the distribution which tells us how statistically significant is the result.

We used the R package to fit the distribution of all pair-wise Gaussian overlap scores to a Gaussian function with $\text{erfc}(x)$.

A p-value was calculated analytically from the scores distribution for any given pair-wise score using the $\text{erfc}(x)$.

A "p-value" for a given score, s , is the probability of finding at random from the distribution some other score, S , which is greater than s .

7. MDDR polypharmacology interaction matrix

MDDR polypharmacology interaction matrix for the top 50 ligand set relationships found from the all against all comparison of the 270 MDDR specific annotations

8. Nuclear hormone receptors

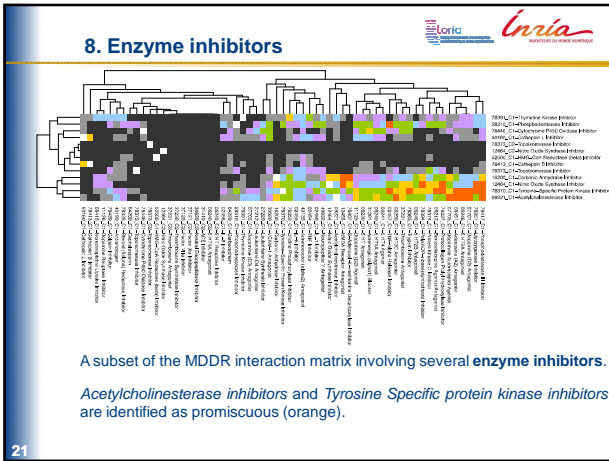
A subset of the MDDR interaction matrix involving several nuclear hormone receptors.

Antiglucocorticoids and progesterone antagonists are identified as promiscuous (orange).

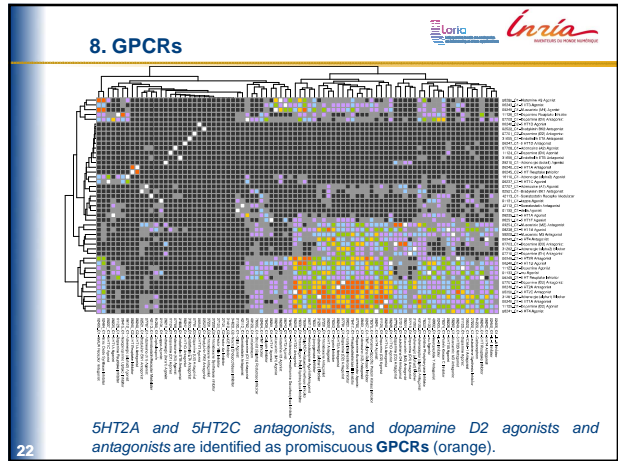
8. Serine proteases

A subset of the MDDR interaction matrix involving several serine proteases.

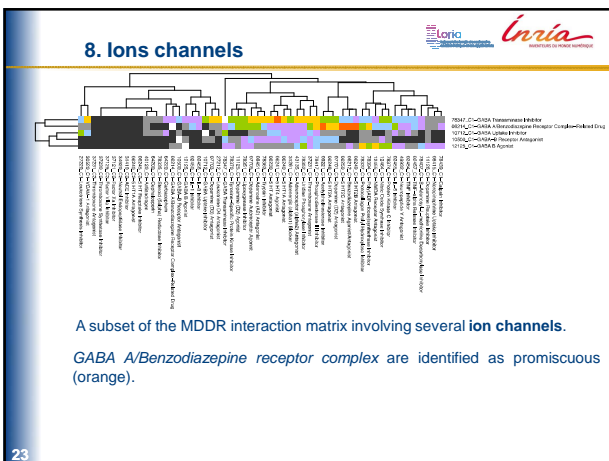
Coagulation factors Xa and VIIa inhibitors are identified as promiscuous, as well as trypsin inhibitors (orange).



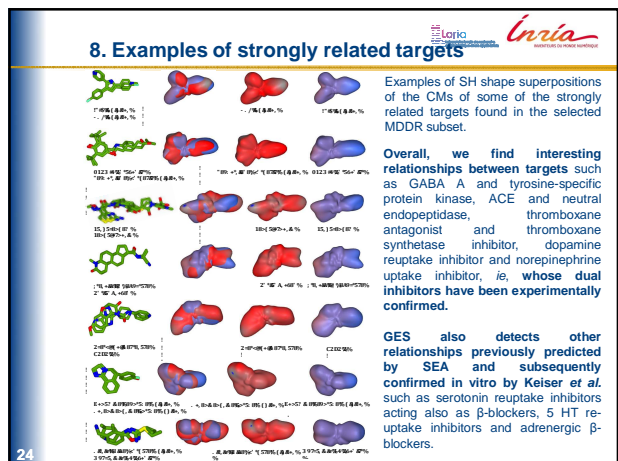
21



22



23



24

Conclusions

- We have presented a new 3D shape-based approach for predicting and quantifying drug promiscuity by correlating Gaussian clusters of ligand SH shapes.
- The method has been validated using drug ligand sets of the MDDR and has been demonstrated to be effective in identifying drug families which are known to have related MDDR activity classes.
- Our results show that GES provides an efficient way to measure the similarity between clusters of arbitrary numbers of members.
- The examples shown in this study demonstrate that GES is a useful way to study polypharmacology relationships, and it could provide a novel way to propose new targets for drug repositioning.

25

Acknowledgements

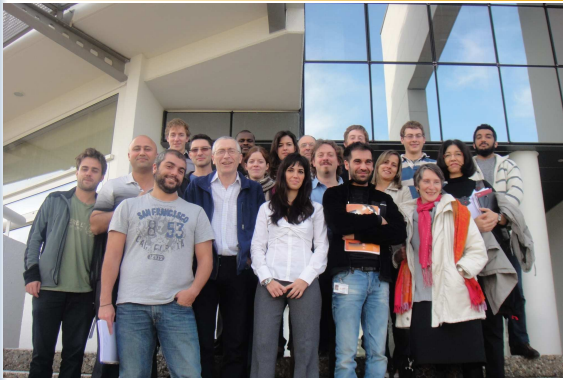
- INRIA Nancy - Grand Est
- FP7 Marie Curie IEF Fellowship (DOVSA 254128)

Papers: <http://www.loria.fr/~pereznie/>
<http://www.loria.fr/~ritchied/>

ParaSurf + ParaFit: <http://www.ceposinsilico.de/>

26

Thank you!



27

Comparison with SEA approach

Cluster	SEA	GES	SEA	GES	SEA	GES
1	1	1	1	1	1	1
2	1	1	1	1	1	1
3	1	1	1	1	1	1
4	1	1	1	1	1	1
5	1	1	1	1	1	1
6	1	1	1	1	1	1
7	1	1	1	1	1	1
8	1	1	1	1	1	1
9	1	1	1	1	1	1
10	1	1	1	1	1	1
11	1	1	1	1	1	1
12	1	1	1	1	1	1
13	1	1	1	1	1	1
14	1	1	1	1	1	1
15	1	1	1	1	1	1
16	1	1	1	1	1	1
17	1	1	1	1	1	1
18	1	1	1	1	1	1
19	1	1	1	1	1	1
20	1	1	1	1	1	1
21	1	1	1	1	1	1
22	1	1	1	1	1	1
23	1	1	1	1	1	1
24	1	1	1	1	1	1
25	1	1	1	1	1	1
26	1	1	1	1	1	1
27	1	1	1	1	1	1
28	1	1	1	1	1	1
29	1	1	1	1	1	1
30	1	1	1	1	1	1
31	1	1	1	1	1	1
32	1	1	1	1	1	1
33	1	1	1	1	1	1
34	1	1	1	1	1	1
35	1	1	1	1	1	1
36	1	1	1	1	1	1
37	1	1	1	1	1	1
38	1	1	1	1	1	1
39	1	1	1	1	1	1
40	1	1	1	1	1	1
41	1	1	1	1	1	1
42	1	1	1	1	1	1
43	1	1	1	1	1	1
44	1	1	1	1	1	1
45	1	1	1	1	1	1
46	1	1	1	1	1	1
47	1	1	1	1	1	1
48	1	1	1	1	1	1
49	1	1	1	1	1	1
50	1	1	1	1	1	1
51	1	1	1	1	1	1
52	1	1	1	1	1	1
53	1	1	1	1	1	1
54	1	1	1	1	1	1
55	1	1	1	1	1	1
56	1	1	1	1	1	1
57	1	1	1	1	1	1
58	1	1	1	1	1	1
59	1	1	1	1	1	1
60	1	1	1	1	1	1
61	1	1	1	1	1	1
62	1	1	1	1	1	1
63	1	1	1	1	1	1
64	1	1	1	1	1	1
65	1	1	1	1	1	1
66	1	1	1	1	1	1
67	1	1	1	1	1	1
68	1	1	1	1	1	1
69	1	1	1	1	1	1
70	1	1	1	1	1	1
71	1	1	1	1	1	1
72	1	1	1	1	1	1
73	1	1	1	1	1	1
74	1	1	1	1	1	1
75	1	1	1	1	1	1
76	1	1	1	1	1	1
77	1	1	1	1	1	1
78	1	1	1	1	1	1
79	1	1	1	1	1	1
80	1	1	1	1	1	1
81	1	1	1	1	1	1
82	1	1	1	1	1	1
83	1	1	1	1	1	1
84	1	1	1	1	1	1
85	1	1	1	1	1	1
86	1	1	1	1	1	1
87	1	1	1	1	1	1
88	1	1	1	1	1	1
89	1	1	1	1	1	1
90	1	1	1	1	1	1
91	1	1	1	1	1	1
92	1	1	1	1	1	1
93	1	1	1	1	1	1
94	1	1	1	1	1	1
95	1	1	1	1	1	1
96	1	1	1	1	1	1
97	1	1	1	1	1	1
98	1	1	1	1	1	1
99	1	1	1	1	1	1
100	1	1	1	1	1	1

28