

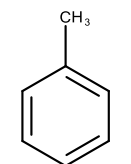
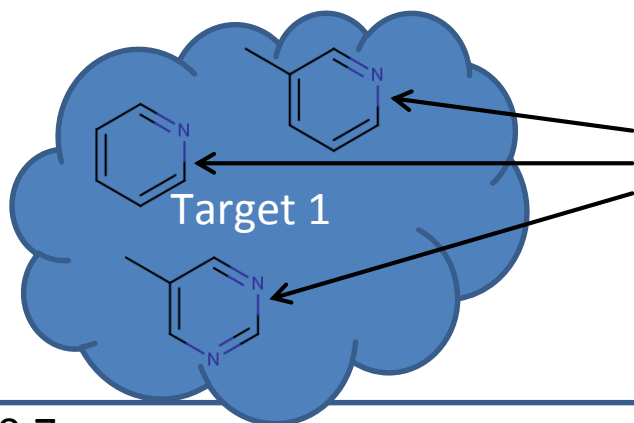
# OCEAN

## Optimized Cross rEActivity estimation

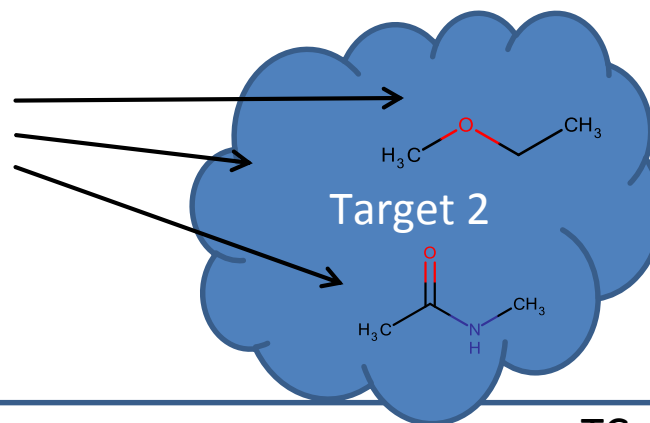
Wolf-Guido Bolick, Paul Czodrowski / RDKit UserGroupMeeting

# OCEAN?

## Re-Implentation of SEA (similarity ensemble approach) by Keiser et al.



## Query



TC:	0.7				
	0.6				T1:
	<u>0.4</u>	z-Score	p-Value	e-Value	
$\Sigma$	1.7	$\rightarrow$ 2.5	$\rightarrow$ 0.01	$\rightarrow$ 0.02	

					TC:	0.6
						0.4
T2:						
e-Value	p-Value	z-Score				0.3
0.22	0.11	1.2	←	Σ		1.3

## How to get to the e-value

$$Score_{Raw} = \sum_{i=0}^n TC[i] > Threshold \quad (1)$$

$$Score_Z = \frac{\mu(Score_{Raw}) - Score_{Raw,expected}}{\sigma(Score_{Raw,expected})} \quad (2)$$

$$x = -\exp\left(\frac{-Score_Z * \pi}{\sqrt{6} - 0,57721}\right) \quad (3)$$

$$Value_P = -\exp\left(\frac{x+x^2}{2} + \frac{x^3}{6}\right) \quad (4)$$

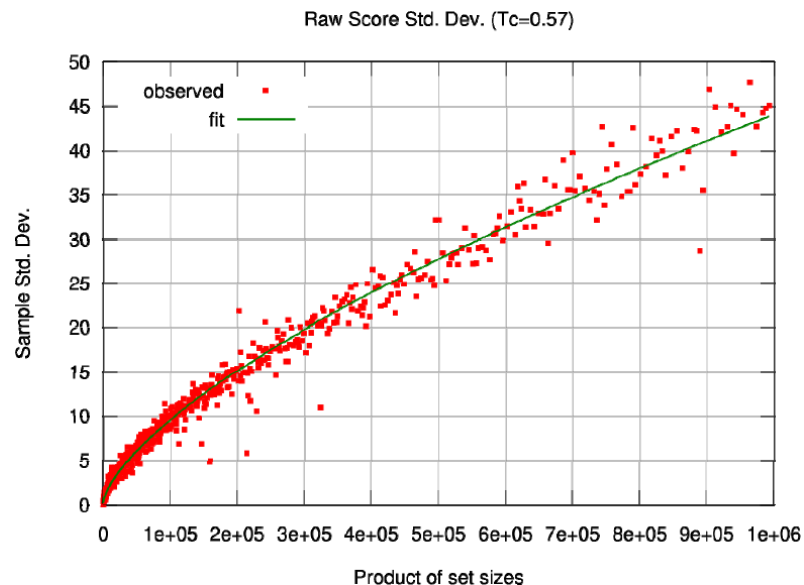
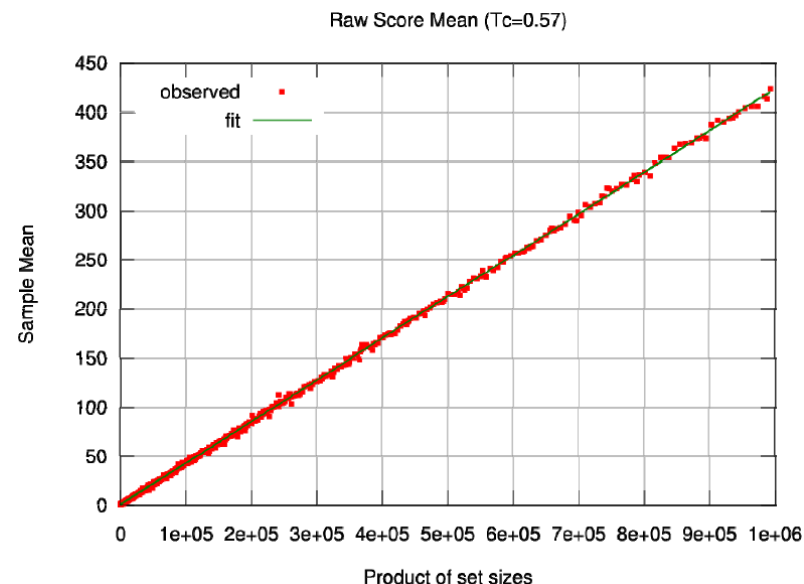
$$Value_E = Value_P * N_{Targets} \quad (5)$$

$$Score_{Raw} = \sum_{i=0}^n TC[i] > Threshold \quad (1)$$

## Background of SEA

Determination of the expected mean average by means of large-scale random molecular comparisons

Standard deviations of the molecular comparison

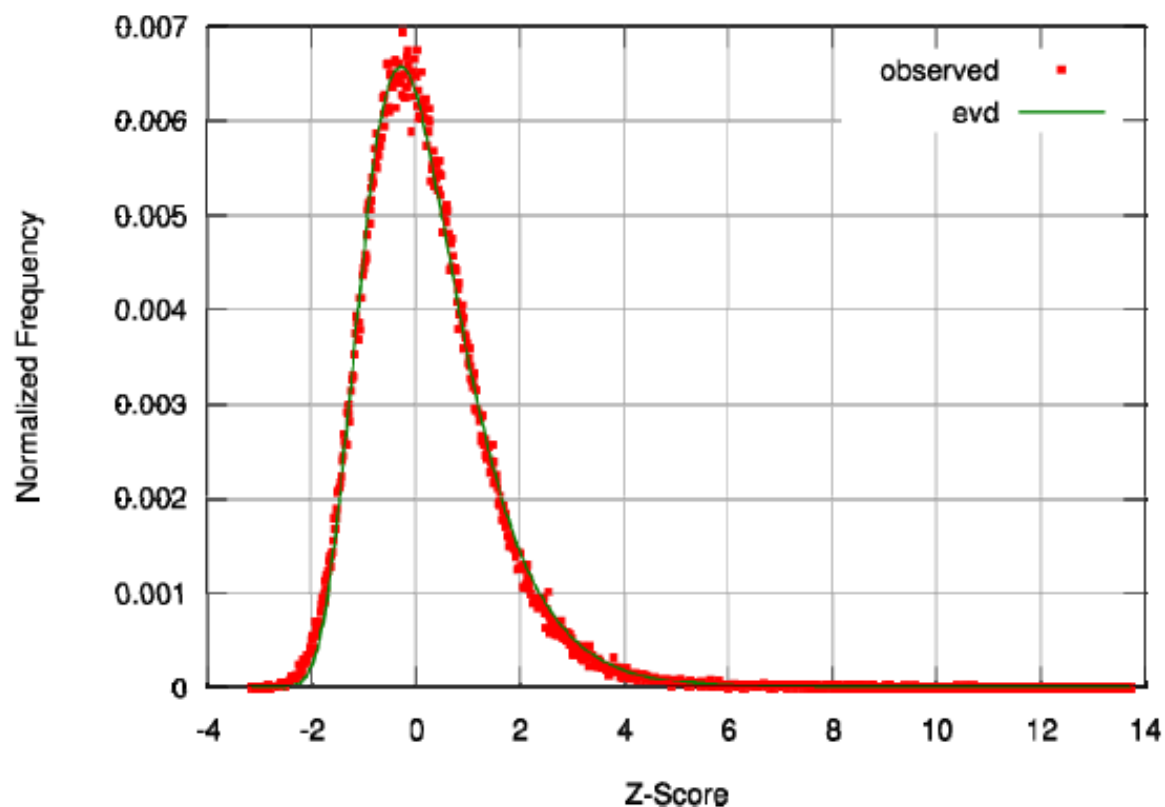


$$Score_Z = \frac{\mu(Score_{Raw}) - Score_{Raw,expected}}{\sigma(Score_{Raw,expected})} \quad (2)$$

## Background of SEA

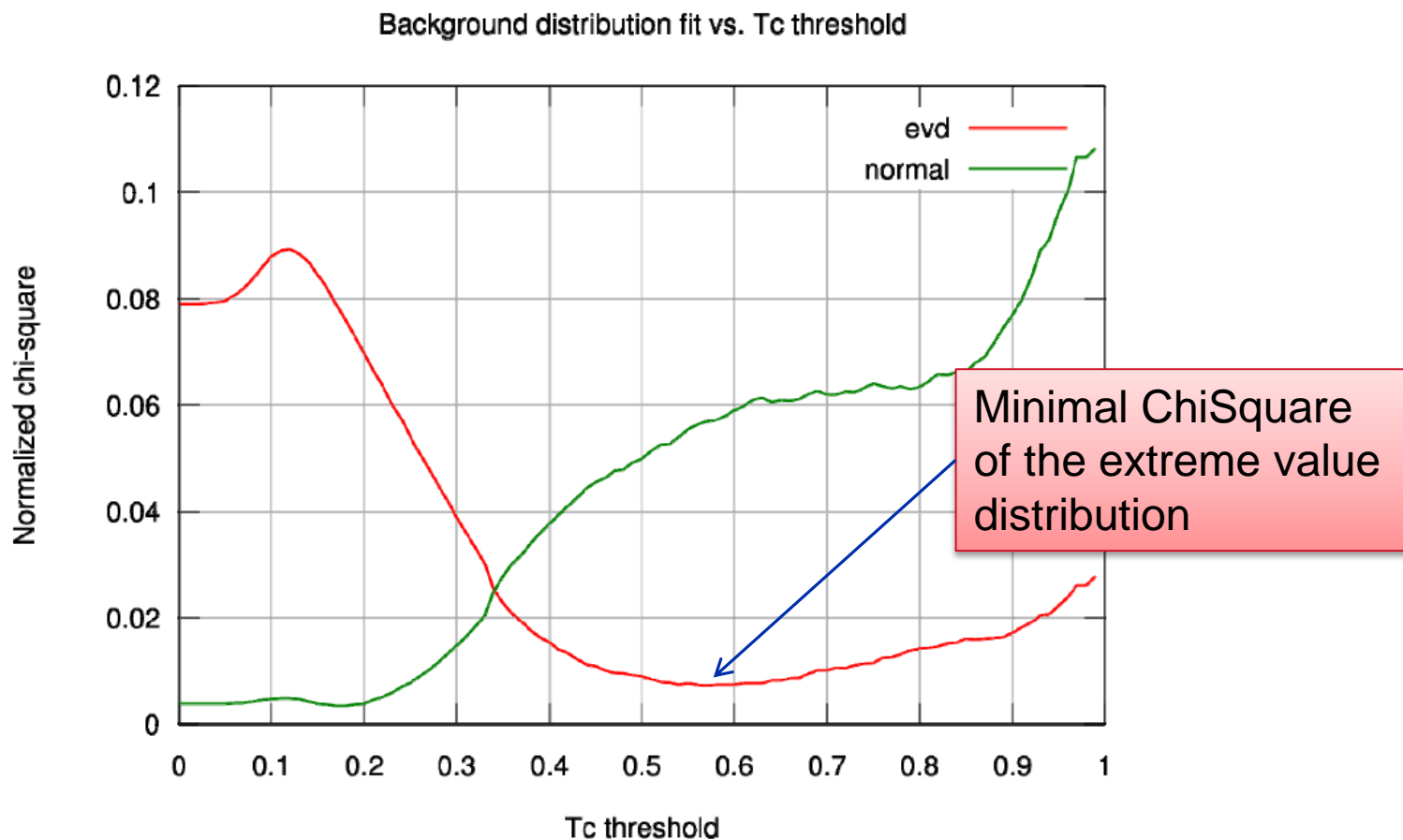
Calculation of z-Scores of all comparisons

z-Score distribution (Tanimoto = 0.57)



# Choice of the optimal Tanimoto coefficient

Fitting of normal and extreme value distribution of the z-scores



# Optimal Tanimoto coefficient

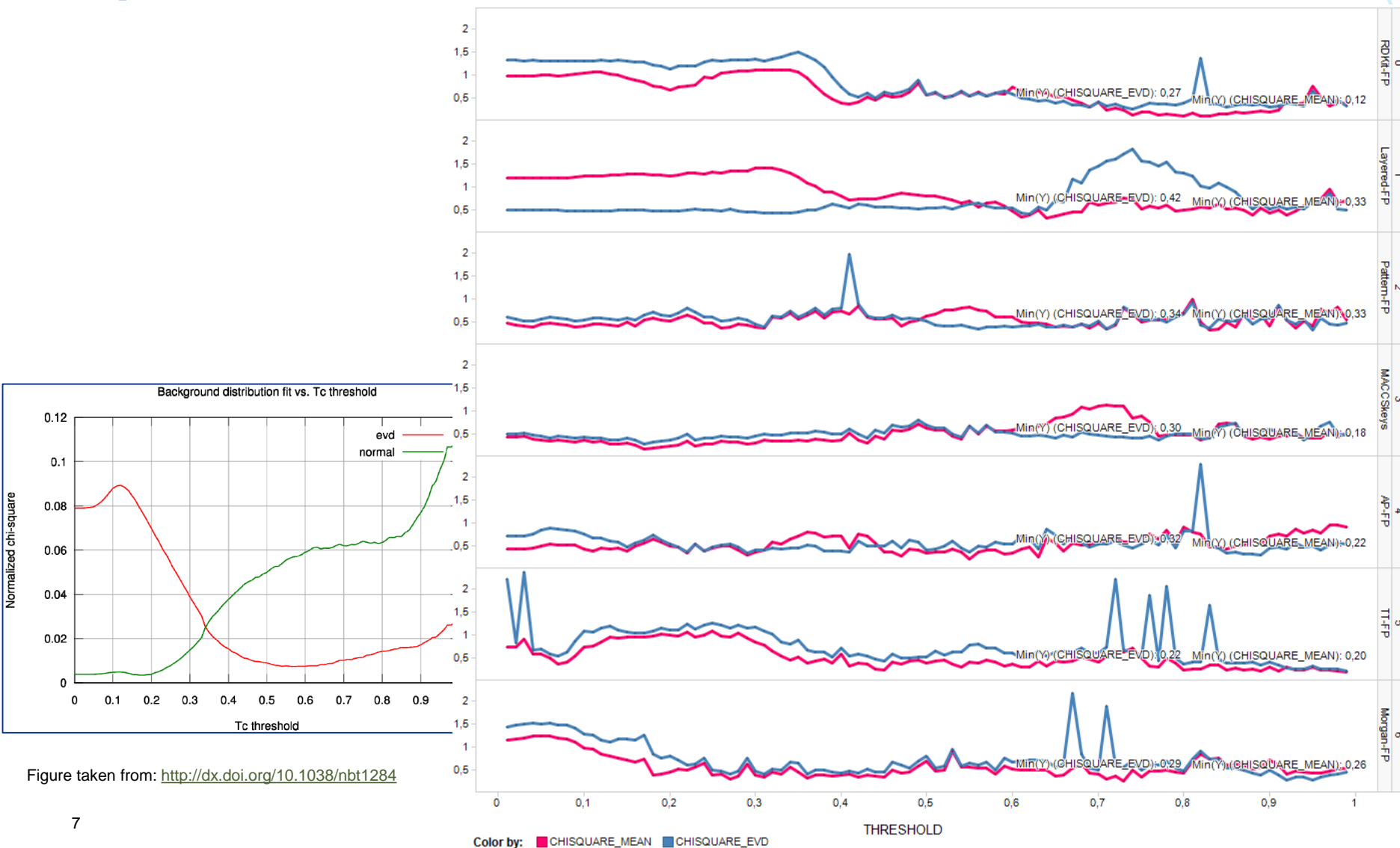


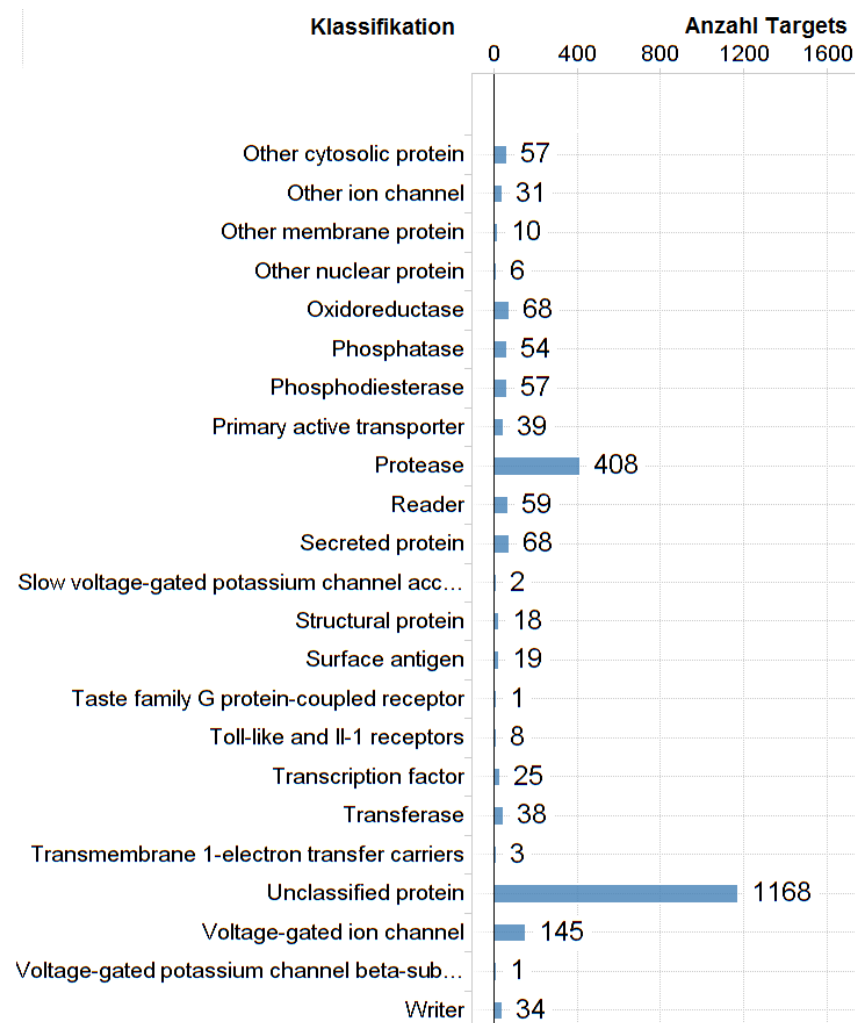
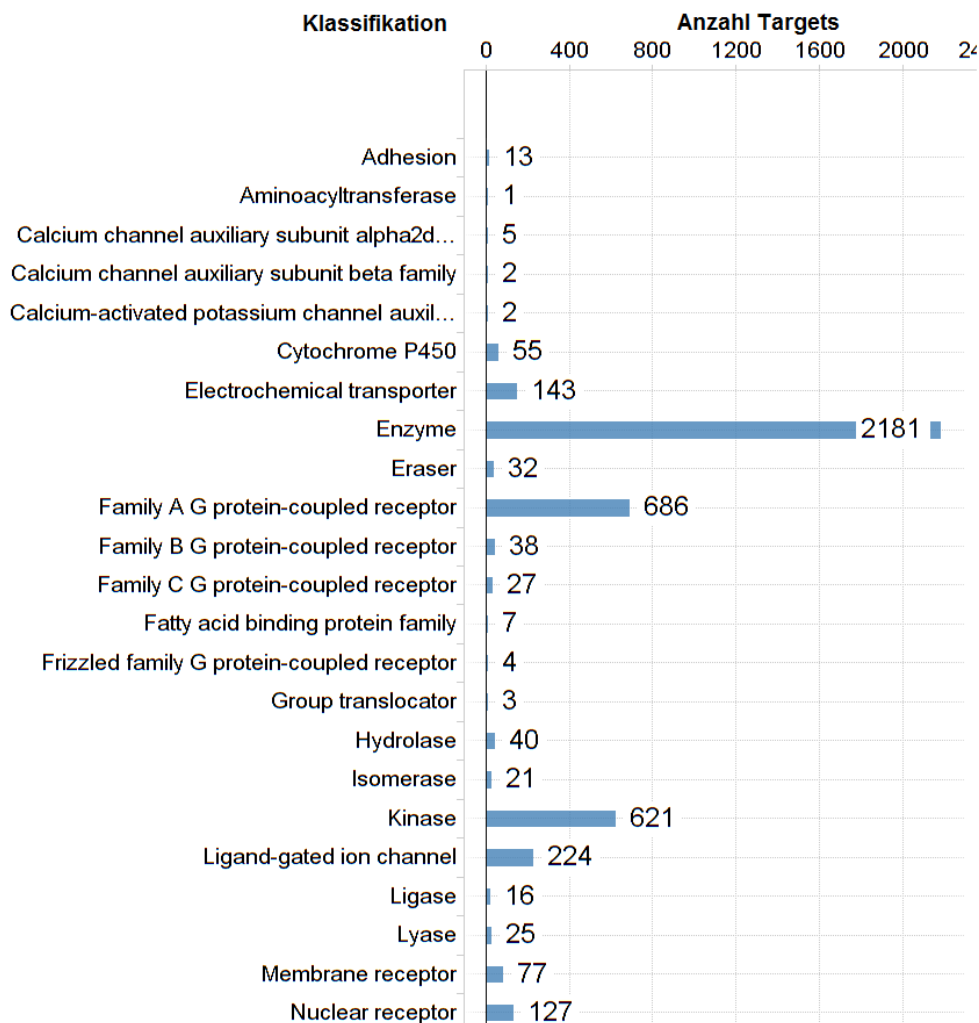
Figure taken from: <http://dx.doi.org/10.1038/nbt1284>

## ChEMBL20 filtering

- $IC_{50}$  and  $K_i$
- Activity  $\leq 10 \mu M$
- $\geq 10$  Compounds per Target
- Target Type: Single Protein or Protein Complex
- Homo sapiens
- If multiple measurements: only if single  $IC_{50}$  within mean  $\pm 3$  std\_dev



# Target Classes



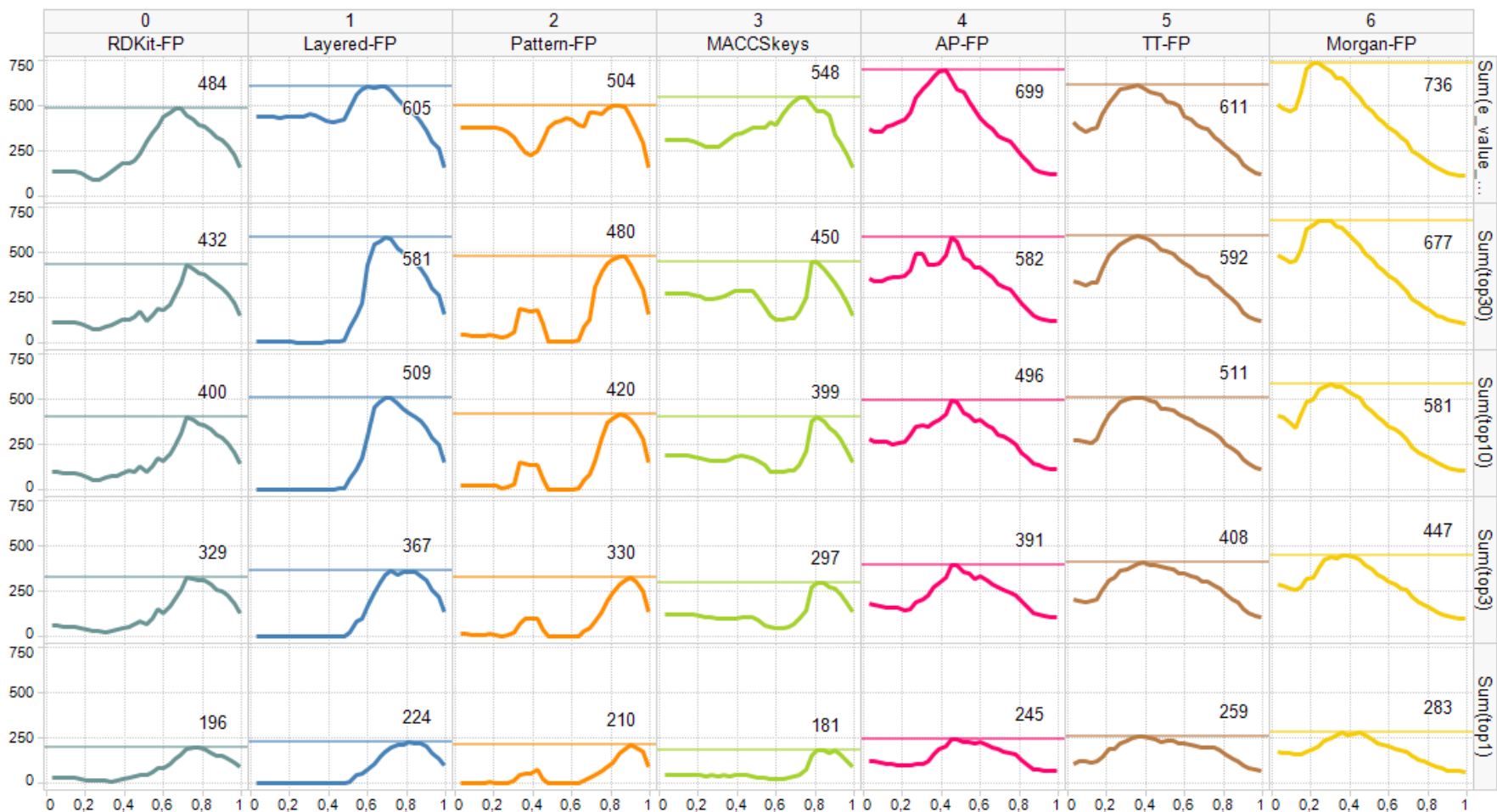
## Validation set

- 1,885 drugs on the market extracted from ChEMBL17
  - For 928 drugs, no drug annotated in OCEAN-DB
    - ➔ 957 drugs with annotated target in OCEAN-DB

How does OCEAN perform in terms of target prediction for this dataset?

[The drug itself was removed for the OCEAN search]

# Tanimoto scan based on the validation set



# Validation

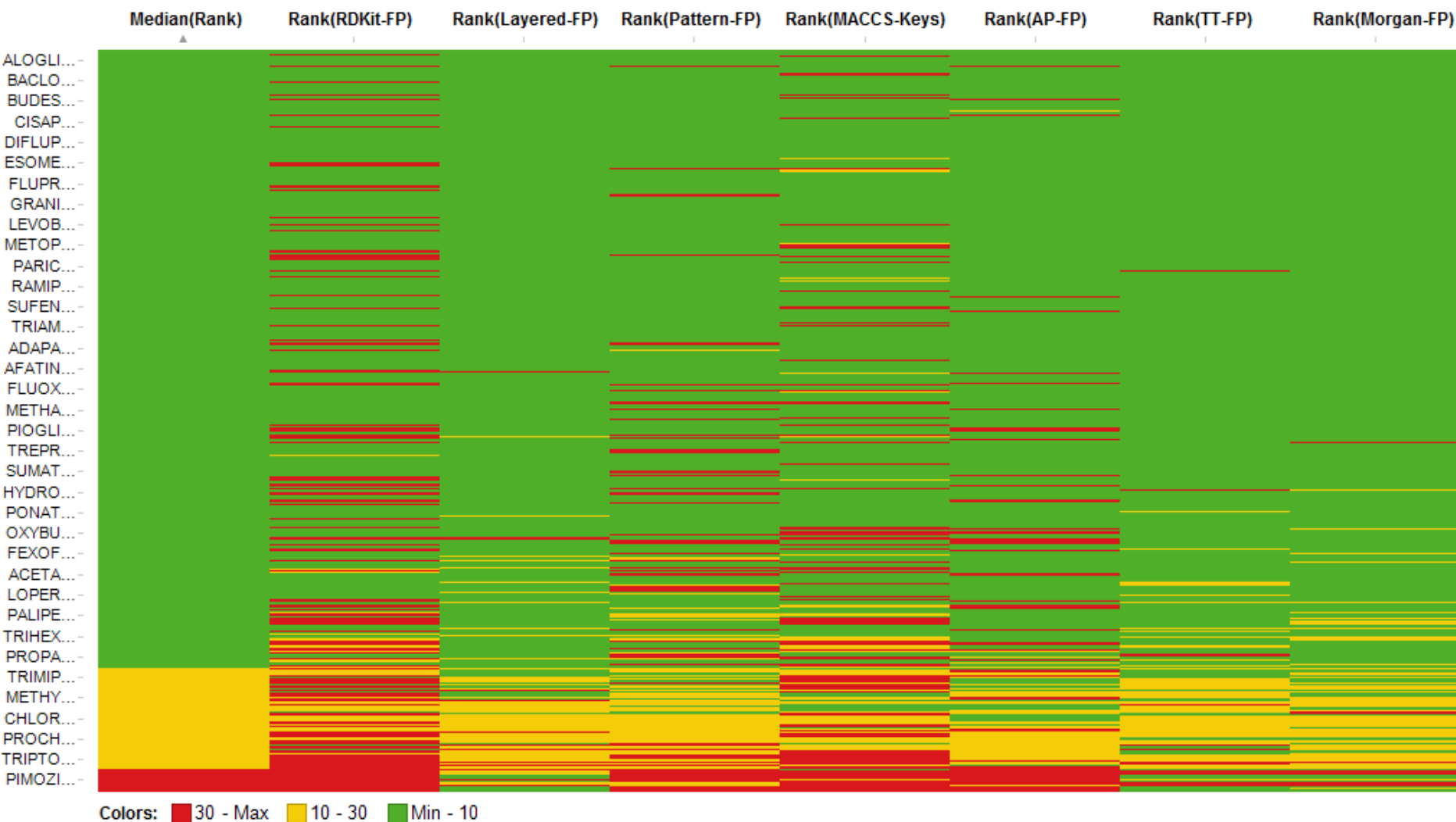
## SEA

Fingerprint	Threshold	Sum(e_value<1)	Sum(Top30)	Sum(Top10)	Sum(Top3)	Sum(Top1)
RDKit-FP	0,75	427	412	385	318	196
Layered-FP	0,33	444	3	1	0	0
Pattern-FP	0,60	423	8	1	1	1
MACCSkeys	0,75	548	250	213	141	74
AP-FP	0,39	693	442	388	298	189
TT-FP	0,99	123	120	115	105	68
Morgan-FP	0,96	115	112	107	100	65

## OCEAN

Fingerprint	Threshold	Sum(e_value<1)	Sum(Top30)	Sum(Top10)	Sum(Top3)	Sum(Top1)
RDKit-FP	0,72	448	432	400	329	187
Layered-FP	0,69	604	581	509	339	177
Pattern-FP	0,84	499	480	420	291	169
MACCSkeys	0,81	469	450	399	297	181
AP-FP	0,45	647	582	496	391	245
TT-FP	0,36	611	592	511	404	259
Morgan-FP	0,30	690	674	581	438	251

# Validation for the drugs

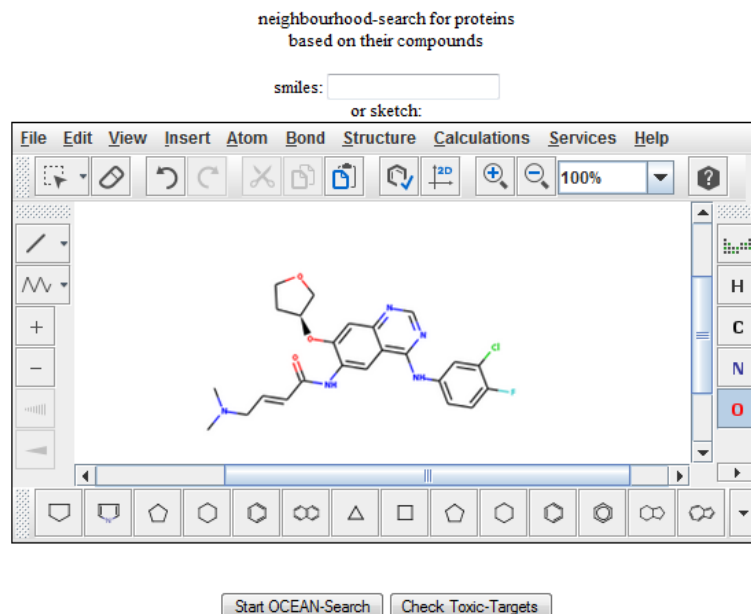


# Validation for the target classes



# Front-End

## Optimized Cross rEActivity estimation



- [ocean Home](#)
- [About ocean](#)

• © Wolf-Guido Bolick, Paul Czodrowski - Global Computational Chemistry - 2015-02-26 (v 0.5)

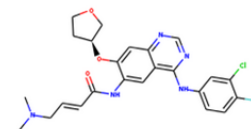
• Questions? Comments? Please send an email to [Paul.Czodrowski@merckgroup.com](mailto:Paul.Czodrowski@merckgroup.com) or [wolf-guido.bolick@external.merckgroup.com](mailto:wolf-guido.bolick@external.merckgroup.com)

# Front-End: Output

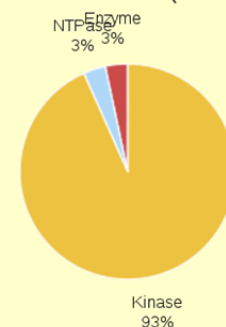
	#	Target_ID	Target_Name	Classification	#Ligands	e-Value	mean-TC
+	0	<a href="#">CHEMBL203</a>	Epidermal growth factor receptor erbB1	Kinase	2834	8.25e-218	0.448
+	1	<a href="#">CHEMBL1824</a>	Receptor protein-tyrosine kinase erbB-2	Kinase	1178	2.60e-184	0.464
+	2	<a href="#">CHEMBL3009</a>	Receptor protein-tyrosine kinase erbB-4	Kinase	131	6.18e-64	0.461
+	3	<a href="#">CHEMBL279</a>	Vascular endothelial growth factor receptor 2	Kinase	3953	1.79e-52	0.400
+	4	<a href="#">CHEMBL2634</a>	Tyrosine-protein kinase CSK	Kinase	38	8.19e-43	0.354
+	5	<a href="#">CHEMBL1868</a>	Vascular endothelial growth factor receptor 1	Kinase	858	1.98e-41	0.388
+	6	<a href="#">CHEMBL3650</a>	Fibroblast growth factor receptor 1	Kinase	870	3.07e-30	0.403
+	7	<a href="#">CHEMBL3935</a>	Serine/threonine-protein kinase Aurora-C	Kinase	23	1.03e-23	0.431
+	8	<a href="#">CHEMBL267</a>	Tyrosine-protein kinase SRC	Kinase	1908	6.25e-23	0.402
+	9	<a href="#">CHEMBL5251</a>	Tyrosine-protein kinase BTK	Kinase	216	1.04e-20	0.424
+	10	<a href="#">CHEMBL4722</a>	Serine/threonine-protein kinase Aurora-A	Kinase	1497	1.62e-20	0.400
+	11	<a href="#">CHEMBL3290</a>	Ephrin type-B receptor 2	Kinase	12	2.31e-17	0.410
+	12	<a href="#">CHEMBL1913</a>	Platelet-derived growth factor receptor beta	Kinase	677	2.79e-17	0.440
+	13	<a href="#">CHEMBL4142</a>	Fibroblast growth factor receptor 2	Kinase	49	4.60e-17	0.407
+	14	<a href="#">CHEMBL1955</a>	Vascular endothelial growth factor receptor 3	Kinase	508	1.06e-16	0.392
+	15	<a href="#">CHEMBL2073</a>	Tyrosine-protein kinase YES	Kinase	44	3.59e-15	0.430
+	16	<a href="#">CHEMBL5699</a>	Serine/threonine-protein kinase SIK2	Kinase	20	1.21e-13	0.401
+	17	<a href="#">CHEMBL3975</a>	Fructose-1,6-bisphosphatase	Enzyme	288	1.55e-11	0.374
+	18	<a href="#">CHEMBL2185</a>	Serine/threonine-protein kinase Aurora-B	Kinase	1329	1.57e-10	0.390
+	19	<a href="#">CHEMBL2250</a>	Tyrosine-protein kinase BLK	Kinase	230	6.46e-10	0.392
+	20	<a href="#">CHEMBL2007</a>	Platelet-derived growth factor receptor alpha	Kinase	302	9.28e-09	0.389
+	21	<a href="#">CHEMBL2041</a>	Tyrosine-protein kinase receptor RET	Kinase	452	4.06e-08	0.388
+	22	<a href="#">CHEMBL3905</a>	Tyrosine-protein kinase Lyn	Kinase	369	3.79e-07	0.415

## Query:

Smiles: CN(C)C/C=C\C(=O)Nc1cc2c(Nc3ccc(F)c(Cl)c3)n cnc2cc1O[C@H]4CCOC4.OC(=O)\C=C/C(=O)O



## Class-Distribution (first 30)

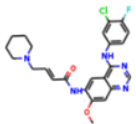
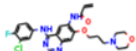
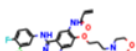
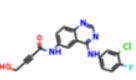




# Front-End: Output

#	Target_ID	Target_Name	Classification	#Ligands	e-Value	mean-TC	threshold	activ
+	0	<a href="#">CHEMBL203</a>	Epidermal growth factor receptor erbB1	Kinase	2834	8.25e-218	0.448	
+	1	<a href="#">CHEMBL1824</a>	Receptor protein-tyrosine kinase erbB-2	Kinase	1178	2.60e-184	0.464	
+	2	<a href="#">CHEMBL3009</a>	Receptor protein-tyrosine kinase erbB-4	Kinase	131	6.18e-64	0.461	

#	molregno	Molecule ID	Smiles	Molecule	Standard Value	Tanimoto Distance to Query ▲
87	1377989	<a href="#">CHEMBL2105719</a>	<chem>O.COC1cc2ncnc(Nc3ccc(F)c(Cl)c3)c2cc1NC(=O)C=C\CN4CCCCC4</chem> <a href="#">u</a> Kinase: 87 % NTPase: 3 % Enzyme: 3 % Hydrolase: 3 %		74	0.890681003584
69	41764	<a href="#">CHEMBL31965</a>	<chem>Fc1ccc(Nc2ncnc3cc(OCCCN4CCOCC4)c(NC(=O)C=C)cc23)cc1Cl</chem> <a href="#">u</a> Kinase: 90 % Enzyme: 3 % Phosphodiesterase: 3 % 7TMI: 3 %		12	0.890070921986
17	547496	<a href="#">CHEMBL545315</a>	<chem>Cl.Cl.Fc1ccc(Nc2ncnc3cc(OCCCN4CCOCC4)c(NC(=O)C=C)cc23)cc1Cl</chem> <a href="#">u</a> <a href="#">Get Hit Profile</a>		31	0.890070921986
103	338755	<a href="#">CHEMBL202411</a>	<chem>OCC#CC(=O)Nc1ccc2ncnc(Nc3ccc(F)c(Cl)c3)c2c1</chem> <a href="#">u</a> <a href="#">Get Hit Profile</a>		7.3	0.672855879752

## Next steps

- Still looking for publically available validation data not contained in ChEMBL
- Validation on internal data
- PCA plot (how well is the query compound represented by the target data)