

Proteochemometrics (PCM) de novo generation (DrugEx)



Willem Jaspers

LACDR

Sequence alignment

Alignment where match is cheap, mismatch is cheap and gap is costly

This is a hard example. -----

|| |||| | |

That is another easy example.

:. ** . * *

Algorithm wants to have maximal points! Aligning correct get's points, misaligning and gaps cost points..!

An * (asterisk) indicates positions which have a single, fully conserved residue.

A : (colon) indicates conservation between groups of strongly similar properties

A . (period) indicates conservation between groups of weakly similar properties

BLOcks SUBstitution (BLOSUM) Matrix

Less frequently occurring or high impact residues aligned correctly increase score

Ala	4																			
Arg	-1	5																		
Asn	-2	0	6																	
Asp	-2	-2	1	6																
Cys	0	-3	-3	-3	9															
Gln	-1	1	0	0	-3	5														
Glu	-1	0	0	2	-4	2	5													
Gly	0	-2	0	-1	-3	-2	-2	6												
His	-2	0	1	-1	-3	0	0	-2	8											
Ile	-1	-3	-3	-3	-1	-3	-3	-4	-3	4										
Leu	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4									
Lys	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5								
Met	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5							
Phe	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6						
Pro	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7					
Ser	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4				
Thr	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5			
Trp	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11		
Tyr	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	
Val	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4
	Ala	Arg	Asn	Asp	Cys	Gln	Glu	Gly	His	Ile	Leu	Lys	Met	Phe	Pro	Ser	Thr	Trp	Tyr	Val

Aligning two sequences

- CALHP & CLLHL
- Gap = -10

	Cys	Ala	Leu	His	Pro
Cys	9	0	-1	-3	-3
Leu	-1	-1	4	-3	-3
Leu	-1	-1	4	-3	-3
His	-3	-2	-3	8	-2
Leu	-1	-1	4	-3	-3

Option 1 : Score -7

C-ALHP

| |

CLL-HL

Option 2 : Score 0

CALH-P

| | |

CLLHL-

Option 3 : Score 17

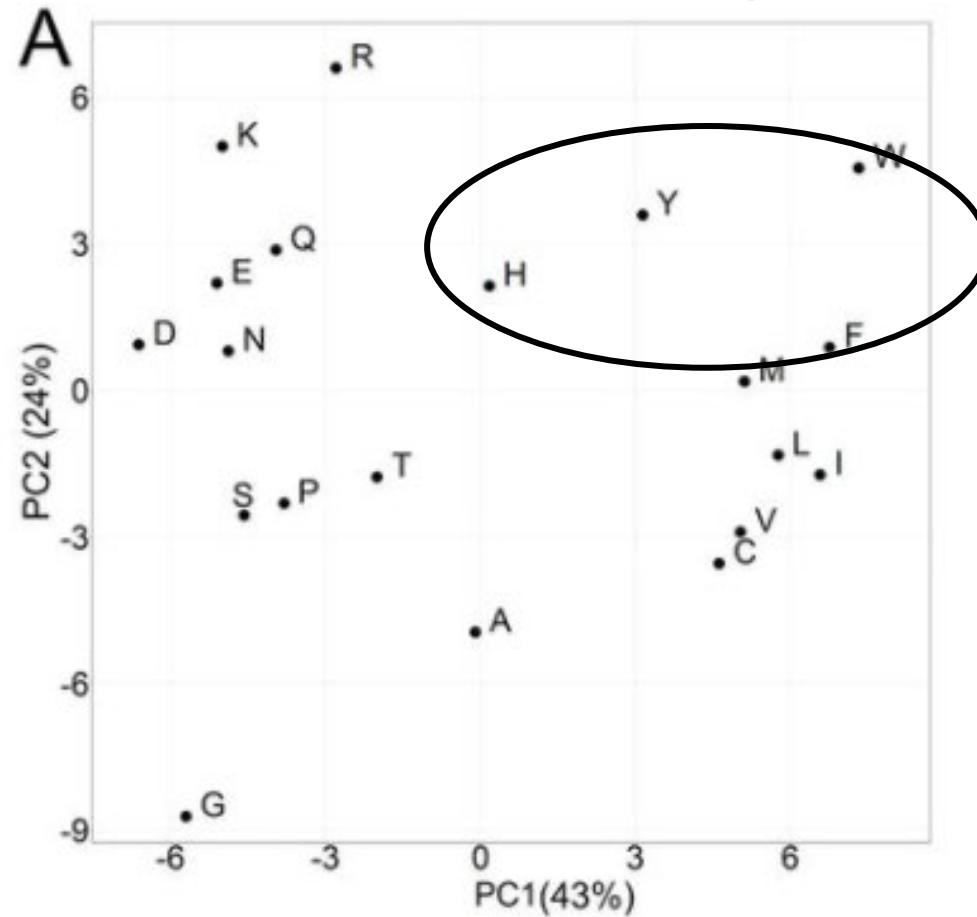
CALHP

| | |

CLLHL

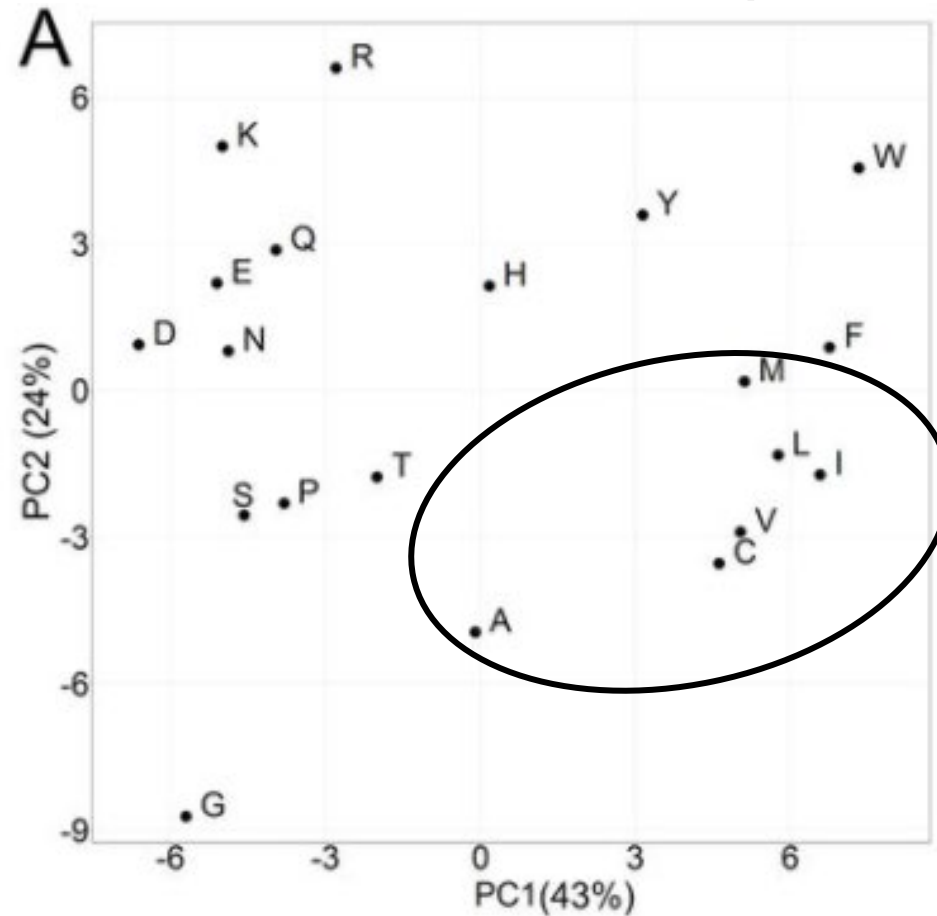
Peptide / binding site descriptors

- Similar to the fingerprints used for small molecules, we can convert amino acids to descriptors.



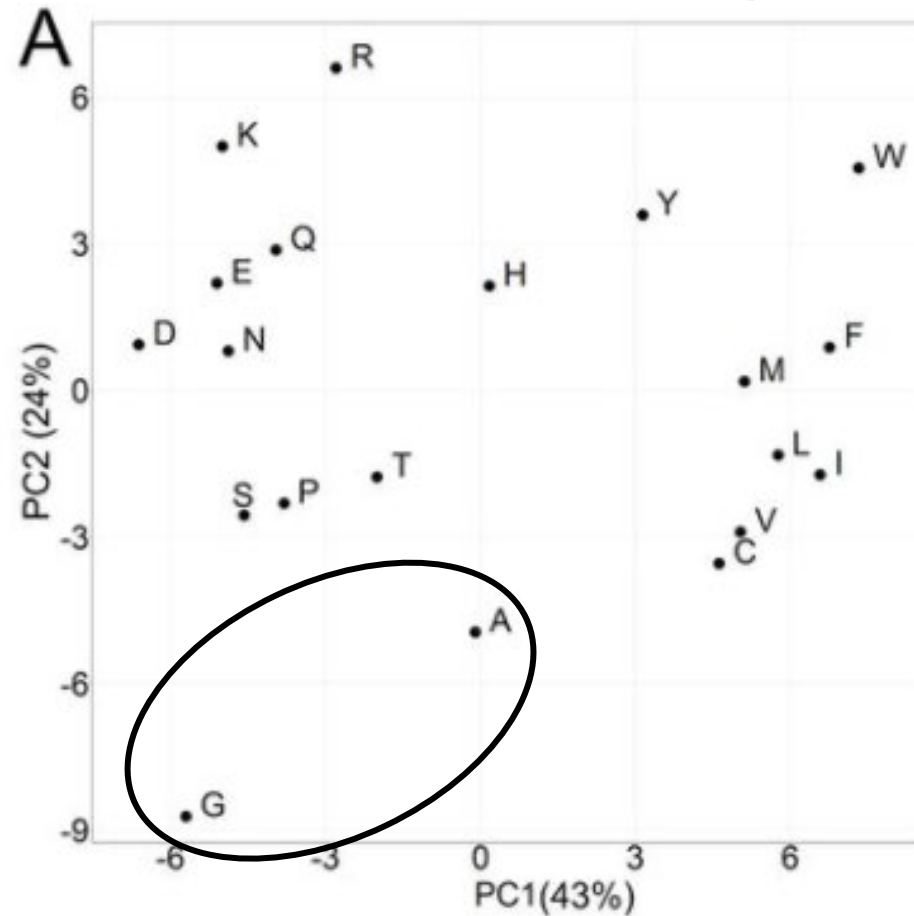
Peptide / binding site descriptors

- Similar to the fingerprints used for small molecules, we can convert amino acids to descriptors.



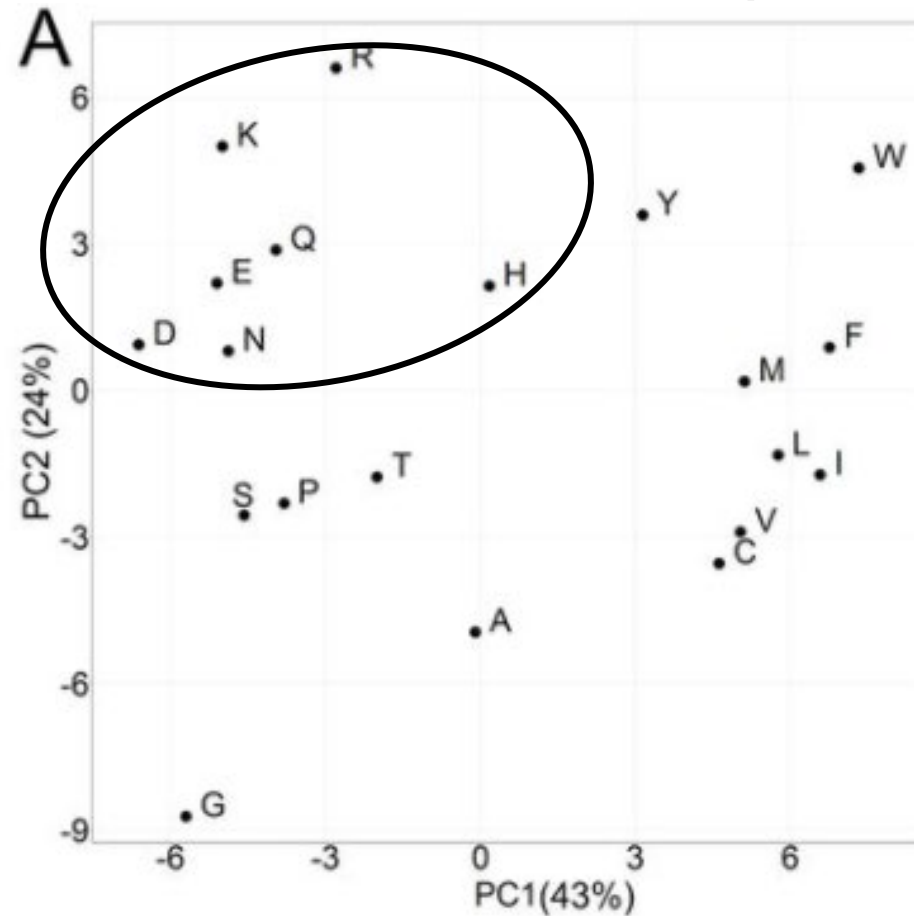
Peptide / binding site descriptors

- Similar to the fingerprints used for small molecules, we can convert amino acids to descriptors.



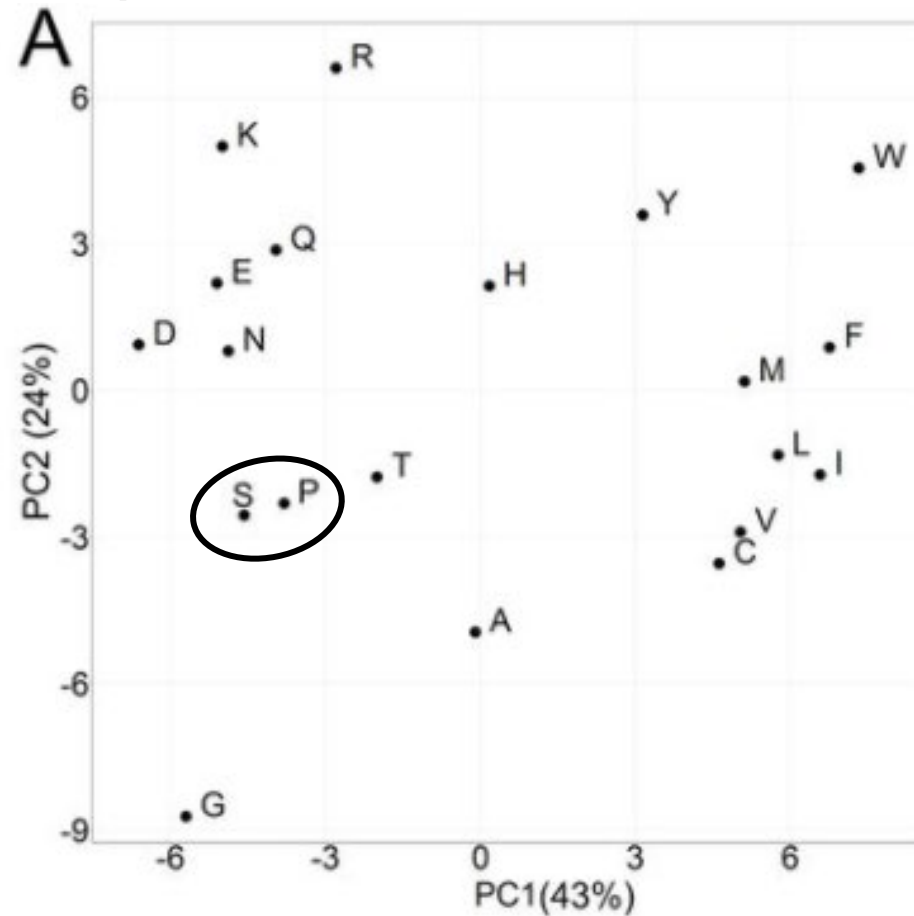
Peptide / binding site descriptors

- Similar to the fingerprints used for small molecules, we can convert amino acids to descriptors.



Peptide / binding site descriptors

- Similar to the fingerprints used for small molecules, we can convert amino acids to descriptors.



How to quantify sequence similarity

- Take the examples from just now and add 4 more sequences
 - CALHP
 - CLLHL
 - CLWHL
 - CLLYP
 - GLLWT
 - GLLYT

How to quantify sequence similarity

- Take the examples from just now and add 4 more sequences

CLUSTAL O(1.2.4) multiple sequence alignment

seq_E	GLLWT-
seq_F	GLLYT-
seq_D	-CLLYP
seq_A	-CALHP
seq_B	-CLLHL
seq_C	-CLWHL

Converting these to descriptors

Receptor	Amino_acid_sequence	Amino_acid_number	Z1	Z2	Z3
A	-	1	0	0	0
	C	2	0.84	-1.67	3.71
	A	3	0.24	-2.32	0.6
	L	4	-4.28	-1.30	-1.49
	H	5	2.47	1.95	0.26
	P	6	-1.66	0.27	1.84
B	-	1	0	0	0
	C	2	0.84	-1.67	3.71
	L	3	-4.28	-1.30	-1.49
	L	4	-4.28	-1.30	-1.49
	H	5	2.47	1.95	0.26
	L	6	-4.28	-1.30	-1.49

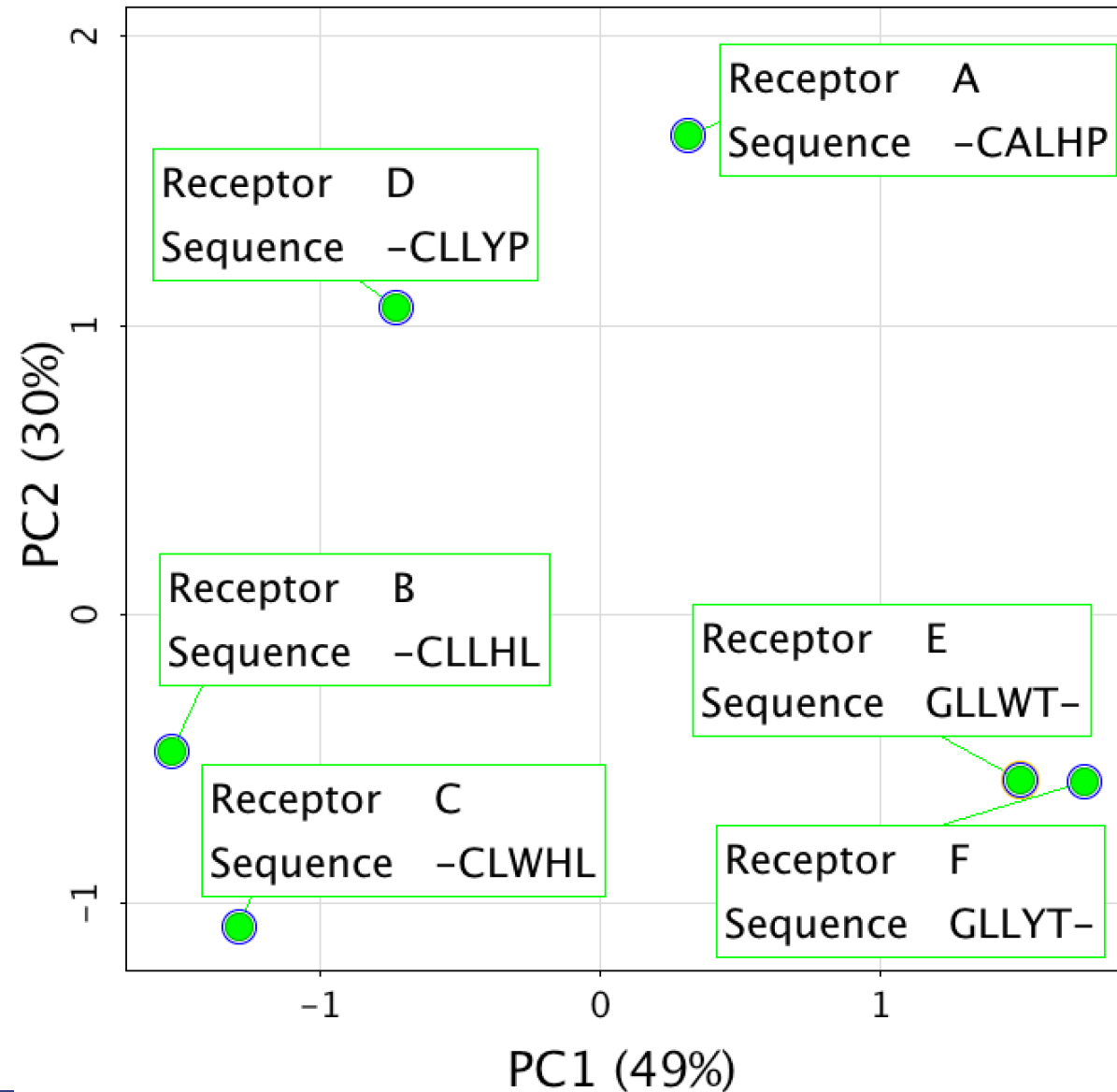
Converting these to descriptors

Receptor	Amino acid sequence	Amino acid number	Z1	Z2	Z3
A	-	1	0	0	0
	C	2	0.84	-1.67	3.71
	A	3	0.24	-2.32	0.6
	L	4	-4.28	-1.30	-1.49
	H	5	2.47	1.95	0.26
	P	6	-1.66	0.27	1.84
B	-	1	0	0	0
	C	2	0.84	-1.67	3.71
	L	3	-4.28	-1.30	-1.49
	L	4	-4.28	-1.30	-1.49
	H	5	2.47	1.95	0.26
	L	6	-4.28	-1.30	-1.49

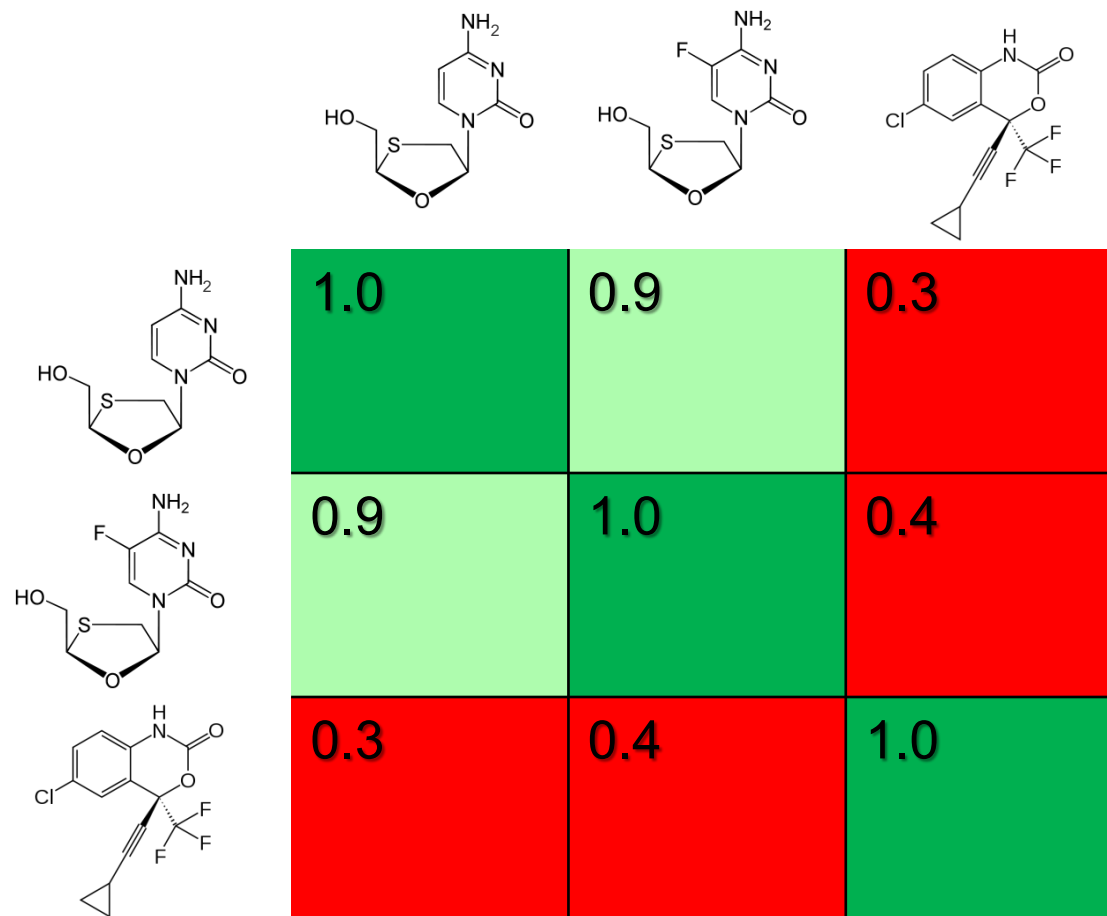
Converting these to descriptors

Receptor	Amino_acid_sequence	Amino_acid_number	Z1	Z2	Z3
A	-	1	0	0	0
	C	2	0.84	-1.67	3.71
	A	3	0.24	-2.32	0.6
	L	4	-4.28	-1.30	-1.49
	H	5	2.47	1.95	0.26
	P	6	-1.66	0.27	1.84
B	-	1	0	0	0
	C	2	0.84	-1.67	3.71
	I	3	-4.28	-1.30	-1.49
	L	4	-4.28	-1.30	-1.49
	H	5	2.47	1.95	0.26
	L	6	-4.28	-1.30	-1.49

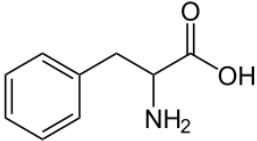
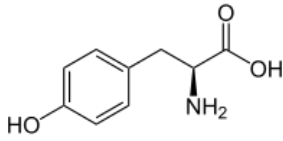
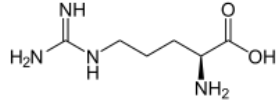
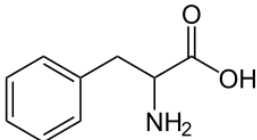
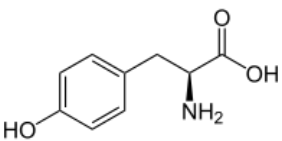
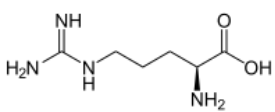
Quantify distance between sequences



Molecular Similarity



Sequence Similarity

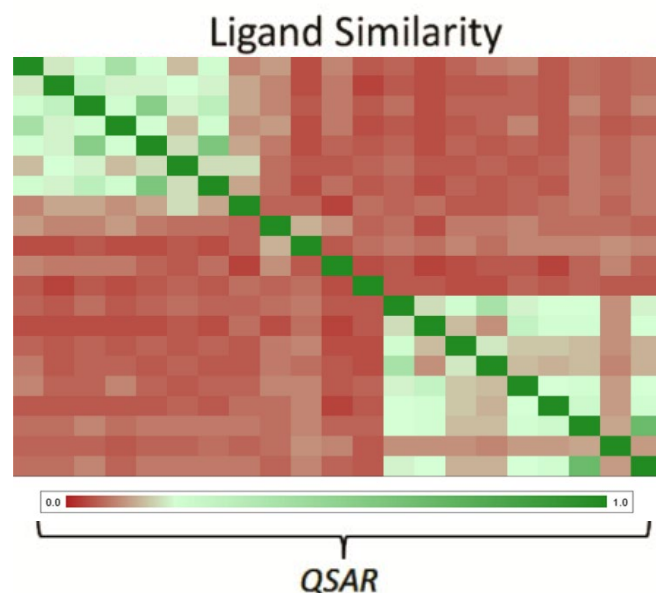
			
Phenylalanine	 1.0	0.9	0.3
Tyrosine	0.9	 1.0	0.4
Arginine	0.3	0.4	 1.0

Sequence Similarity

	FYI	IYF	WTF
FYI	1.0	0.9	0.3
IYF	0.9	1.0	0.4
WTF	0.3	0.4	1.0

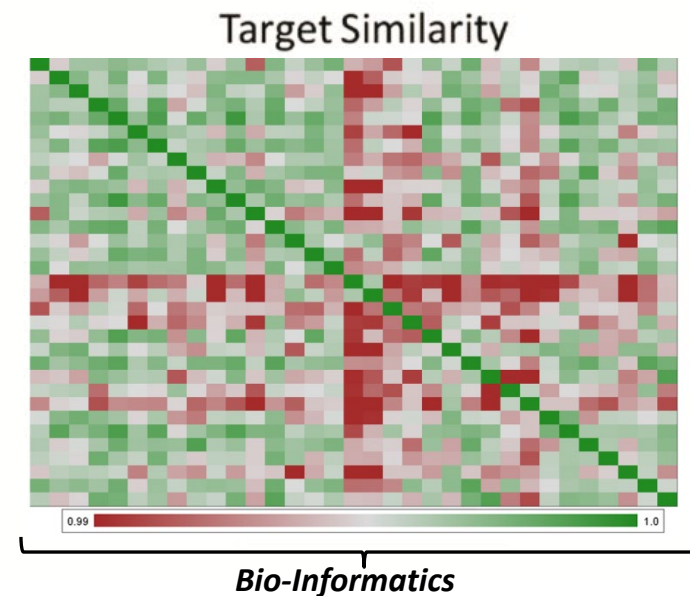
The how... *what is PCM* ?

- Proteochemometric modeling combines both a ligand descriptor and target descriptor



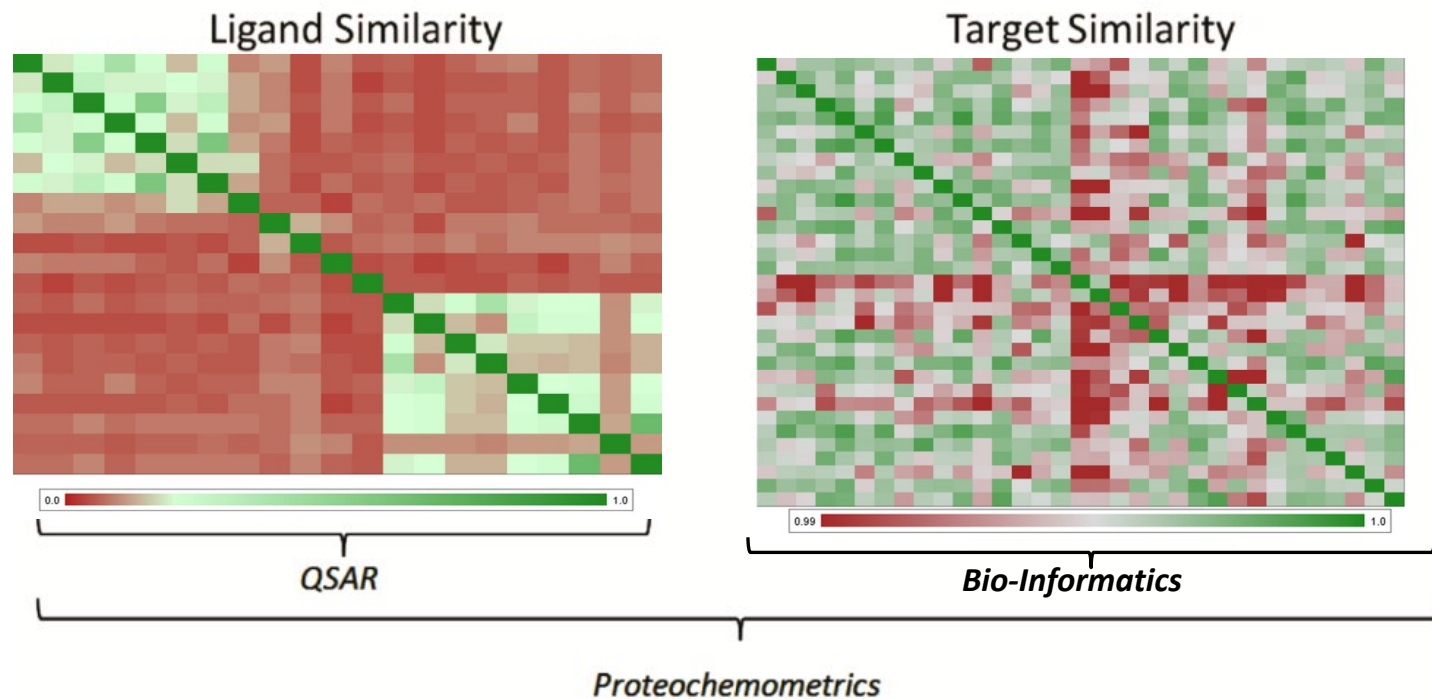
What is PCM ?

- Proteochemometric modeling combines both a ligand descriptor and target descriptor



What is PCM ?

- Proteochemometric modeling combines both a ligand descriptor and target descriptor

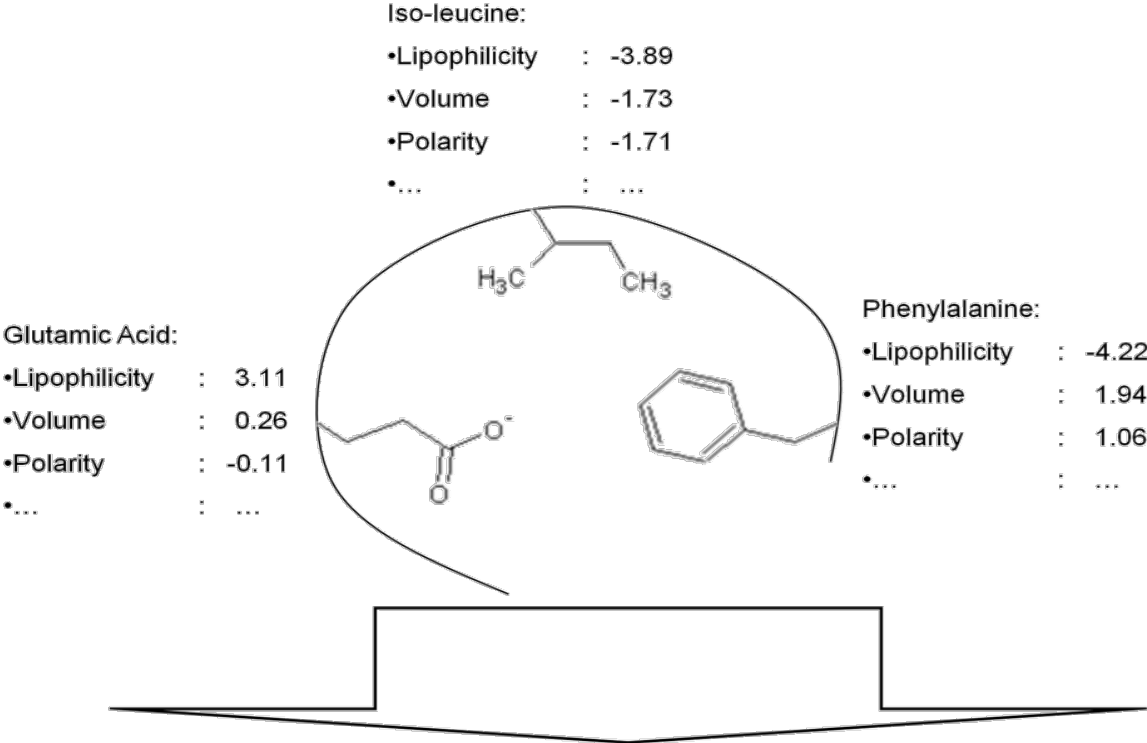


Describe protein properties

Describe protein properties

- Simple way to derive protein descriptors
 1. Align the relevant residues
 2. Convert to physicochemical properties

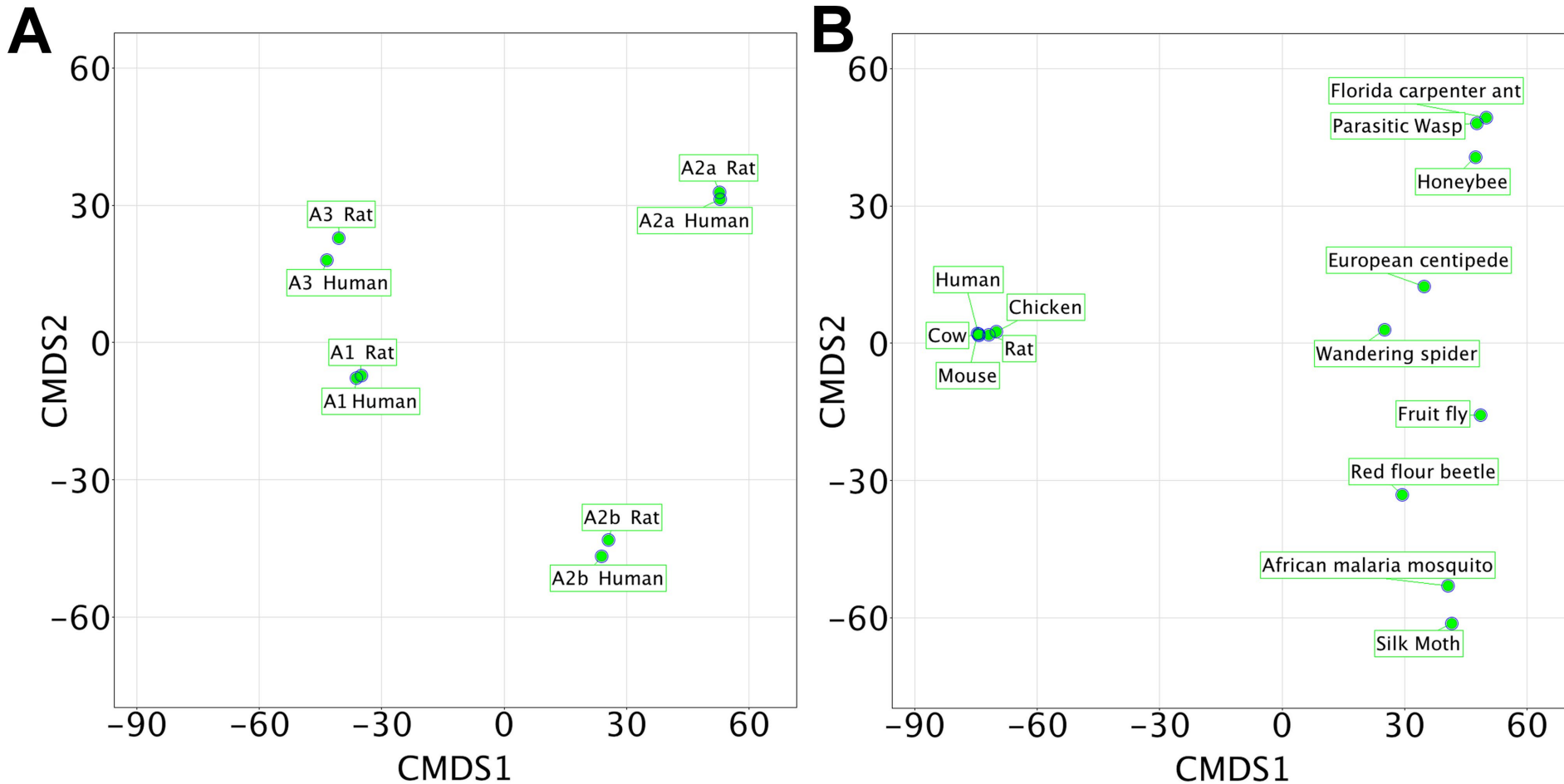
Describe protein properties



Sequence	Descriptor		
	Z1	Z2	Z3
E	3.11	0.26	-0.11
I	-3.89	-1.73	-1.71
F	-4.22	1.94	1.06
...

So what can we use this method for?

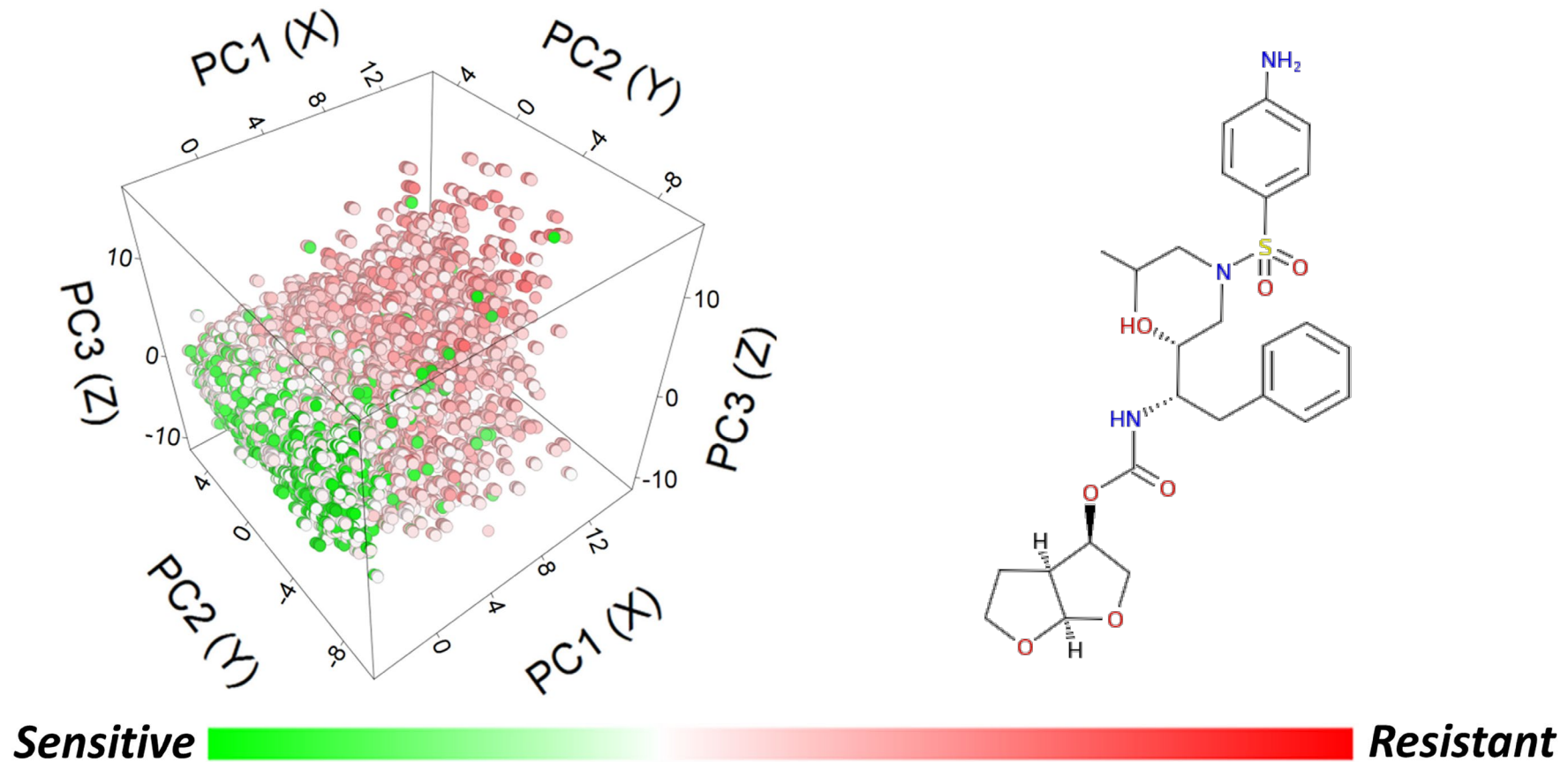
Target descriptors



a Similarity between the different adenosine receptors in human and rat.

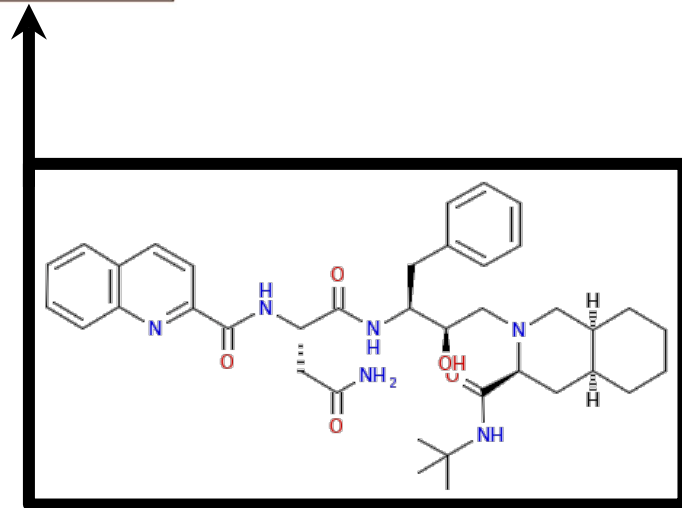
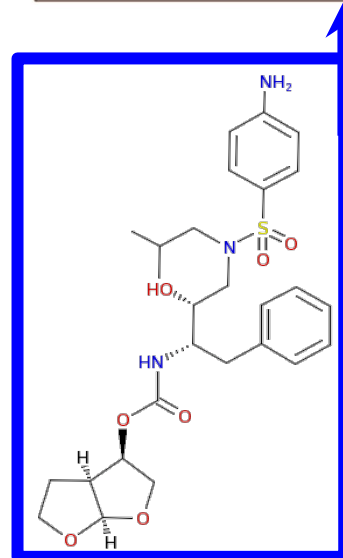
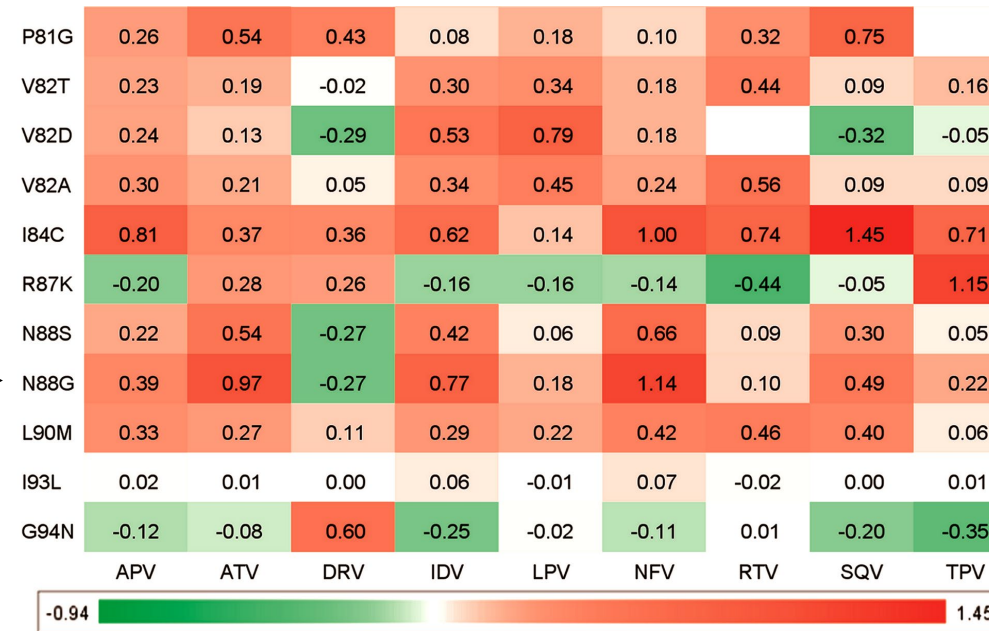
b Similarity between the GABA-A ligand-gated ion channels in mammals arthropods.

Prevent anti-HIV treatment failure



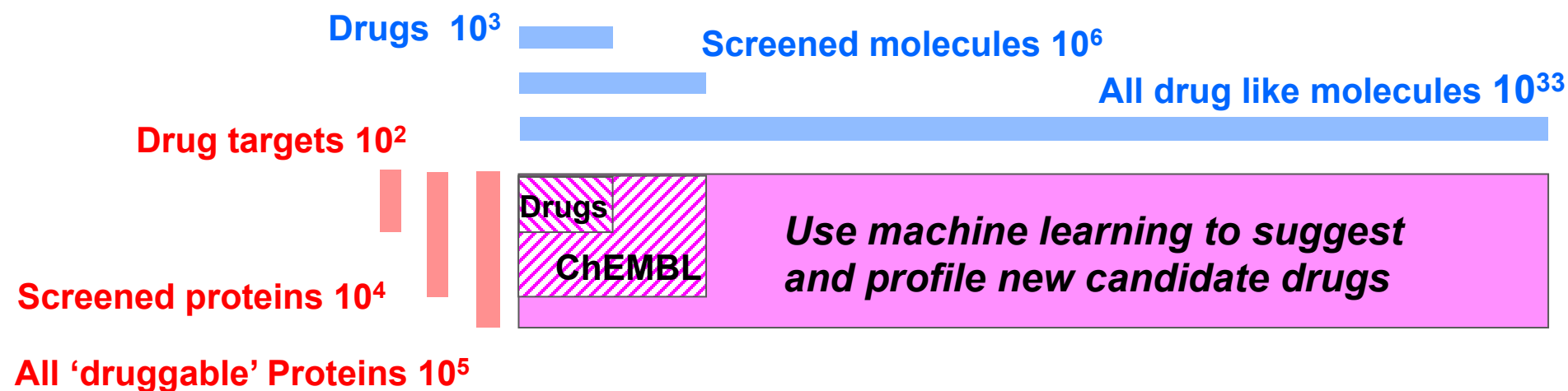
Prevent anti-HIV treatment failure

Patient
(Viral Genotype)

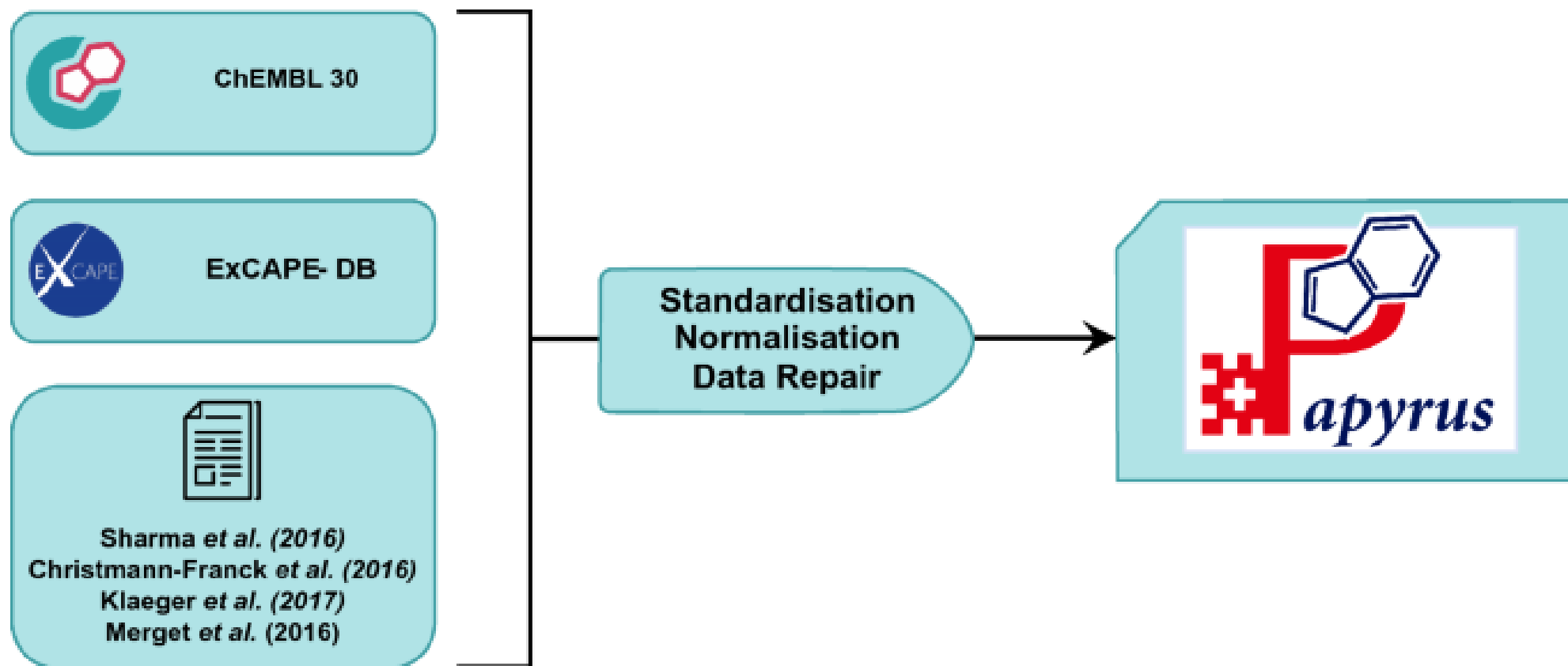


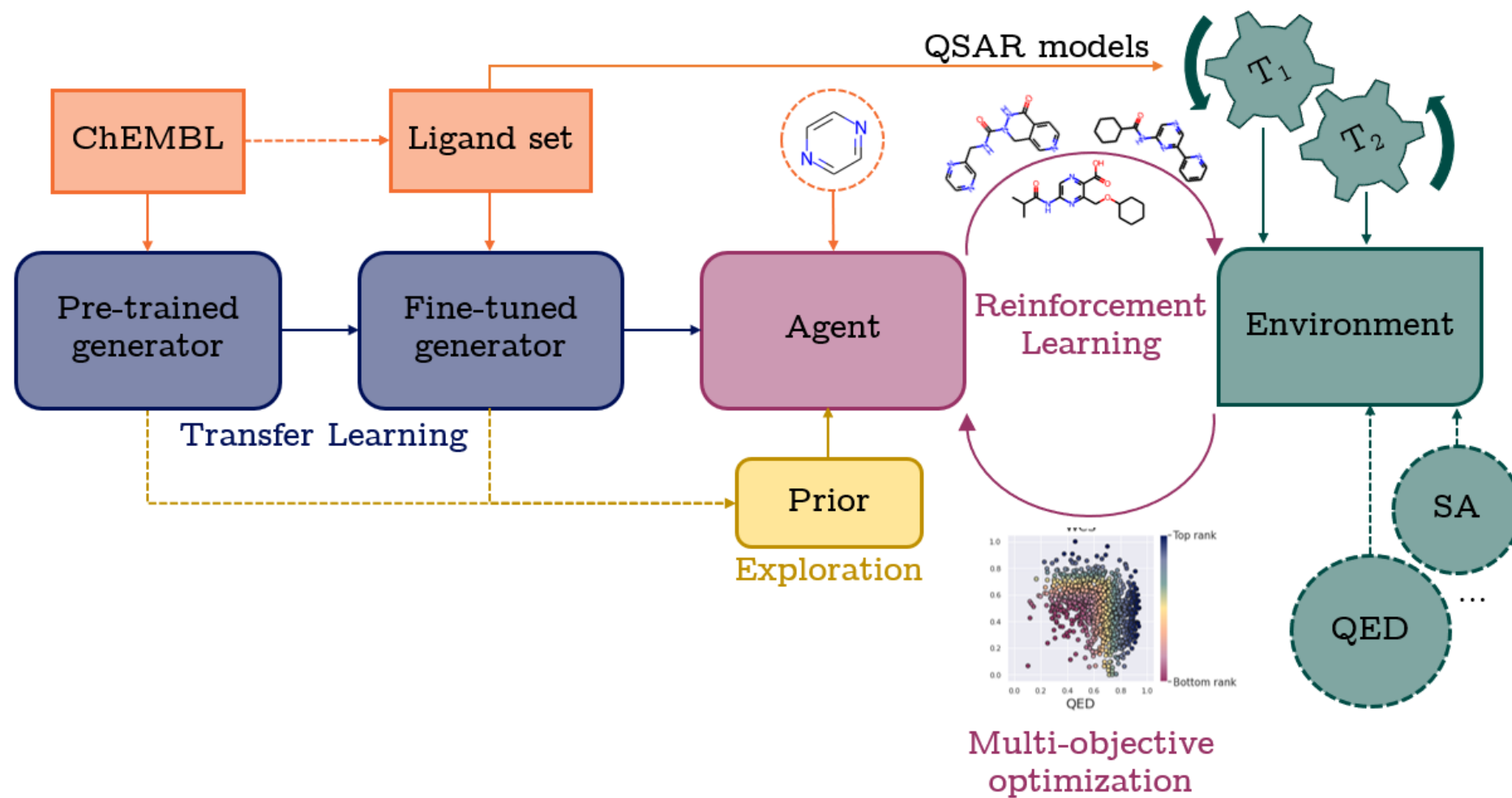
de novo generation

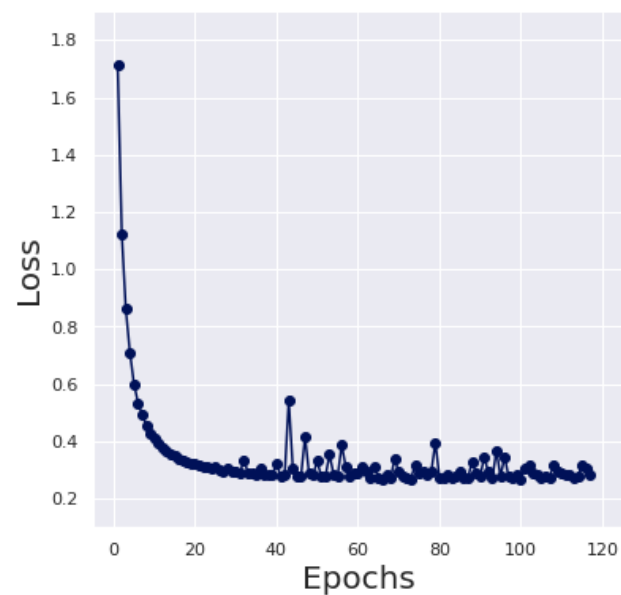
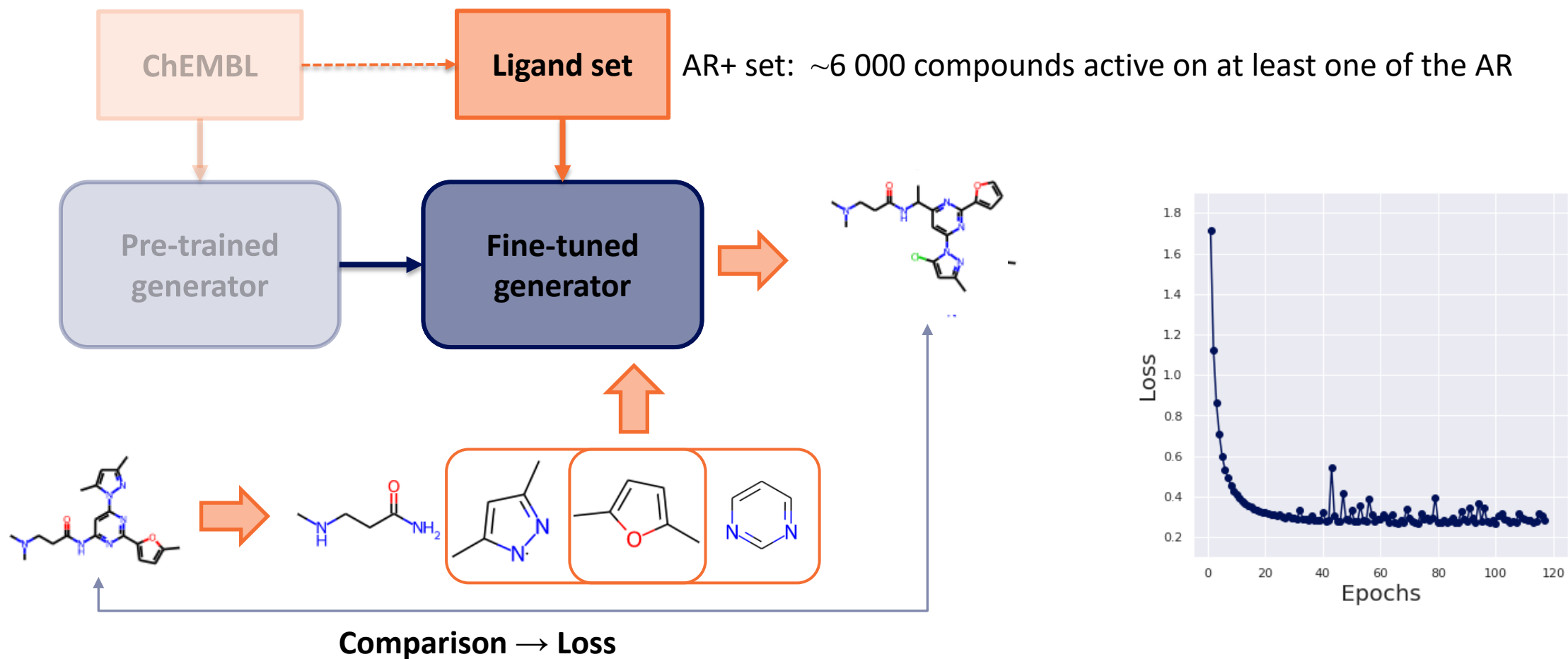
Learning a machine the grammar of molecules..



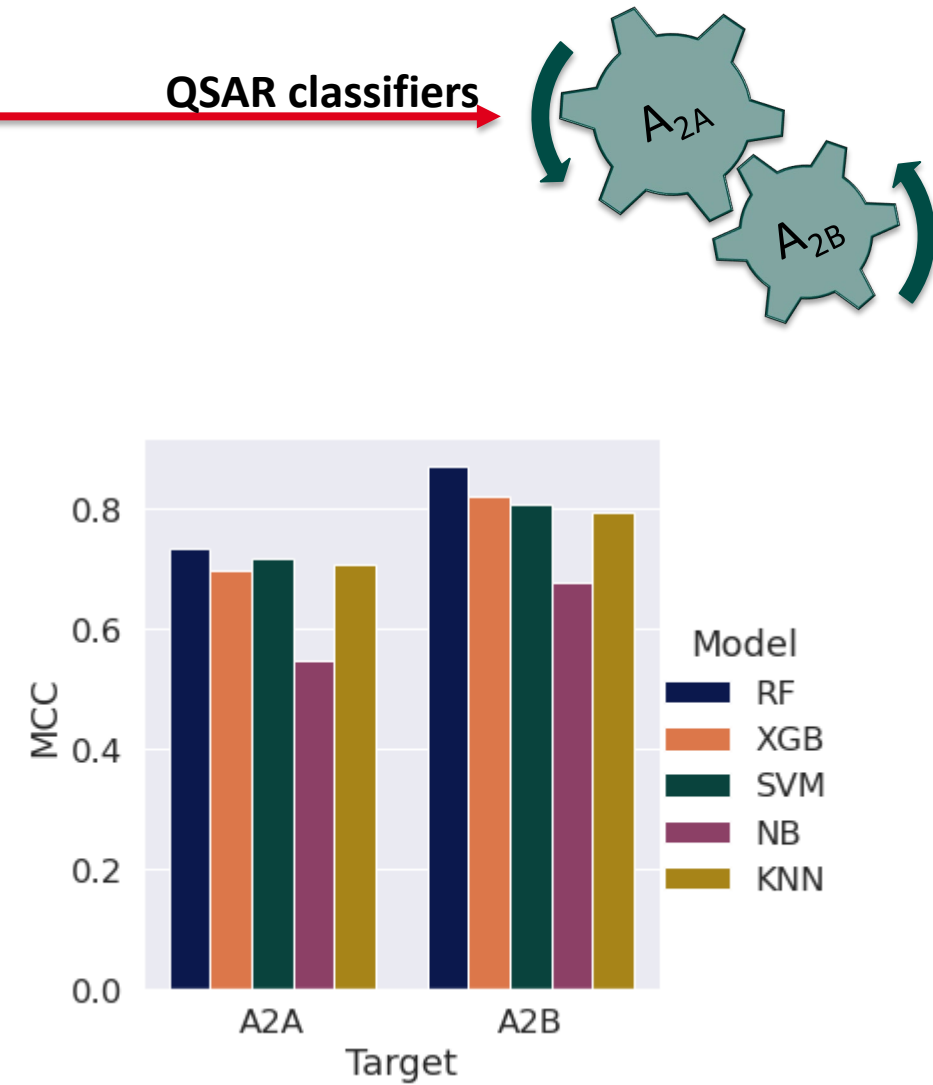
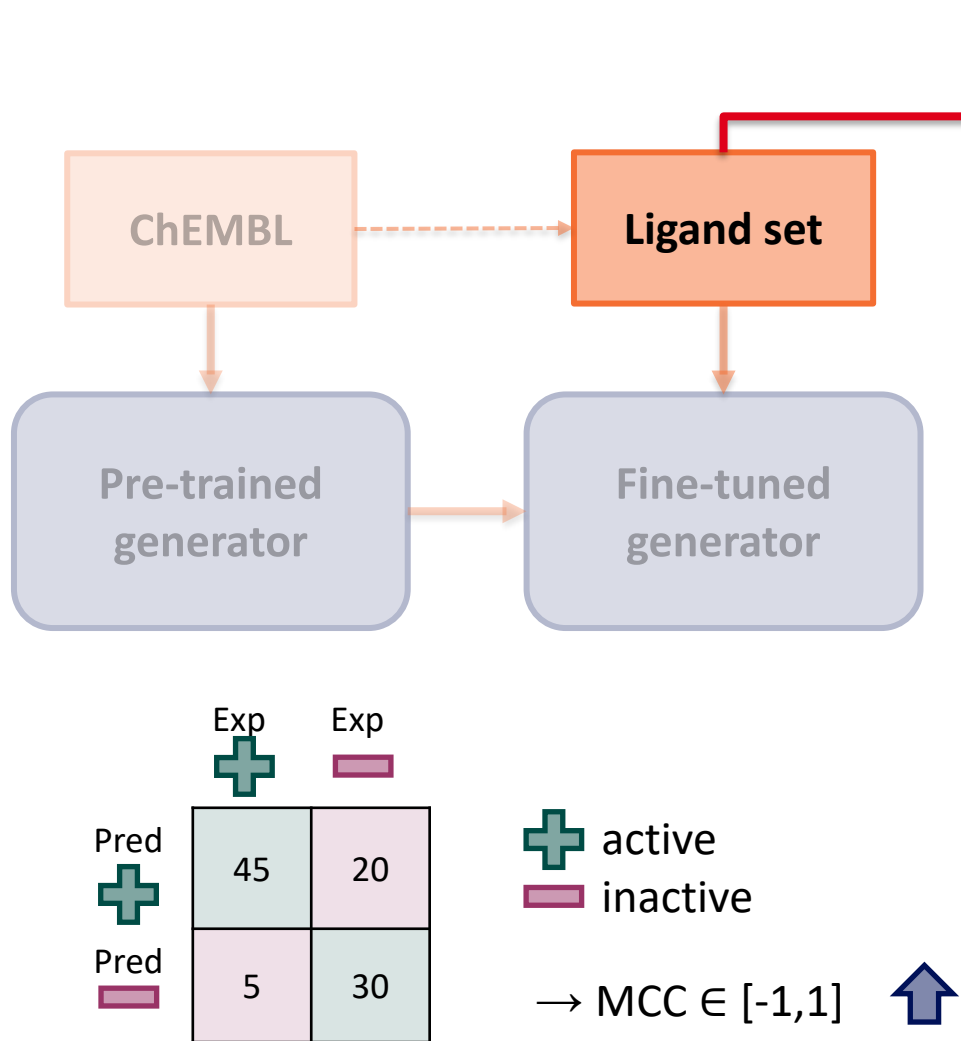
Input data: Papyrus

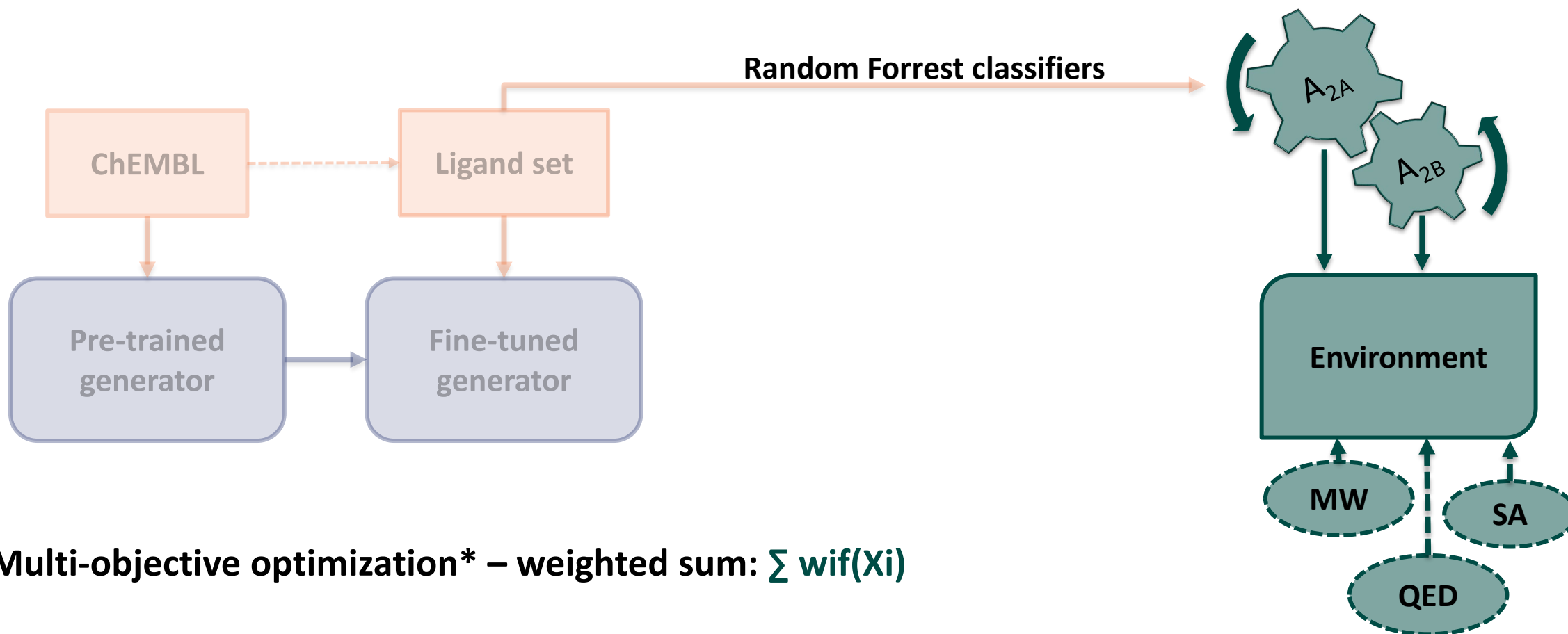






QSAR: quantitative structure-activity relationship





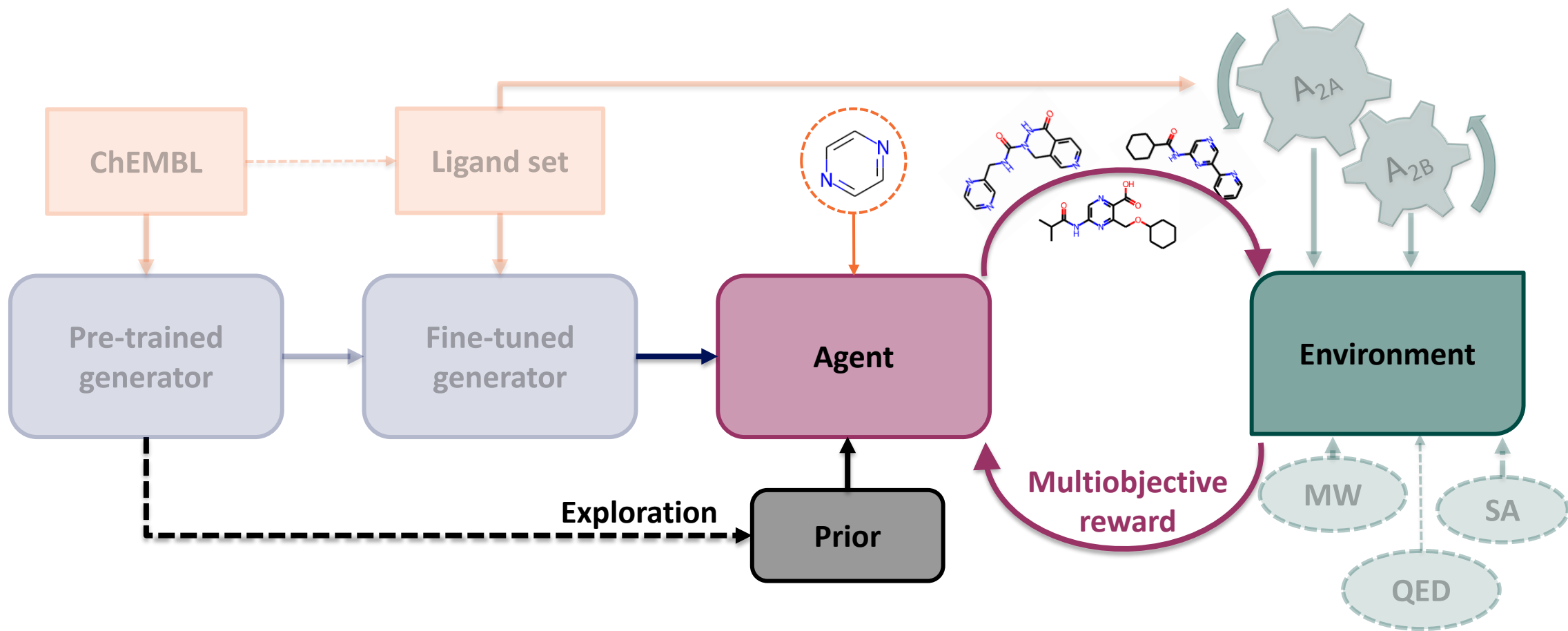
Multi-objective optimization* – weighted sum: $\sum w_i f(X_i)$

$f(X)$ – clips values between 0 (✖) and 1 (✓)

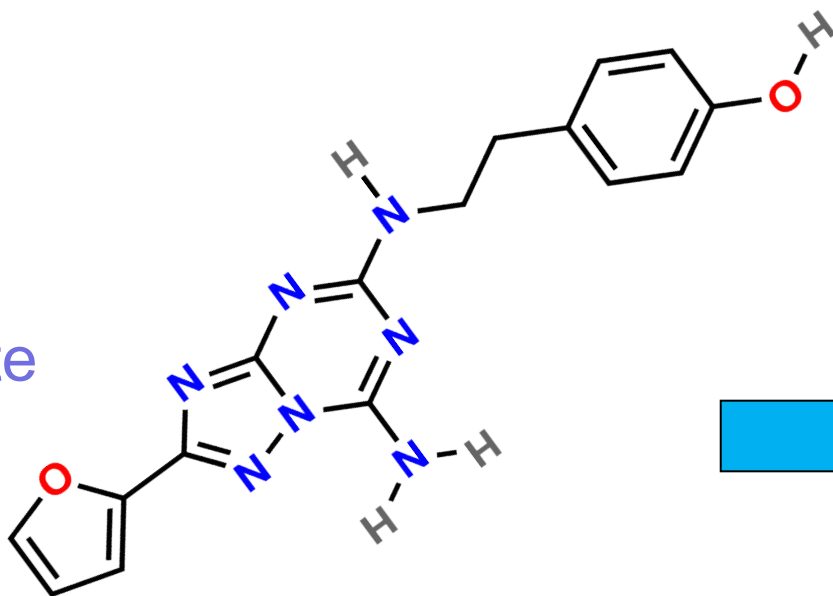
MW – Molecular weight : < 400 Da

SA – Synthetic accessibility

QED – Quantitative estimate of drug-likeness



1 - Translate



ZM241385

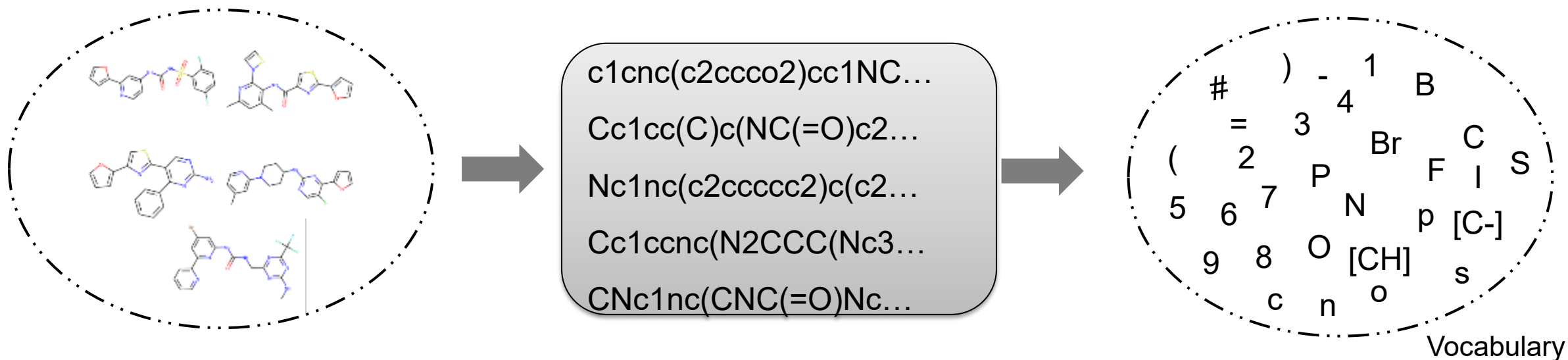


Nc1nc(NCCc2ccc(O)cc2)nc3nc(nn13)c4occc4

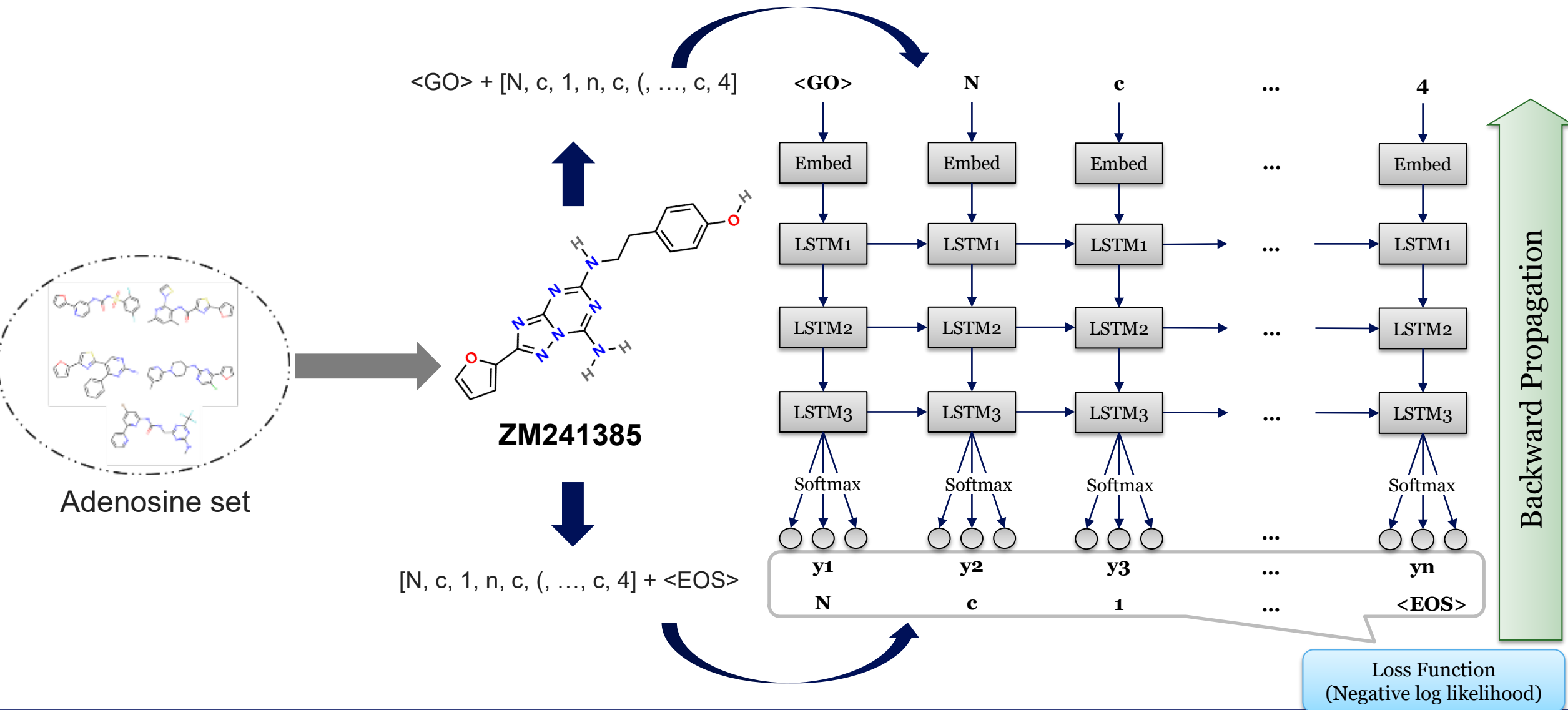
Dataset

Adenosine dataset

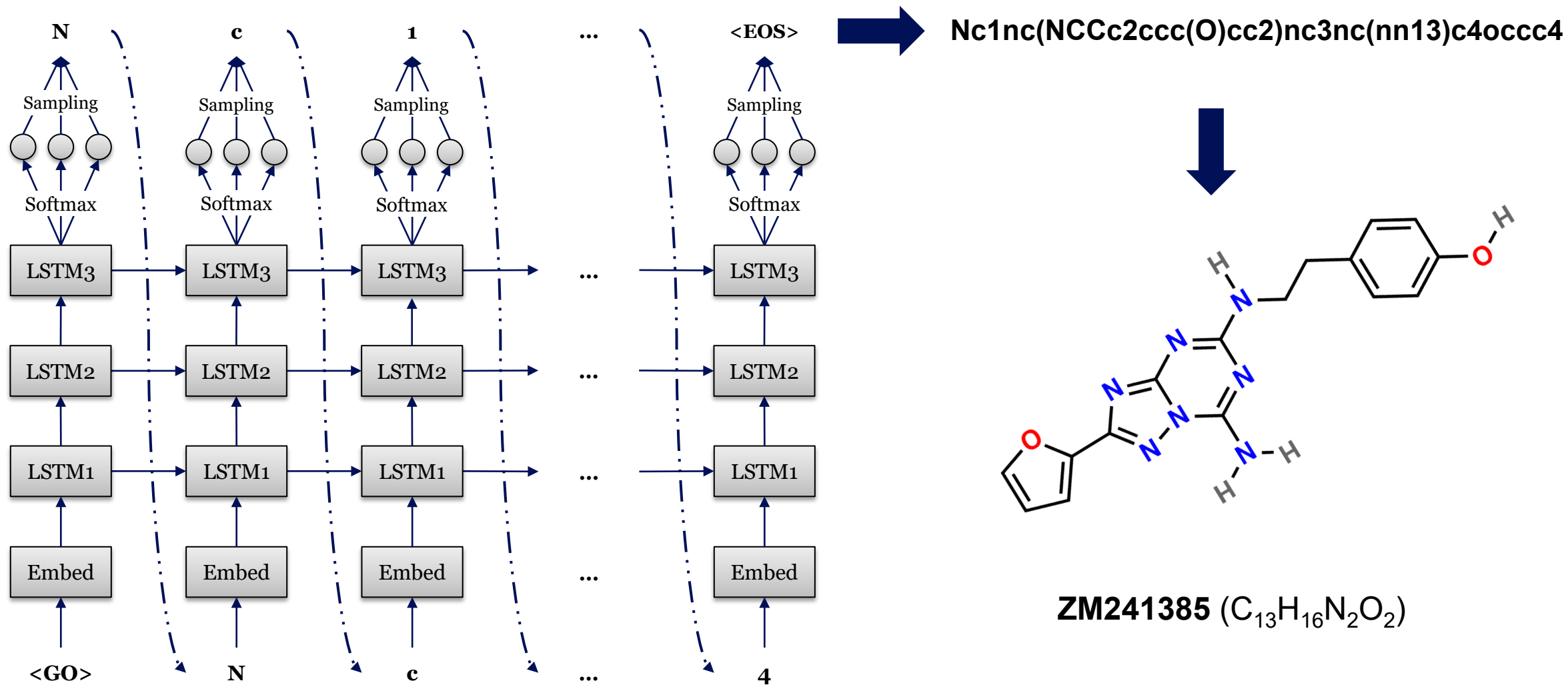
- All public compounds tested on the adenosine receptors ChEMBL (v 24).



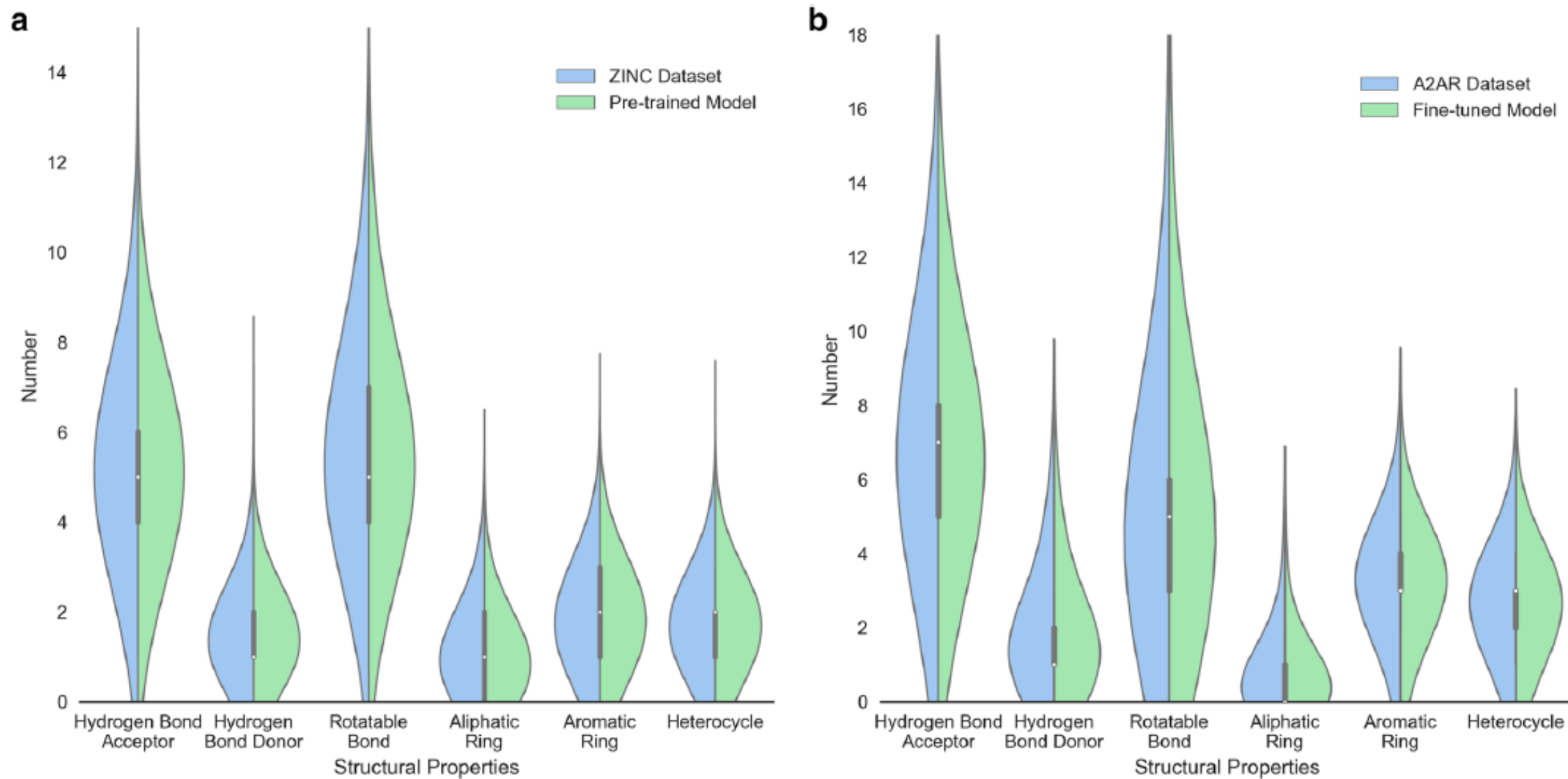
RNN Training



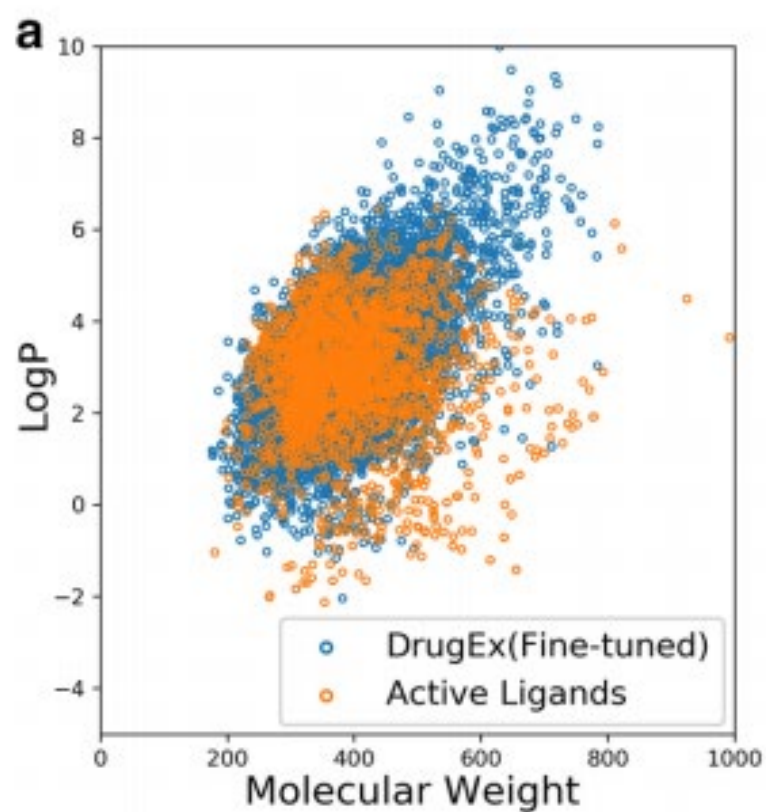
Molecule Generation



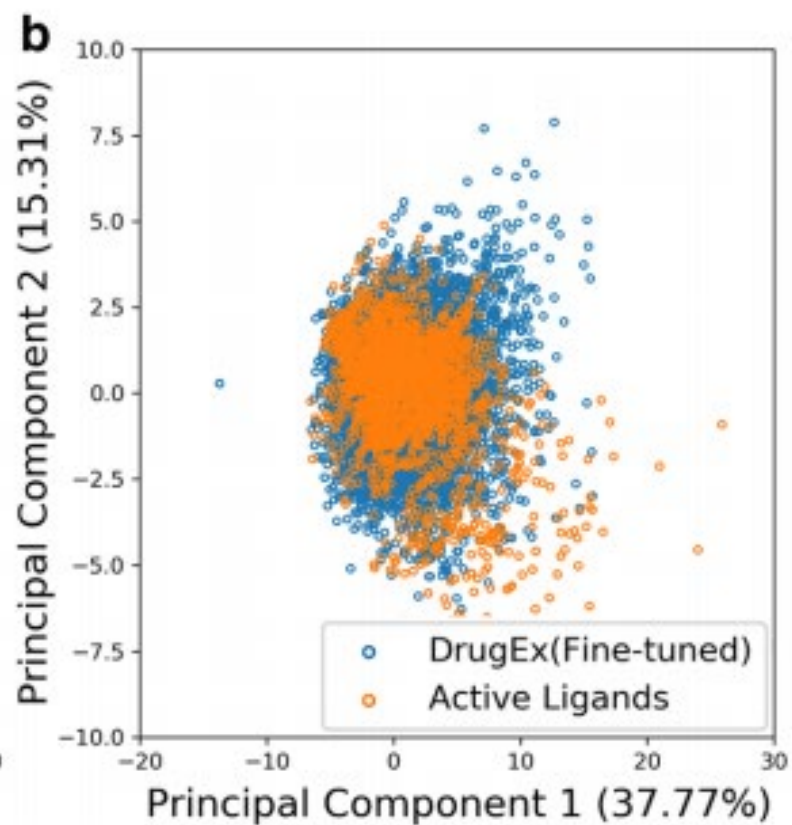
Generated molecules chemically similar



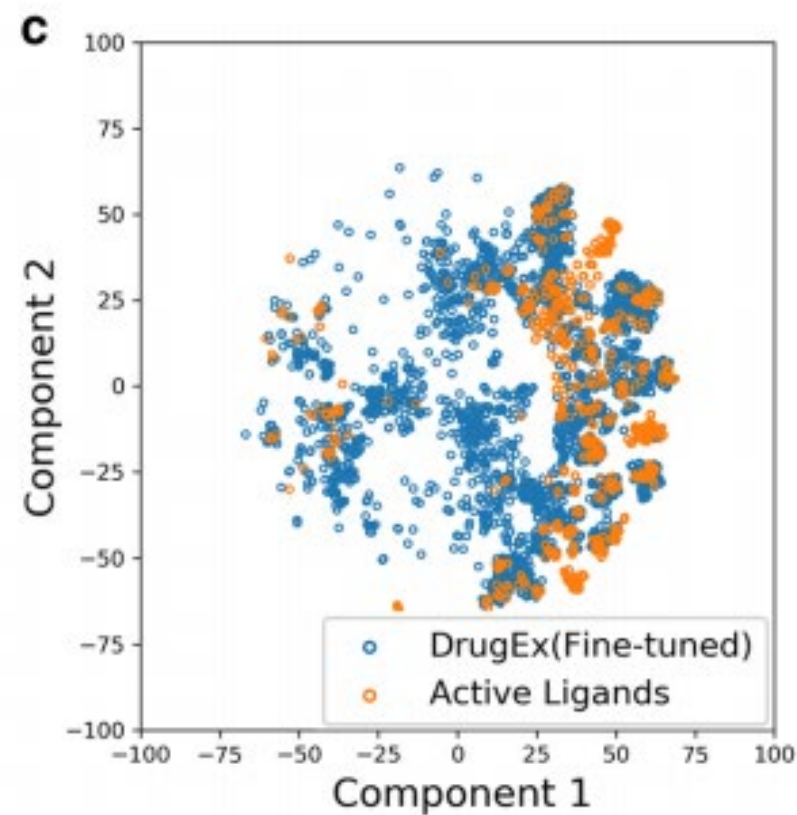
New A2A ligands



logP~MW



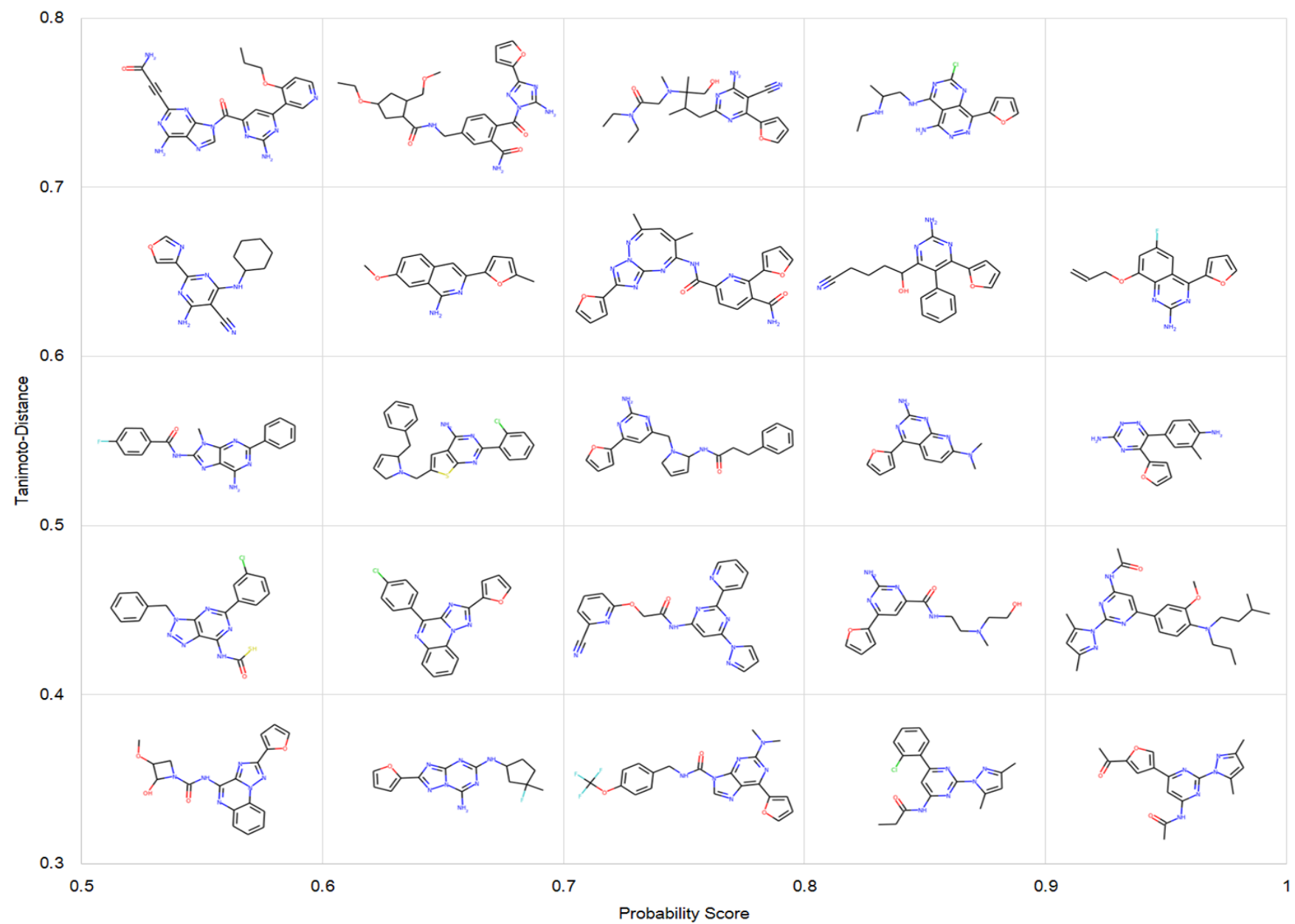
PCA (PhysChem)



t-SNE
(Fingerprints)

Also more complex chemical features are generated

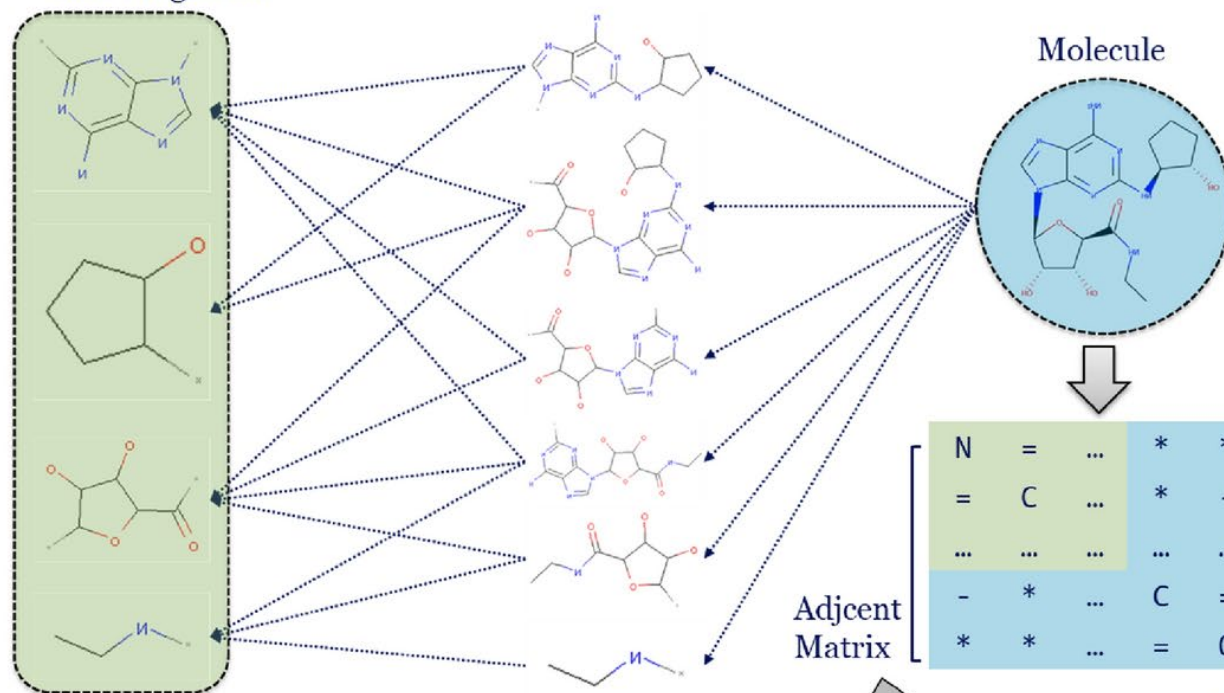
		Fused Ring	Furan Ring	Benzene Ring
DrugEx (Pre-trained)		9.12%	82.32%	61.48%
DrugEx (Fine-tuned)		60.69%	66.35%	65.62%
REINVENT		0.20%	95.26%	61.98%
ORGANIC		0.02%	99.96%	39.45%
Pre-trained		24.22%	4.51%	63.31%
Fine-tuned		76.33%	23.82%	72.85%
ZINC		26.66%	3.86%	63.97%
A2AR	Active	79.09%	40.29%	75.33%
	Inactive	76.73%	9.33%	70.88%



Graph based (this just in)

A

Leaf Fragments



Input (Scaffold)

NCC.C1CCCC1O
C(=O)C1OC(n2cnc3c(N)nc...
C1CCCC1O.C(=O)C1...
C(=O)C1OC(n2cnc3c(N)nc...
C(=O)C1OCC(O)C1O
c1nc(N)c2cnc2n1

Output (Molecule)

CCNC(=O)C1OC(n2cnc3c(N)nc...
CCNC(=O)C1OC(n2cnc3c(N)nc...
CCNC(=O)C1OC(n2cnc3c(N)nc...
CCNC(=O)C1OC(n2cnc3c(N)nc...
CCNC(=O)C1OC(n2cnc3c(N)nc...
CCNC(=O)C1OC(n2cnc3c(N)nc...

B

C

Graph Matrix

Atom Type	1	8	...	5	28	...	5	0	2	...	2
Bond Type	0	0	...	2	1	...	3	0	1	...	1
Connected atom index	0	0	...	5	0	...	5	0	2	...	13
Atom Index	0	0	...	10	11	...	43	0	10	...	20
Fragment Index	0	1	...	2	0	...	0	0	0	...	0

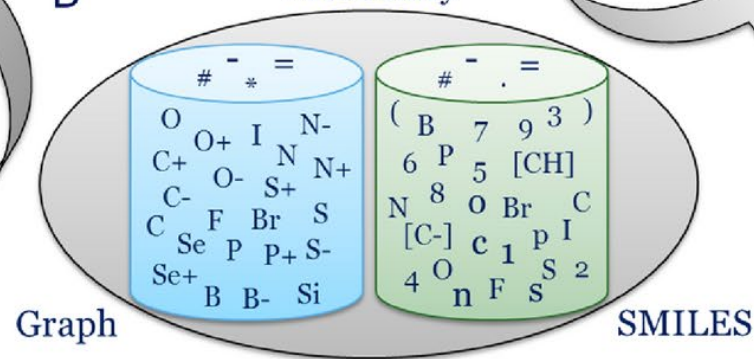
Scaffolds Growing Linking

Adjcent Matrix

Encoding

D

Vocabulary



Encoding

Input Matrix

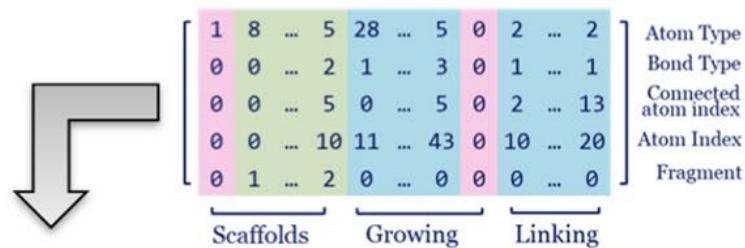
1	8	3	9	5	4	...	0
1	31	47	24	8	6	...	0
...
1	11	50	4	6	20	...	0

Output Matrix

1	9	28	5	44	29	...	0
1	9	28	5	44	29	...	0
...
1	9	28	5	44	29	...	0

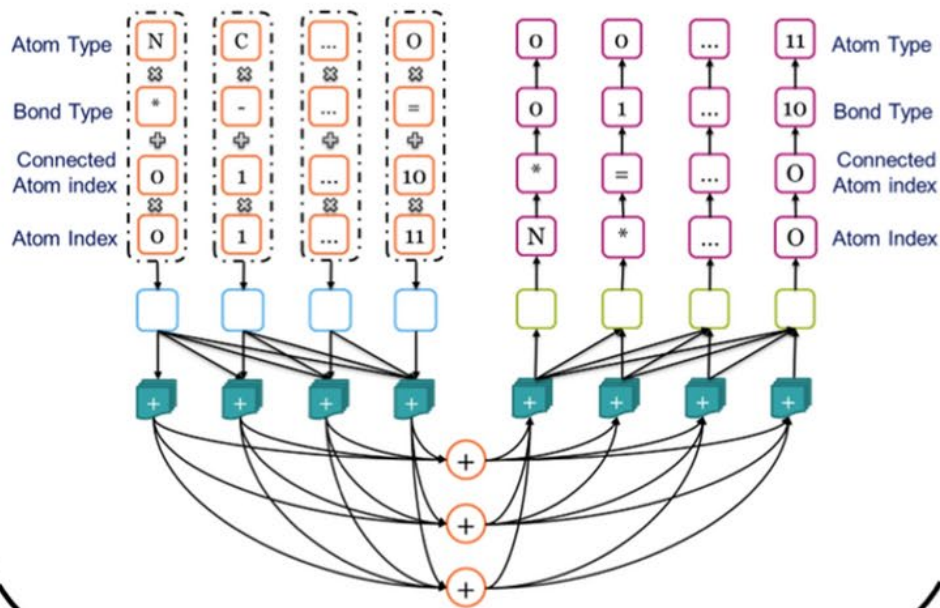
E

A. Graph Transformer



Word Input: $W = Type_{Atom} * 4 + Type_{Bond}$

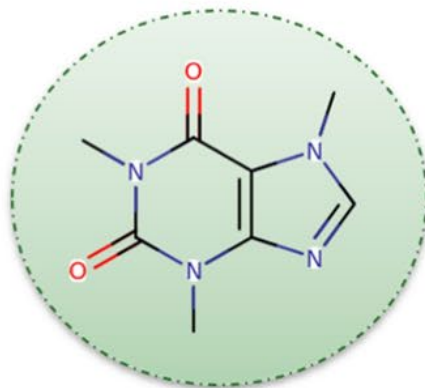
Position Input: $P = Index_{Atom} * L_{max} + Index_{Connected}$



Graph-based Generator

$$\begin{bmatrix} N & = & \dots & * & * \\ = & C & \dots & * & - \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ - & * & \dots & C & = \\ * & * & \dots & = & O \end{bmatrix}$$

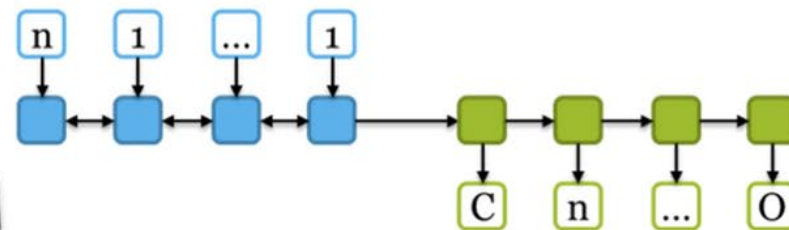
Adjacency Matrix



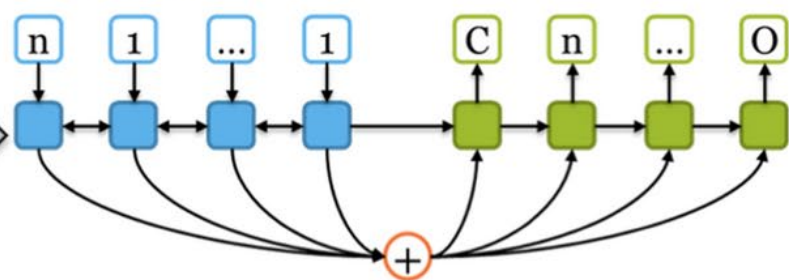
SMILES

Cn1c(=O)c2c(ncn2C)n(C)c1=O

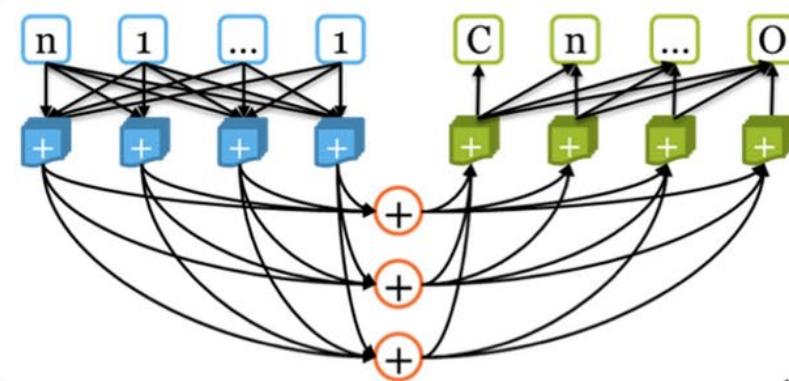
B. LSTM-BASE



C. LSTM+ATTN



D. Sequential Transformer



SMILES-based Generator

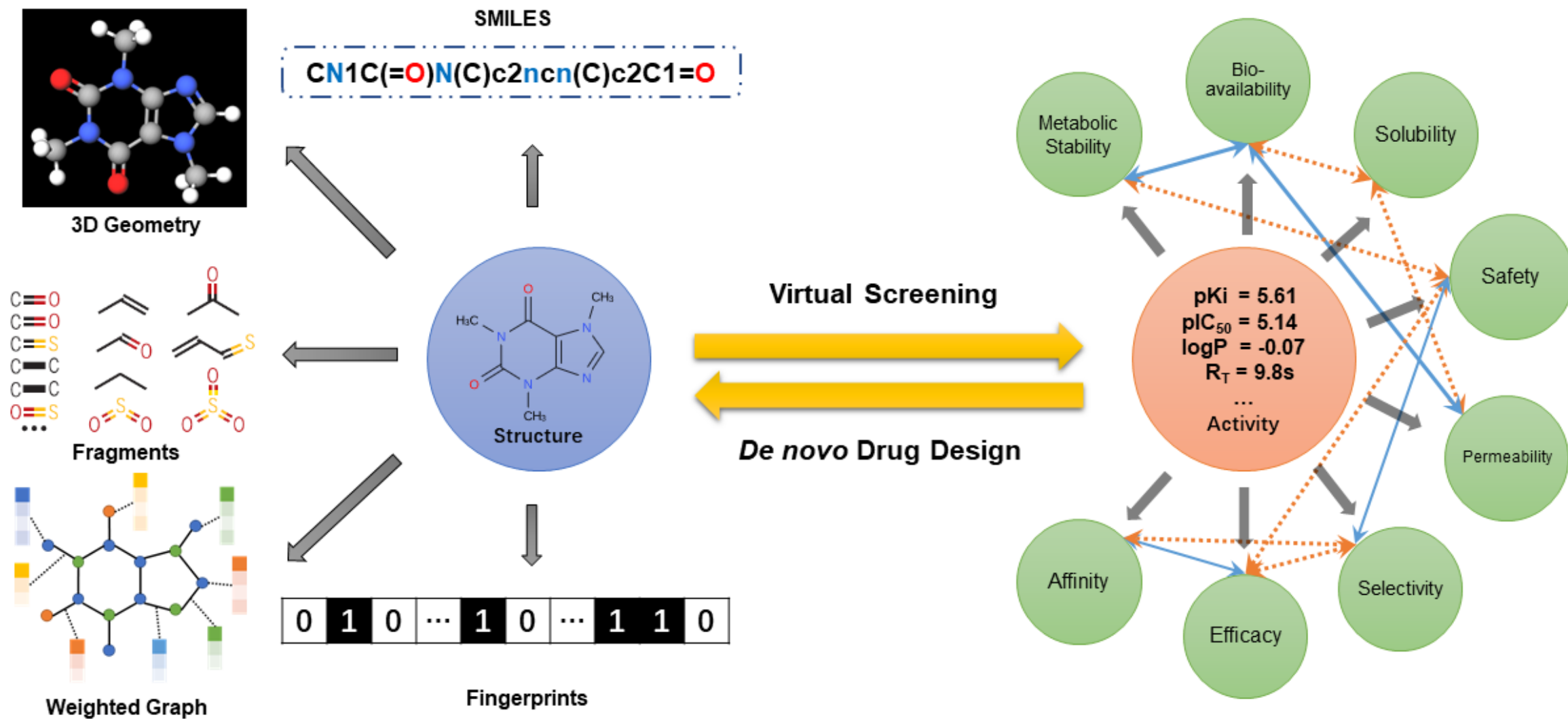
Ongoing work

- Select molecules for follow up by synthesis
- Multi-objective ranking to identify best molecules from batch of ~ 10.000
 - Maximize on-target affinity (e.g. adenosine A2A)
 - Minimize off-target affinity (e.g. hERG, other adenosine receptors)
 - Maximize Quantitative Estimation of Drug-likeness (QED)
 - Estimate synthetic accessibility (currently SA Score)

Take home messages de novo generation

- Machine learning can be used to learn the grammar of molecules in SMILES. It does this by learning the probabilities that in SMILES certain atoms follow other atoms in a sequence.
- After training the algorithm can be used to suggest new molecules that resemble the training set but are not the same.
- The power is in the numbers, 10, 100, 1000000 molecules can be generated

AI approaches in a ligand based world..



Proteochemometrics (PCM) de novo generation (DrugEx)



Willem Jaspers

LACDR