



Chemography approach to Chemical Space exploration

Alexandre Varnek

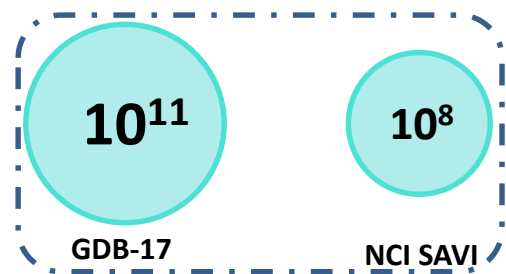
University of Strasbourg, France

ICReDD, Hokkaido University, Japan

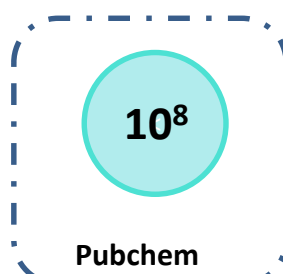
XXVII Symposium on Bioinformatics and Computer-Aided Drugm, 5th April 2021

Sizes of selected chemical data collections

Public

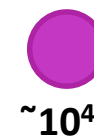


Virtual compounds

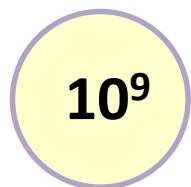


Real compounds

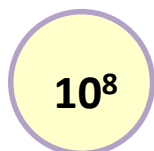
Approved drugs



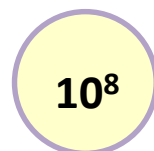
Commercial



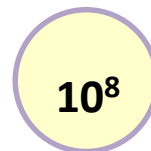
Enamine REAL Space



Enamine
REAL « Database »



ZINC15

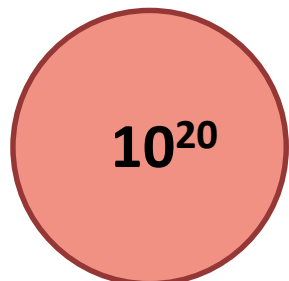


Sigma Aldrich
« in-stock »

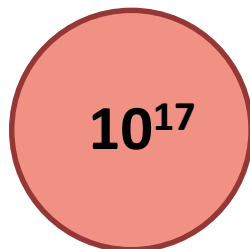


eMolecules Plus

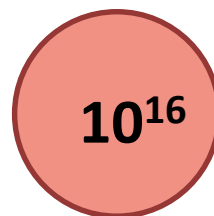
Proprietary



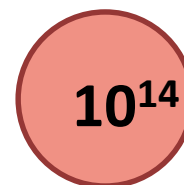
Merck MASSIV 2018



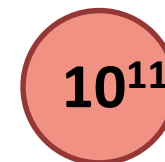
AstraZeneca 2018



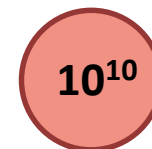
Evotec
EVOSpace 2016



Pfizer
PGVL 2008

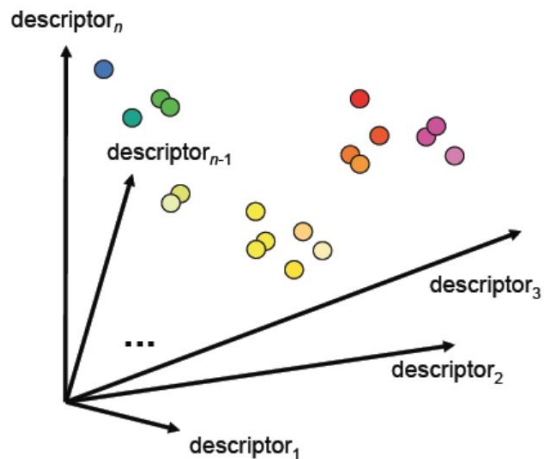


Boehr.-Ing.
BICLAIM 2012

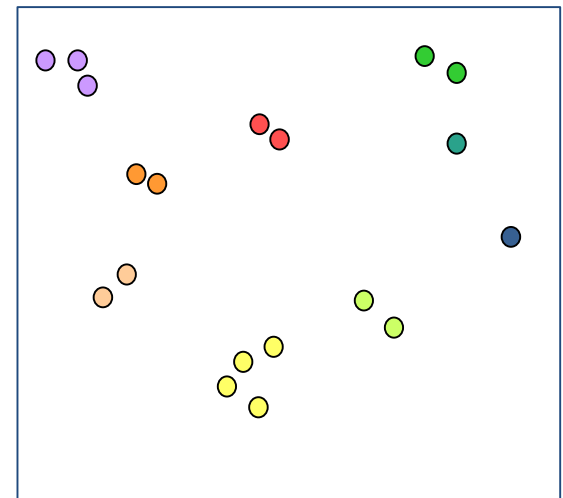


Lilly
LPC 2016

Data visualization: dimensionality reduction problem



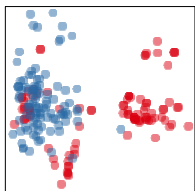
Data space
(N-dimensional)



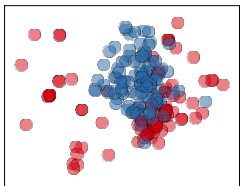
Latent space
(2-dimensional)

Dimensionality reduction methods

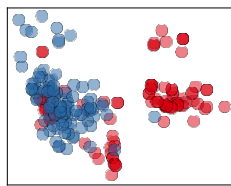
Acetylcholinesterase dataset (DUD) : 100 actives and 100 inactives
ISIDA descriptors



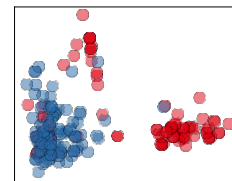
Multi-Dimensional
Scaling



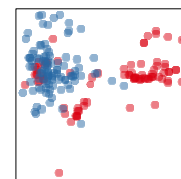
Canonical Correlation
Analysis



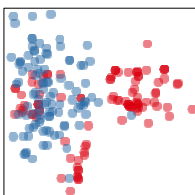
Independent
Component Analysis



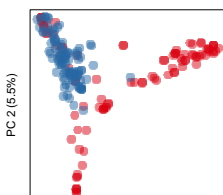
Exploratory Factor
Analysis



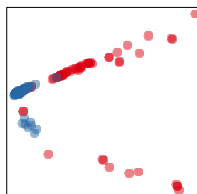
Sammon map



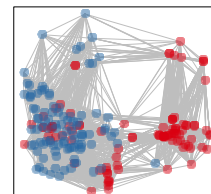
PCA



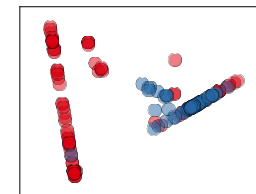
Kernel PCA
(RBF kernel)



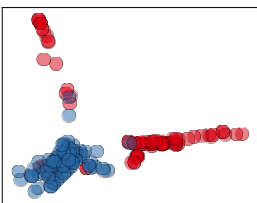
Kernel PCA
(polynomial kernel)



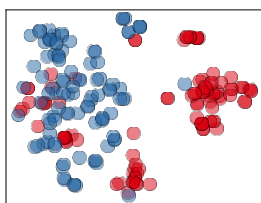
Isomap



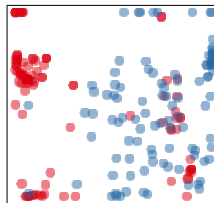
Locally Linear
Embedding



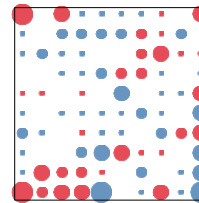
Laplacian Eigenmaps



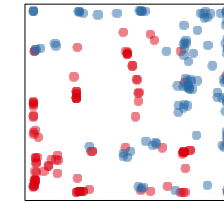
t-SNE



Autoencoder dimensionality
reduction



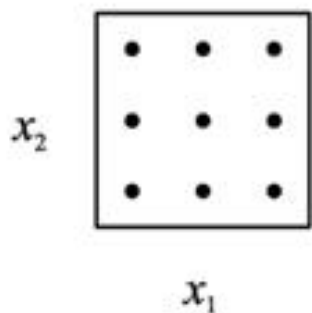
SOM



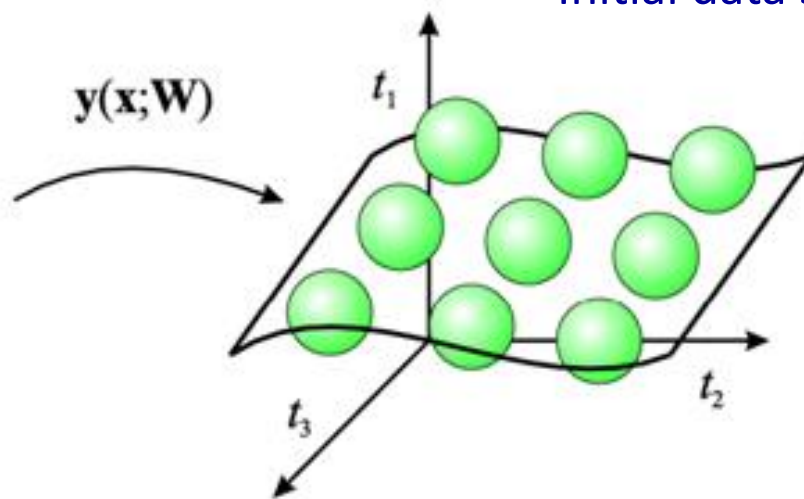
GTM

Generative Topographic Mapping

latent space



$y(x; \mathbf{W})$



Initial data space

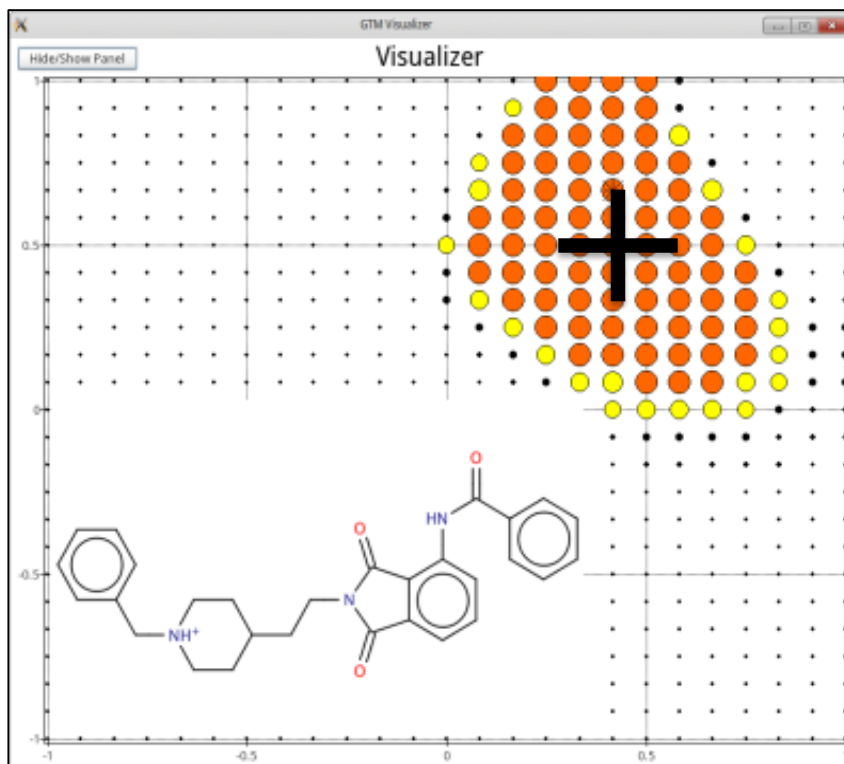
GTM generates a data probability distribution in ***both initial and latent data spaces***.

This opens an opportunity to use GTM not only to visualize the data but also for structure-property modeling tasks

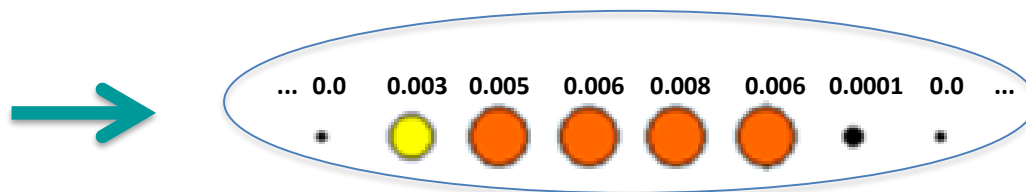
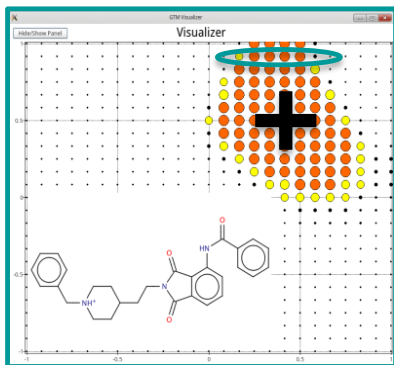
- C. M. Bishop *Pattern Recognition and Machine Learning*. 2006 Springer
- N. Kireeva. I.I. Baskin. H. A. Gaspar. D. Horvath. G. Marcou and A. Varnek. *Mol. Informatics*. 2012. 31. 201-312 5

GTM: Probability Density distribution in the latent space

Projection of an object on GTM is described by the probability distribution (*responsibilities*) over the lattice nodes.



GTM descriptors for molecules and datasets



Map resolution: $N_{nodes} = K * K$

Standard setting: $K = 25$, $N_{grid} = 625$

Molecule \rightarrow responsibilities' vector $\{R_{tk}\}$ of N_{nodes} length

Dataset \rightarrow normalized cumulated responsibilities' vector of N_{nodes} length

GTM landscapes



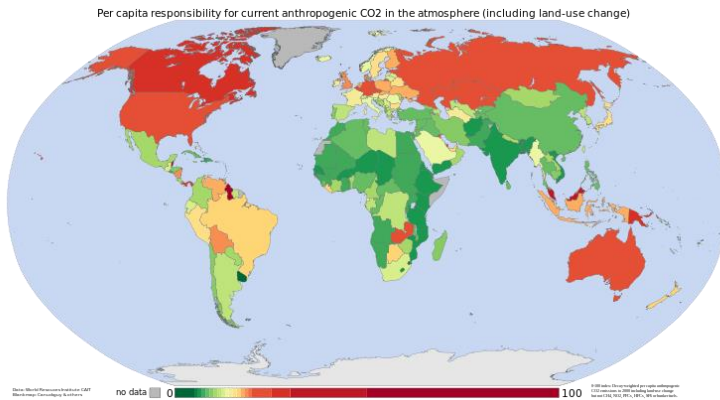
Properties mapping



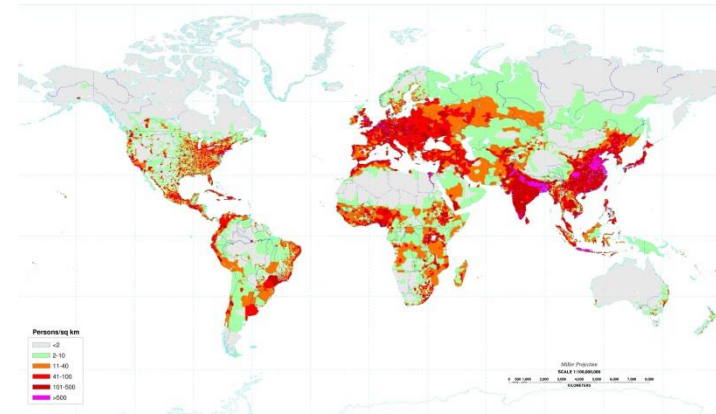
political map



physical map



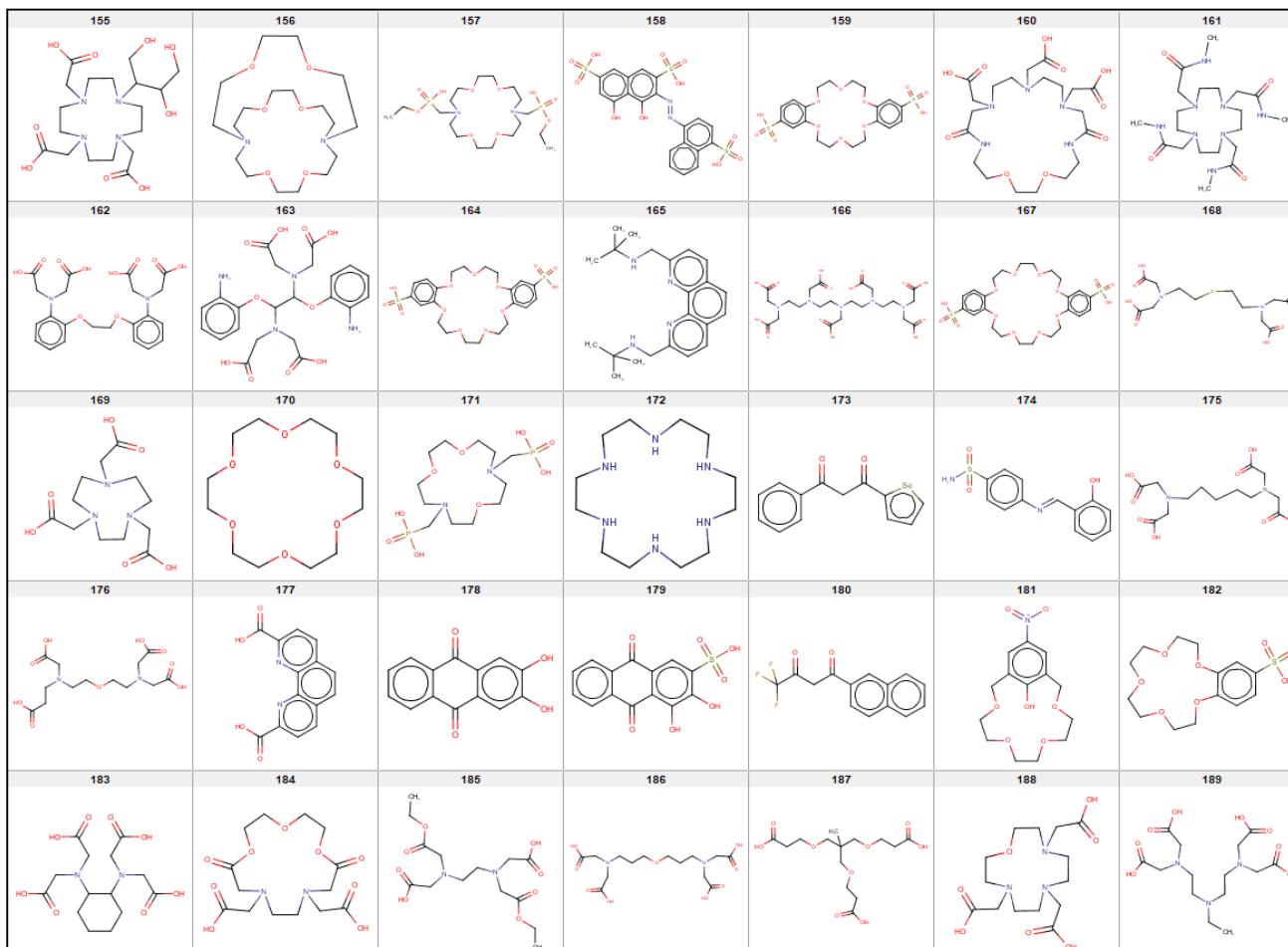
CO₂ emission



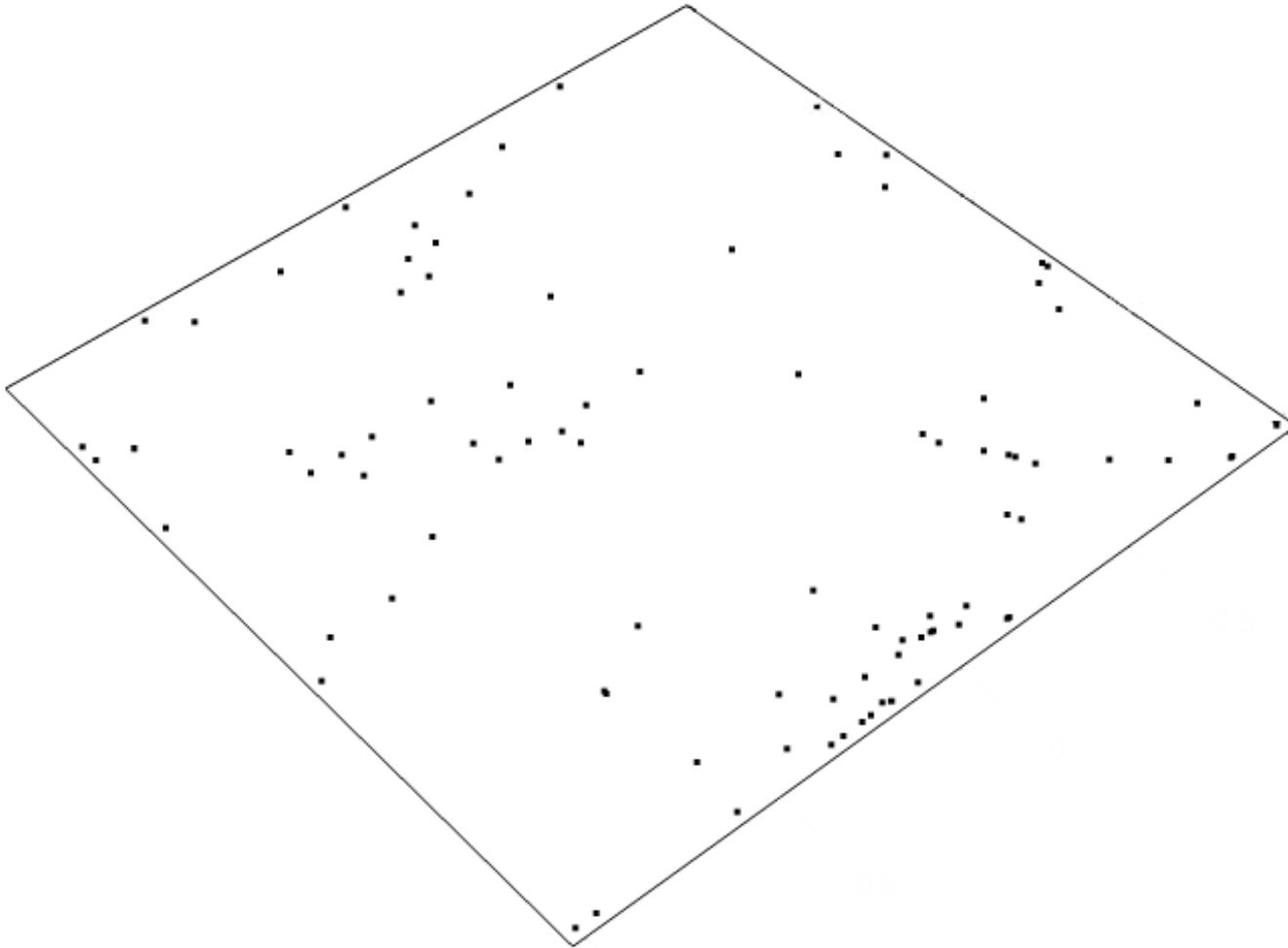
population density

Case study: chemical space of metal binders

Data: 102 organic molecules which complex the Lu^{3+} cation in water

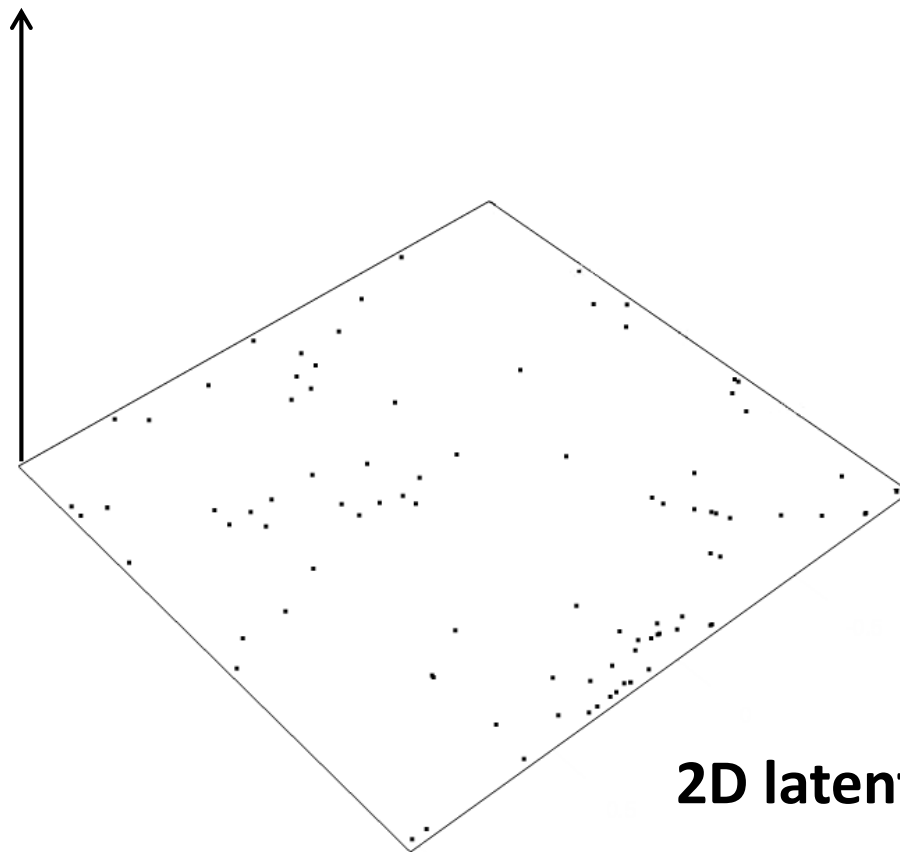


GTM of Lu³⁺ binders



Activity landscape for Lu³⁺ binders

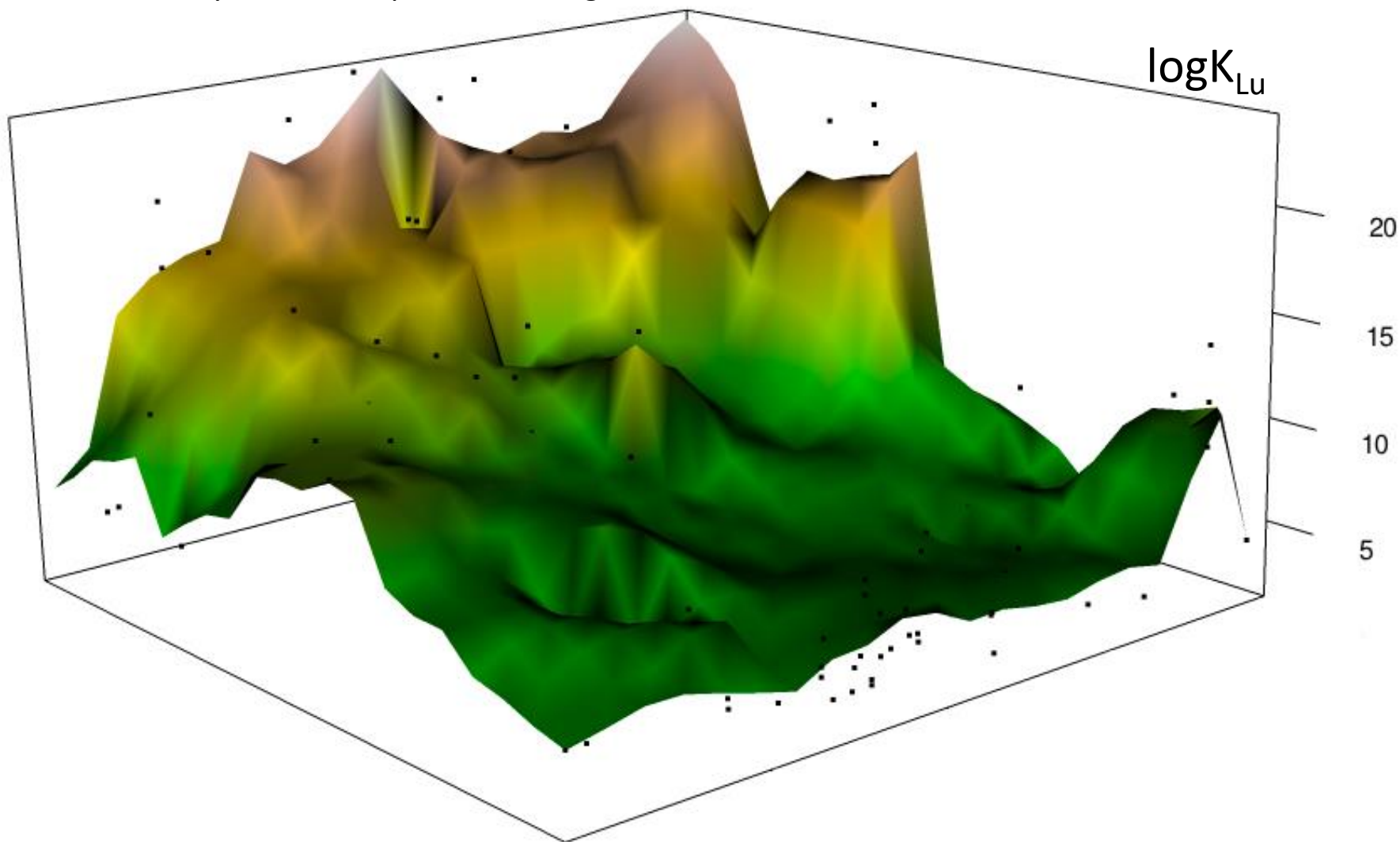
Activity

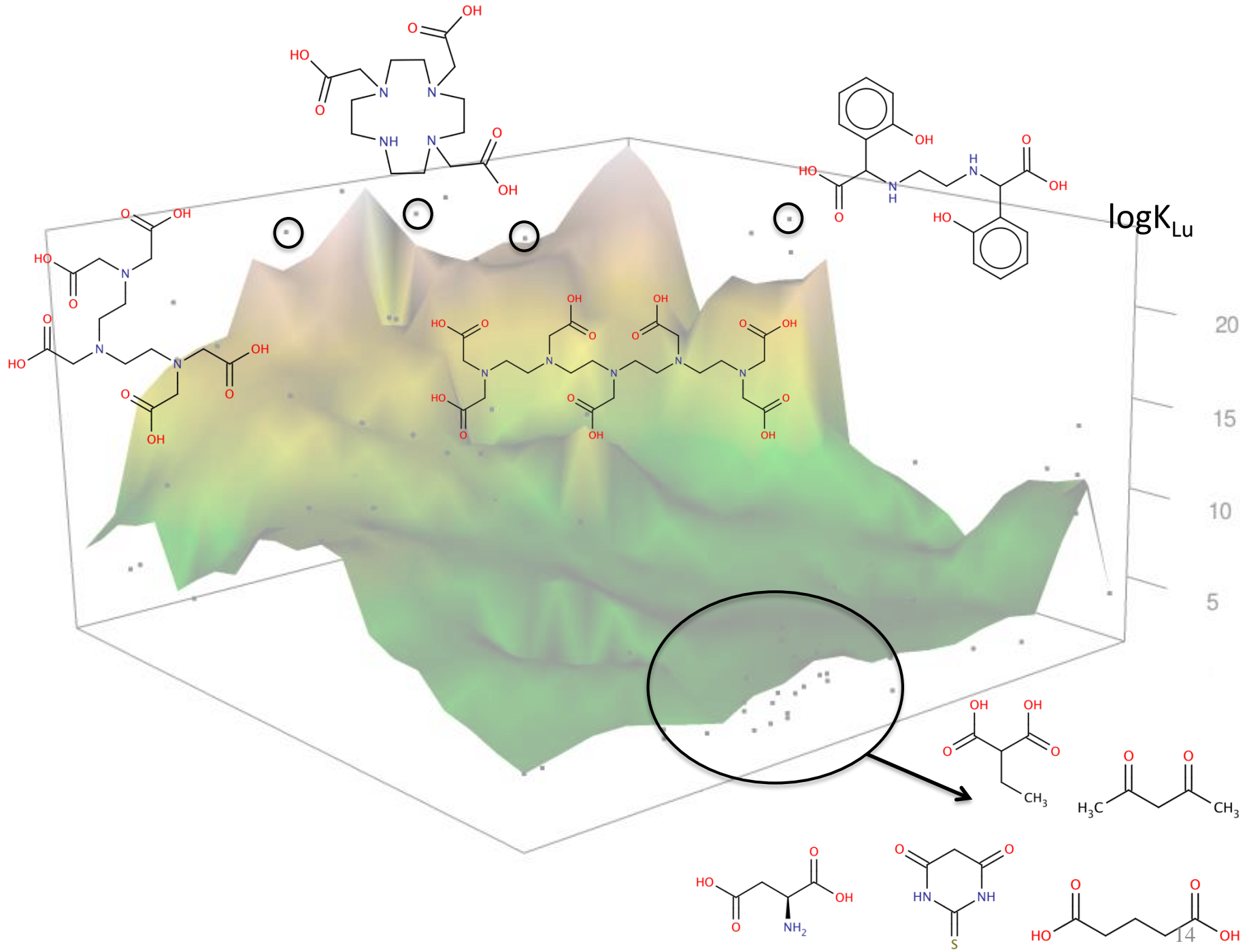


2D latent space

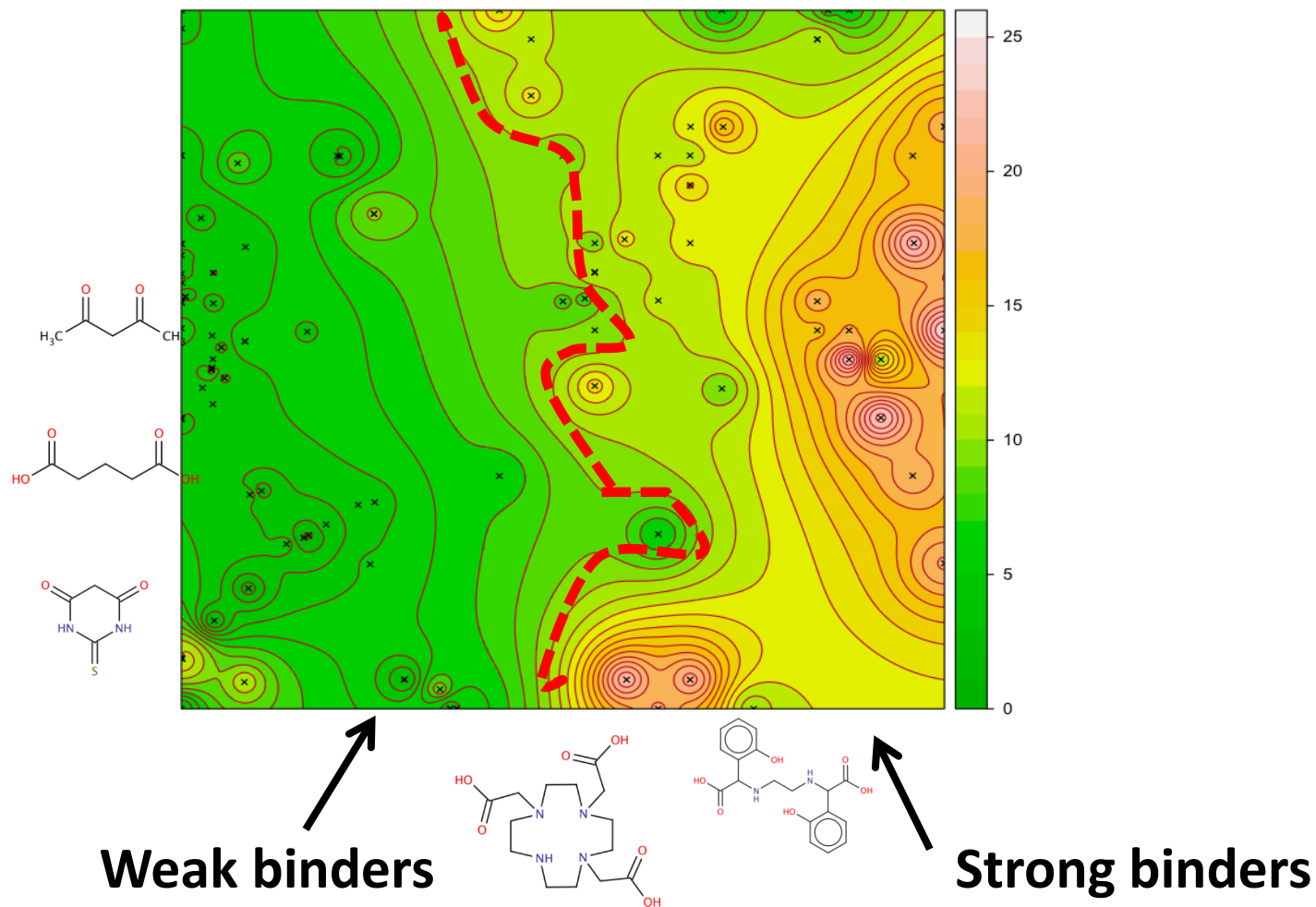
Activity landscape for Lu³⁺ binders

Stability of Lu³⁺ complexes with organic molecules



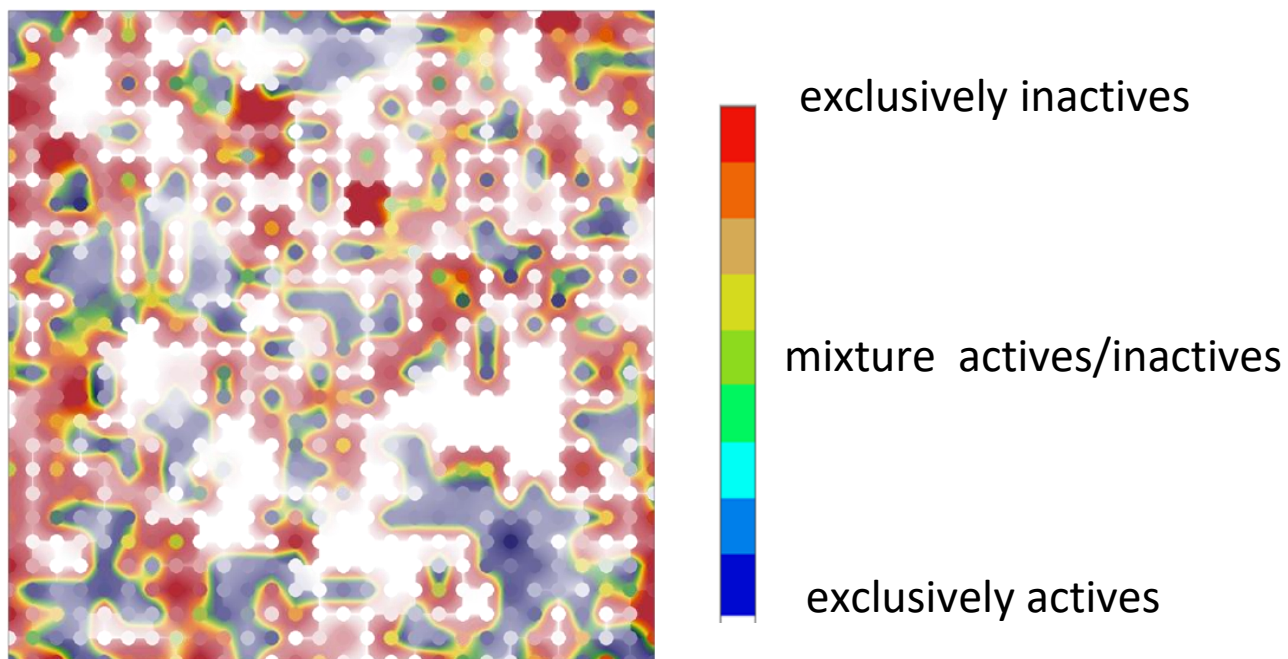


Activity landscape for Lu³⁺ complexation



Class landscapes

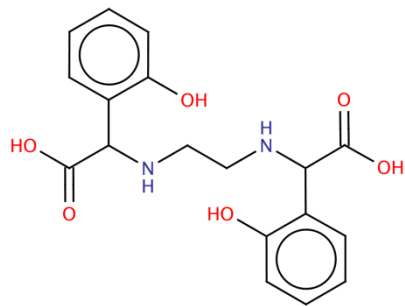
ChEMBL (1.7 M cmds) : class landscape of antiviral compounds



GTM Landscape as predictive tool

Regression

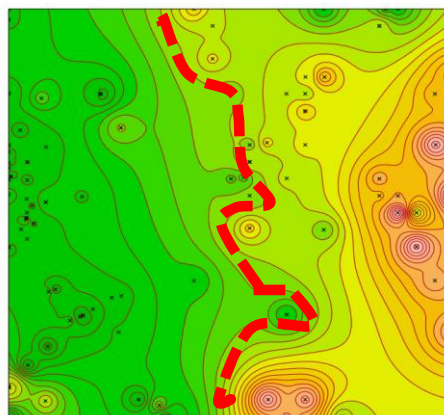
$$\hat{A}_j(\text{test}) = \sum_k \bar{A}_k(\text{training}) R_{jk}(\text{test})$$



new compound

Projection

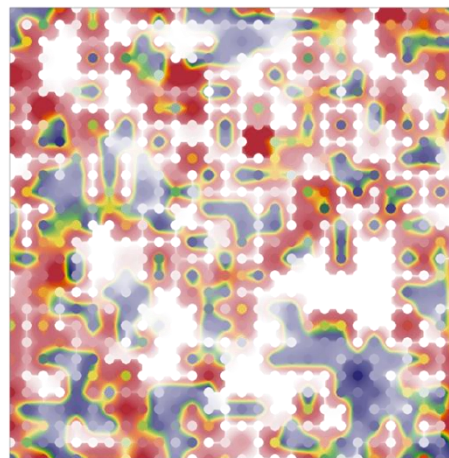
Activity landscape



Predicted activity value

Classification

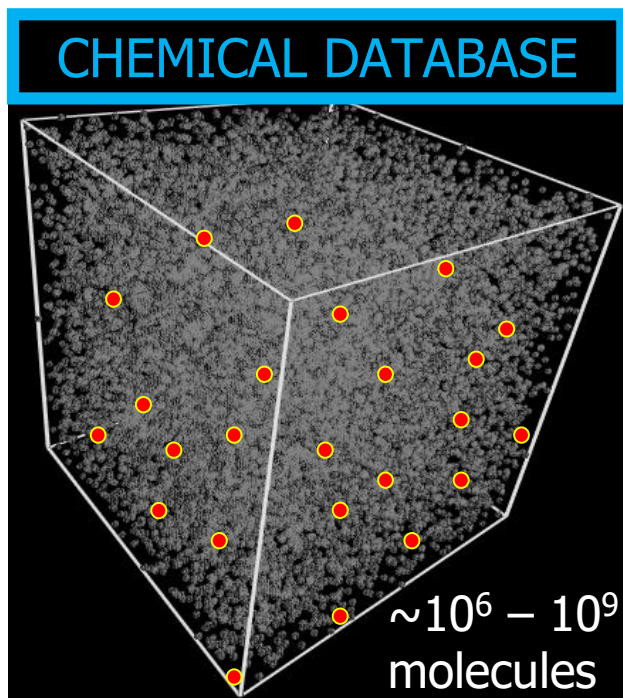
$$P(c_i | \mathbf{x}_k) = \frac{P(\mathbf{x}_k | c_i) \times P(c_i)}{\sum_i P(\mathbf{x}_k | c_i) \times P(c_i)}$$



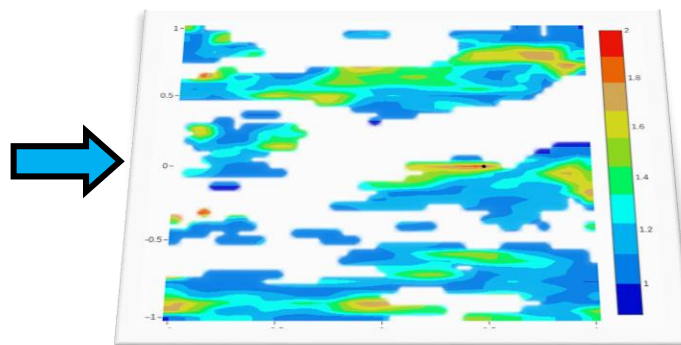
Class landscape

Predicted category
"active" or "inactive"

Universal maps: application to virtual screening



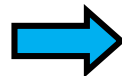
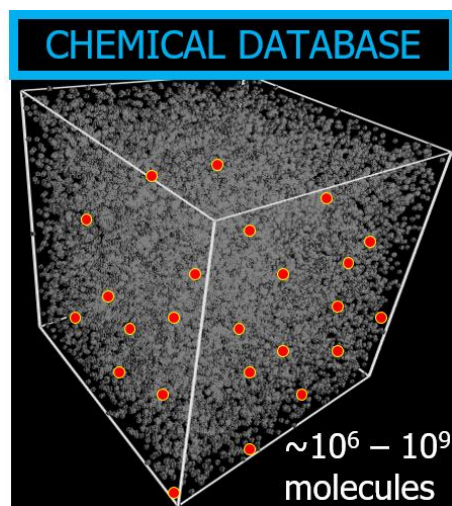
GTM activity or class landscape



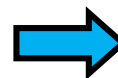
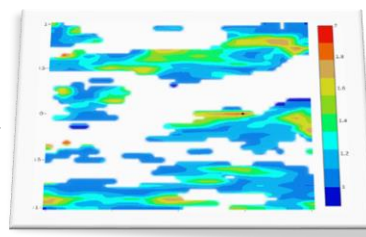
HITS

TRASH

Universal maps: application to virtual screening



GTM



In silico designed with GTM and experimentally validated compounds

- *Antiviral compounds*
- *Antimalarial compounds*
- *Solvents for Li-batteries*
- *Bromodomain (BRD4) inhibitors*

GTM : areas of application

Conformational space analysis

Data visualisation and analysis

Ligand to Protein docking

Libraries comparison

Sequences space analysis

Drugs repurposing

Structure-Activity modeling

Library design

Virtual screening



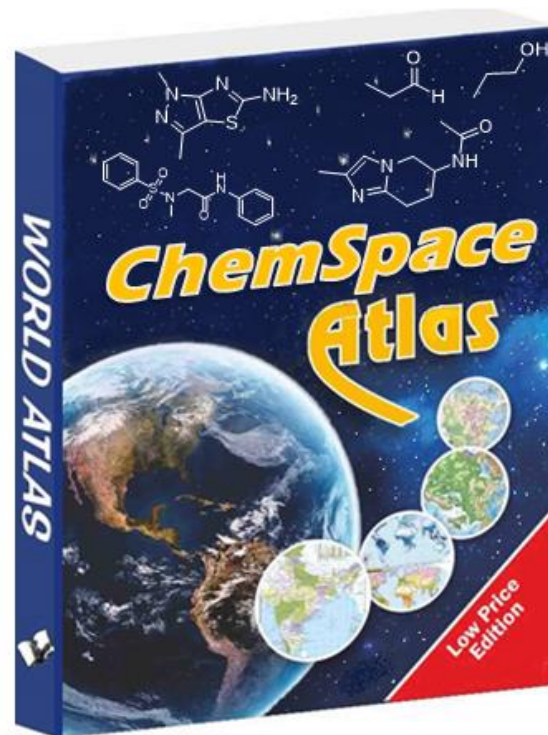
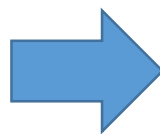
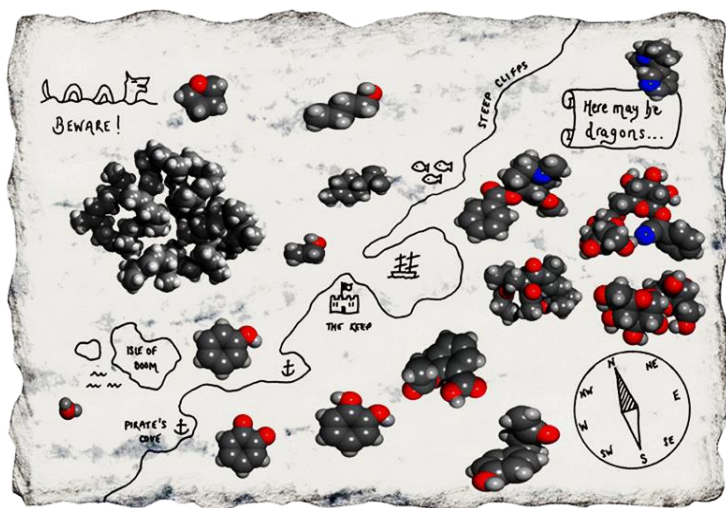
de novo design

GTM : case studies



- **Visualisation and analysis of ultra-large data**
- **Artificial Intelligence driven design of novel molecular structures and reactions**

ChemSpace Atlas



Universal Map of chemical space

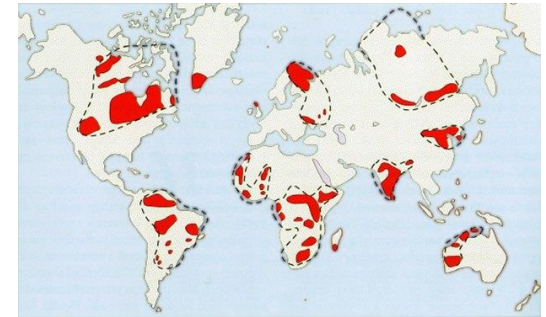
What do we expect from an “universal” map ?

Map of a chemical space is expected:

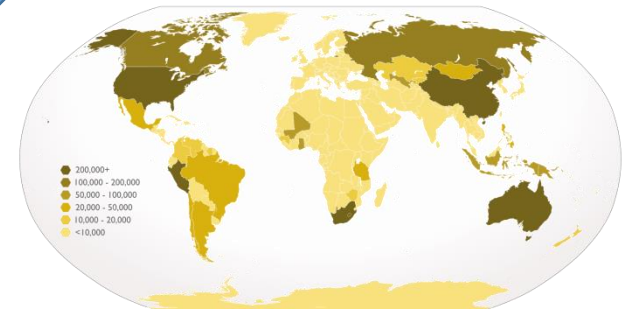
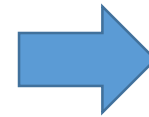
- to accommodate variety of known chemotypes;
- to be able to distinguish different activity classes;
- to separate actives and inactives within a given activity class;
- to be *neighbourhood behaviour (NB)* compliant, e.g., molecules grouped together are expected to display similar activities

«Universal» map

- Defines a frame of a relevant space
- Accommodates different landscapes



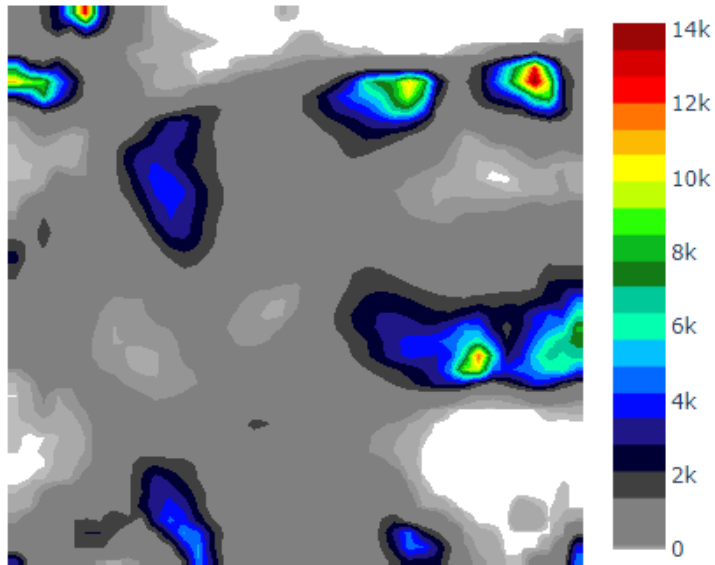
Highly prospective mineral regions



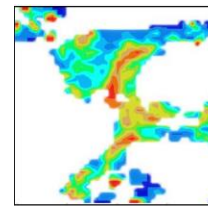
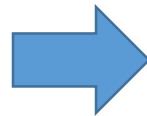
Gold production by country

«Universal» map

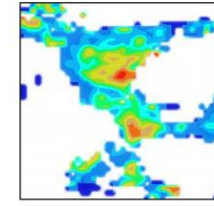
- Defines a frame of biological relevant chemical space
- ISIDA fragment descriptors are used
- Constructed on the basis of ChEMBL database compounds
- Predicts of > 700 biological activities



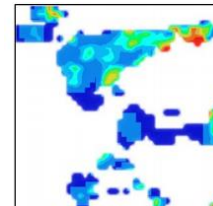
Density landscape of the ChEMBL database (1.7 M cmds)



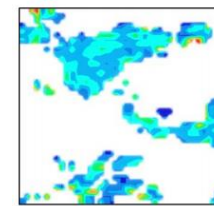
MAP kinase p38 alpha



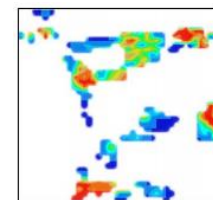
Vascular endothelial growth factor receptor 2



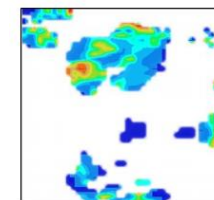
Cyclin-dependent kinase 2



Serine/threonine-protein kinase AKT



Phosphodiesterase 5A



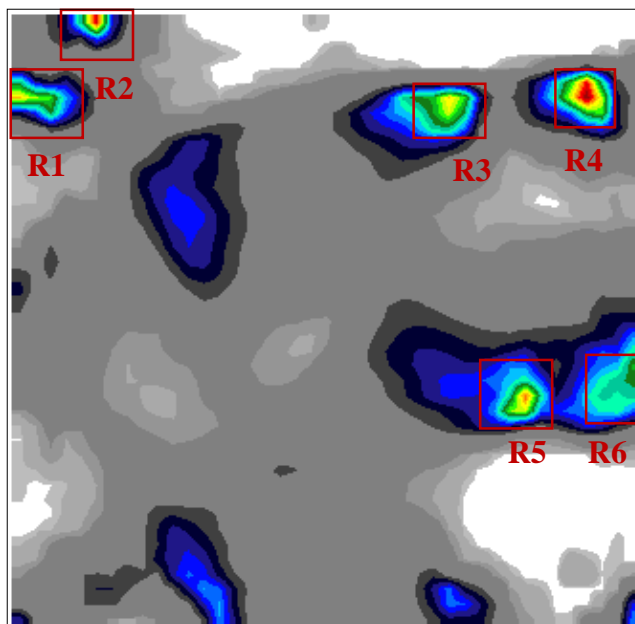
Adenosine A2a receptor

Active

Inactive

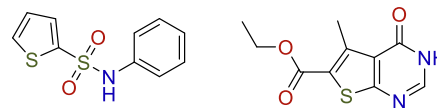
«Universal» map

chemotypes distribution

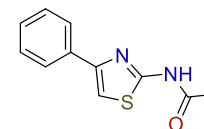


ChEMBL density landscape

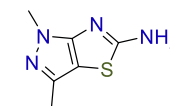
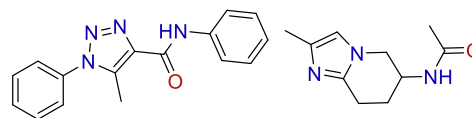
R1: Thiophenes



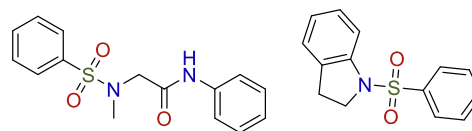
R2: Thiazoles



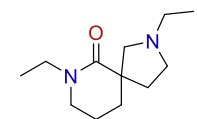
R3: Heterocyclic amides



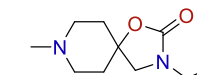
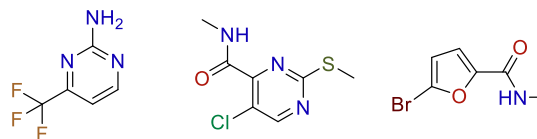
R4: Benzensulphonamides



R5: Spiro-heterocyclic amides and carbamites



R6: Halloginated heterocycles



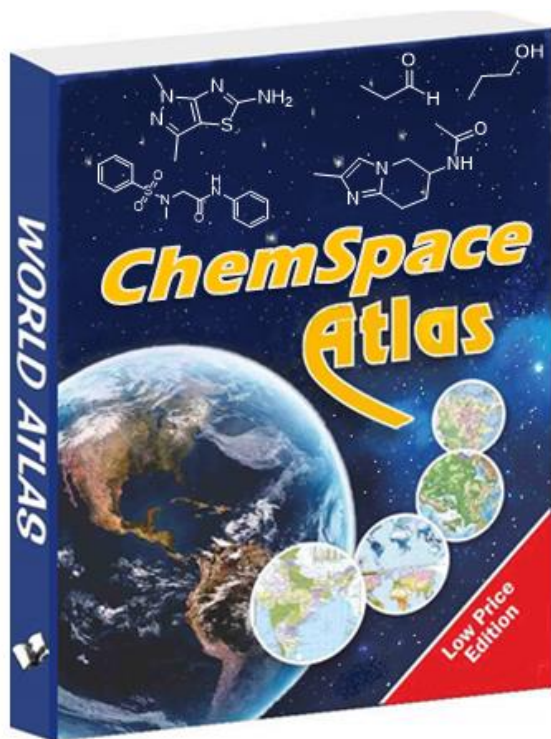
ChemSpace Atlas

Main features

- polyvalent tool based on the GTM Universal Map
- accomodates > 3 billion cmpds
- assembles > 40000 hierarchically related maps of different scale and > 1.5 million activity landscapes

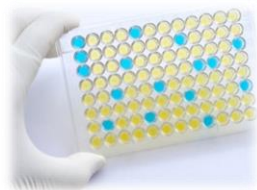
Main options

- Data visualization, search, subsets selection
- Automatized extraction of Maximal Common Structures
- Scaffold analysis
- Projection of new compounds
- Pharmacological profiling with respect to >700 biological targets

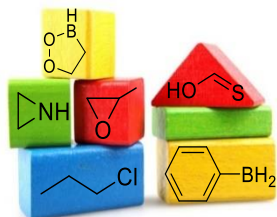


ChemSpace Atlas

The tool consists of 4 main parts:



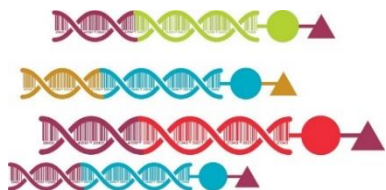
Screening Compounds



Building Blocks



Natural Products



DNA-Encoded Libraries

ChemSpace Atlas

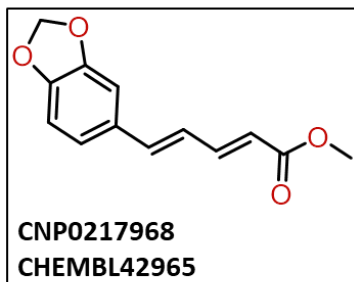
discovery of synthetic analogs of natural products



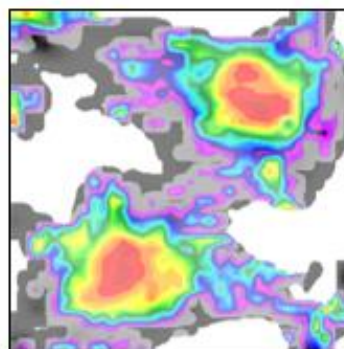
ChemAtlas NP database:

253 893 Natural Products + 586 235 synthetic analogs from ZINC

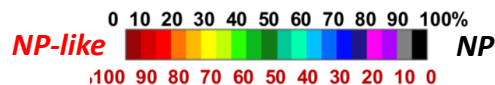
Query



Active against **Monoamine oxidase B**
(*Rattus norvegicus*)



**Structures and activity profiles
of synthetic analogues**





Dashboard

START HERE

Welcome

ANALYSIS

Chemspace tracker

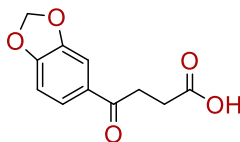
Activity prediction

View: By Targets By Compounds

DOWNLOAD

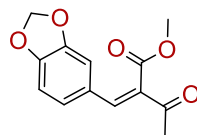
Compounds

Predicted activity (Target ID)



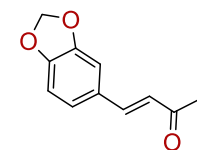
ZINC000000040327

[CHEMBL2039](#) : Monoamine oxidase B (*Homo sapiens*)



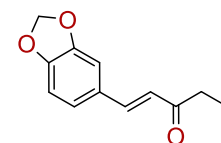
ZINC000016138715

[CHEMBL2993](#): Monoamine oxidase B (*Rattus norvegicus*)



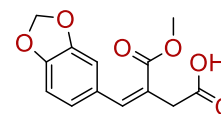
ZINC000001754404
ZINC000012417143

[CHEMBL2039](#) : Monoamine oxidase B (*Homo sapiens*)



ZINC000002015852
ZINC000020232188

[CHEMBL2039](#) : Monoamine oxidase B (*Homo sapiens*)
[CHEMBL4376](#) : Dual-specificity tyrosine-phosphorylation regulated kinase 2 (*Homo sapiens*)



ZINC000005218035

[CHEMBL2993](#): Monoamine oxidase B (*Rattus norvegicus*)



Dashboard

START HERE

Welcome

ANALYSIS

Chemspace tracker

Activity prediction

View: By Targets By Compounds

DOWNLOAD

ChEMBL Target ID	Name of the target	Number of predicted hits	
CHEMBL2039	Monoamine oxidase B <i>Homo sapiens</i>	256	See the hit list
CHEMBL2993	Monoamine oxidase B <i>Rattus norvegicus</i>	76	See the hit list
CHEMBL312	Arachidonate 5-lipoxygenase <i>Rattus norvegicus</i>	33	See the hit list
CHEMBL2003	Retinoic acid receptor gamma <i>Homo sapiens</i>	29	See the hit list
CHEMBL242	Estrogen receptor beta <i>Homo sapiens</i>	24	See the hit list
CHEMBL4376	Dual-specificity tyrosine-phosphorylation regulated kinase <i>Homo sapiens</i>	19	See the hit list
CHEMBL2186	Carbonic anhydrase XIII <i>Mus musculus</i>	14	See the hit list
CHEMBL1860	Thyroid hormone receptor alpha <i>Homo sapiens</i>	11	See the hit list
CHEMBL5339	G-protein coupled receptor 120 <i>Homo sapiens</i>	10	See the hit list
CHEMBL324	Serotonin 2c (5-HT2c) receptor <i>Rattus norvegicus</i>	10	See the hit list

Chemography: Searching for Hidden Treasures

Yuliana Zabolotna, Arkadii Lin, Dragos Horvath, Gilles Marcou, Dmitriy M. Volochnyuk,
and Alexandre Varnek*

J. Chem. Inf. Model. 2021, 61, 1, 179–188



Initial gold-bearing ore



Gold-enriched ore



Pure gold

Commercial vs Biologically relevant data

**Commercially available
chemotypes**

ZINC
DATABASE

>1.3 billion cmpds

**Biologically
relevant
chemotypes**

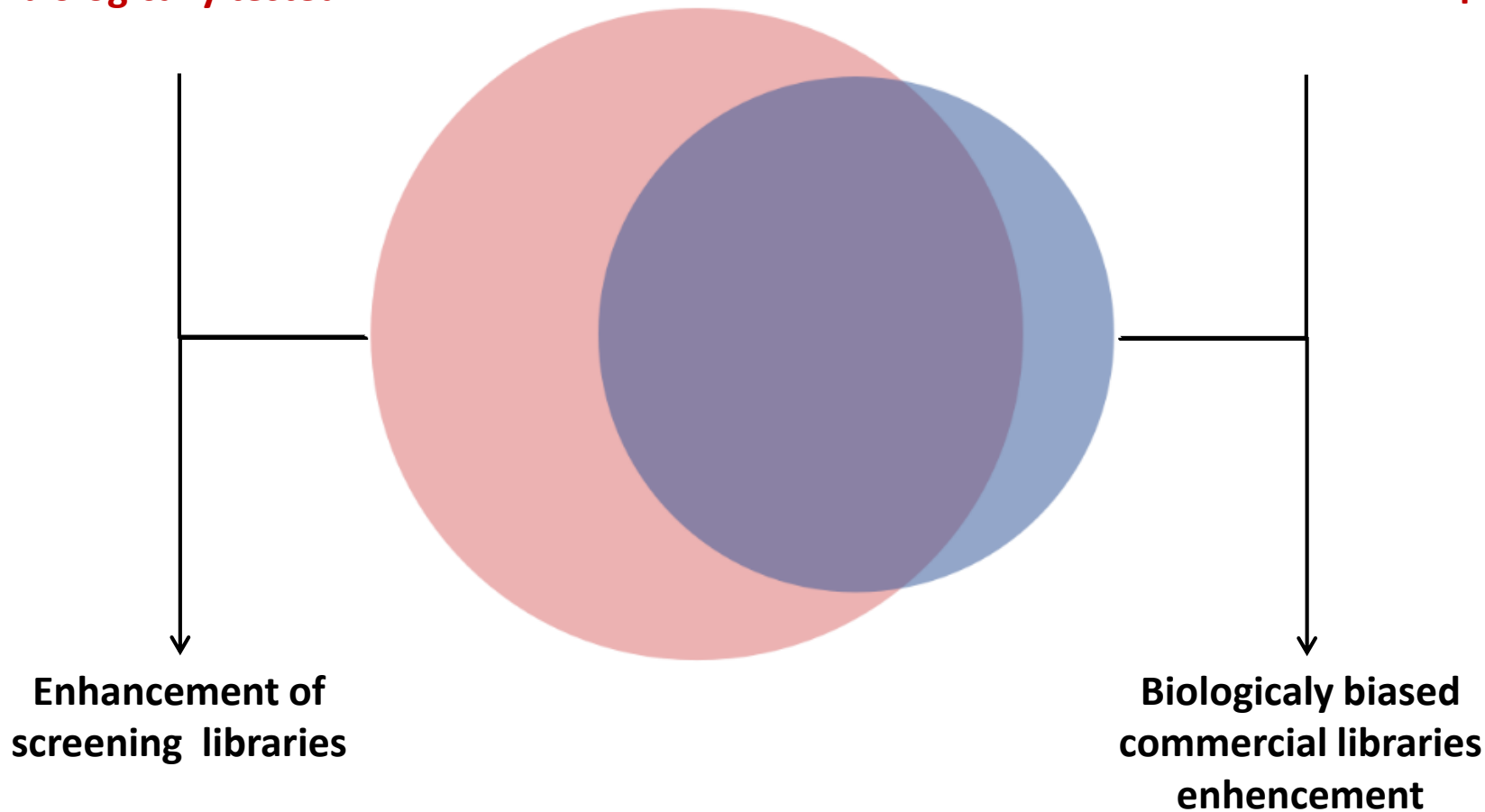
ChEMBL 

>1.8 M cmpds

Commercial vs Biologically relevant data

**Chemotypes never been
biologically tested**

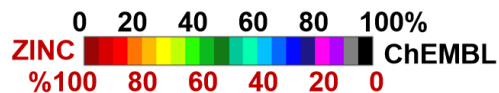
**Chemotypes missing in the
commercial chemical space**



**Enhancement of
screening libraries**

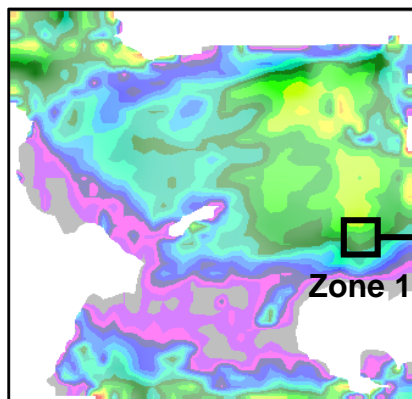
**Biologically biased
commercial libraries
enhancement**

Hierarchical GTM navigation of the chemical space



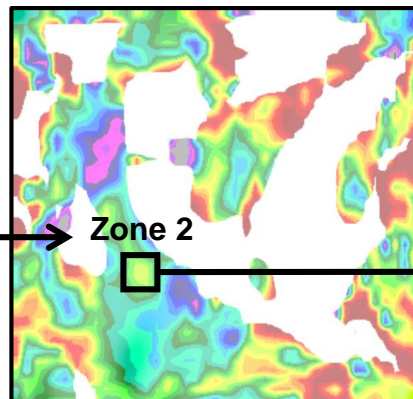
Maximum Common Substructure (MCS)

Universal map



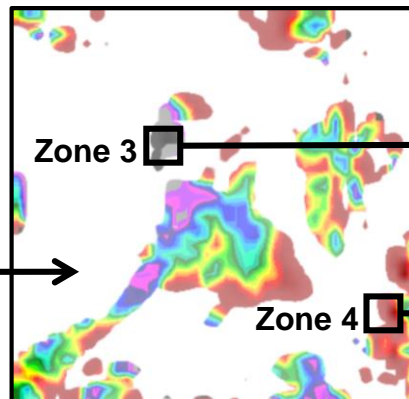
#compounds 3 614 394

Level1



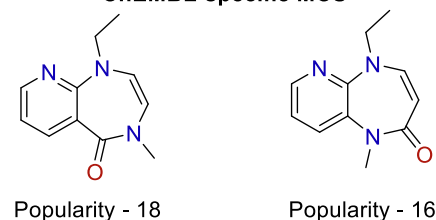
82 246

Level2

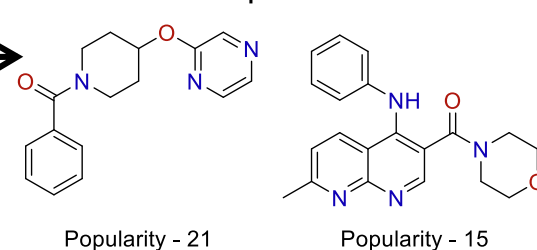


4 230

ChEMBL-specific MCS



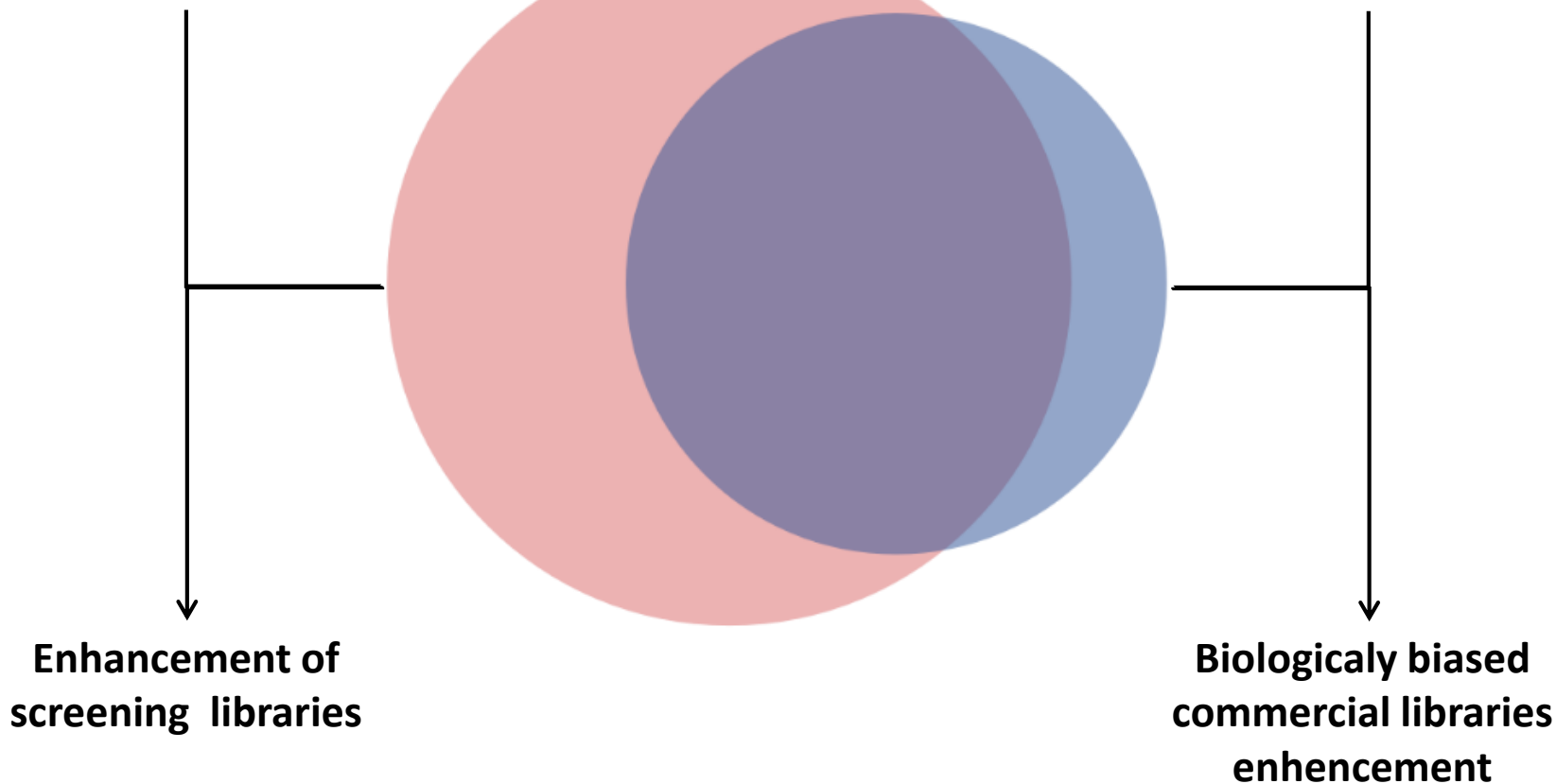
ZINC-specific MCS



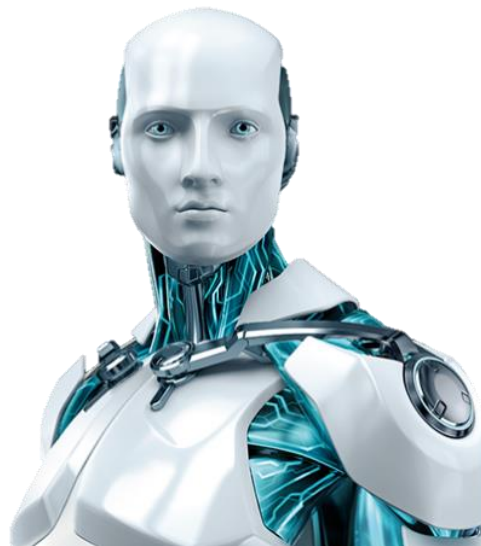
Commercial vs Biologically relevant data

> 100K chemotypes
never been biologically tested

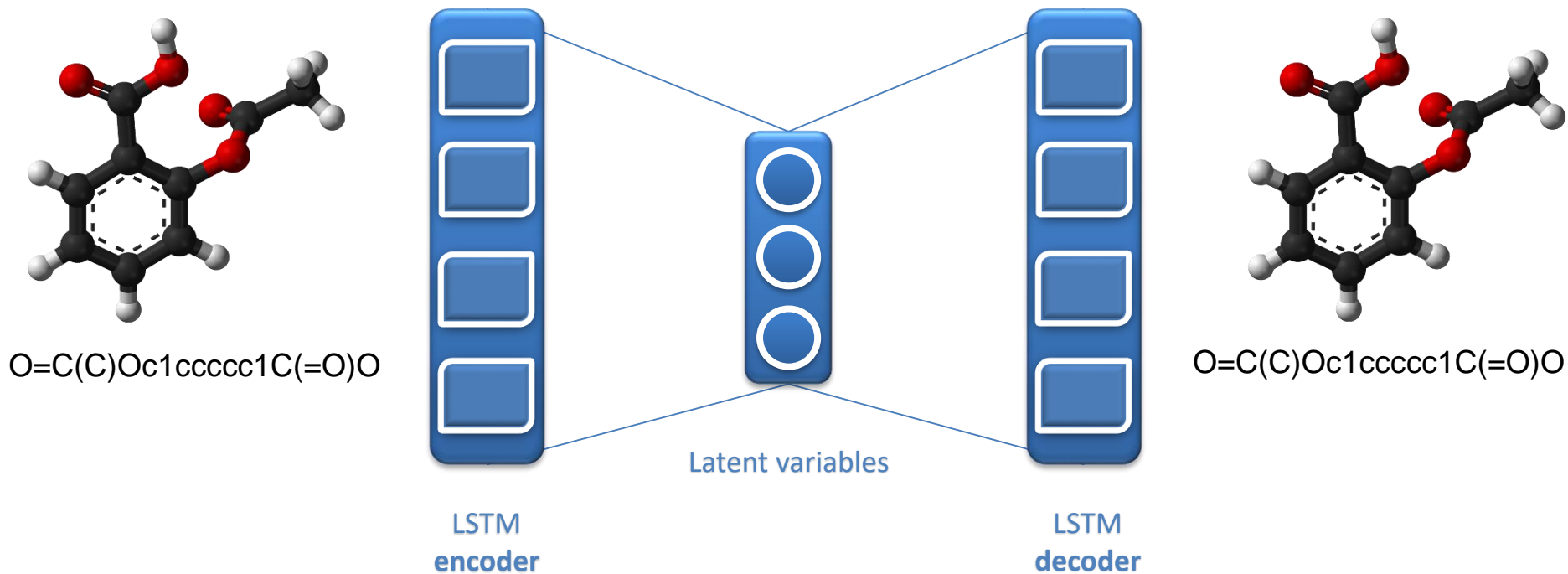
>20 K chemotypes missing in
the commercial chemical space



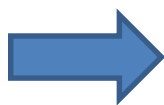
De novo design of molecules
possessing desirable biological activity



Autoencoder performing SMILES reconstruction



**Chemical
structure**



**Real numbers
encoding**

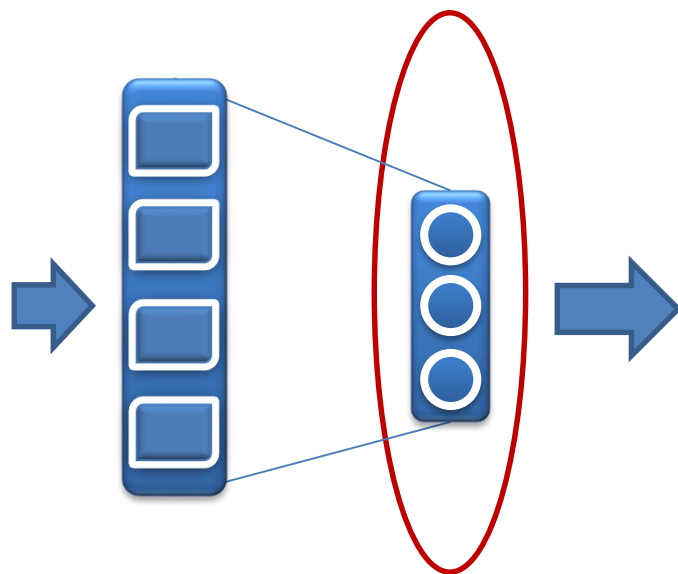


**Chemical
structure**

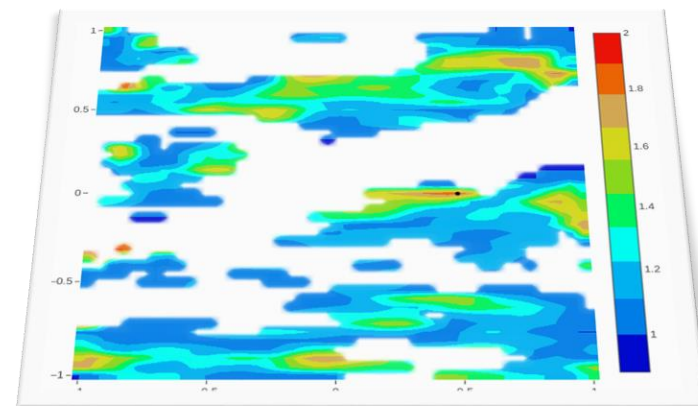
Building GTM on latent variables of autoencoder

Latent variables
(*vector on real numbers*)

Chemical
Database
(SMILES)

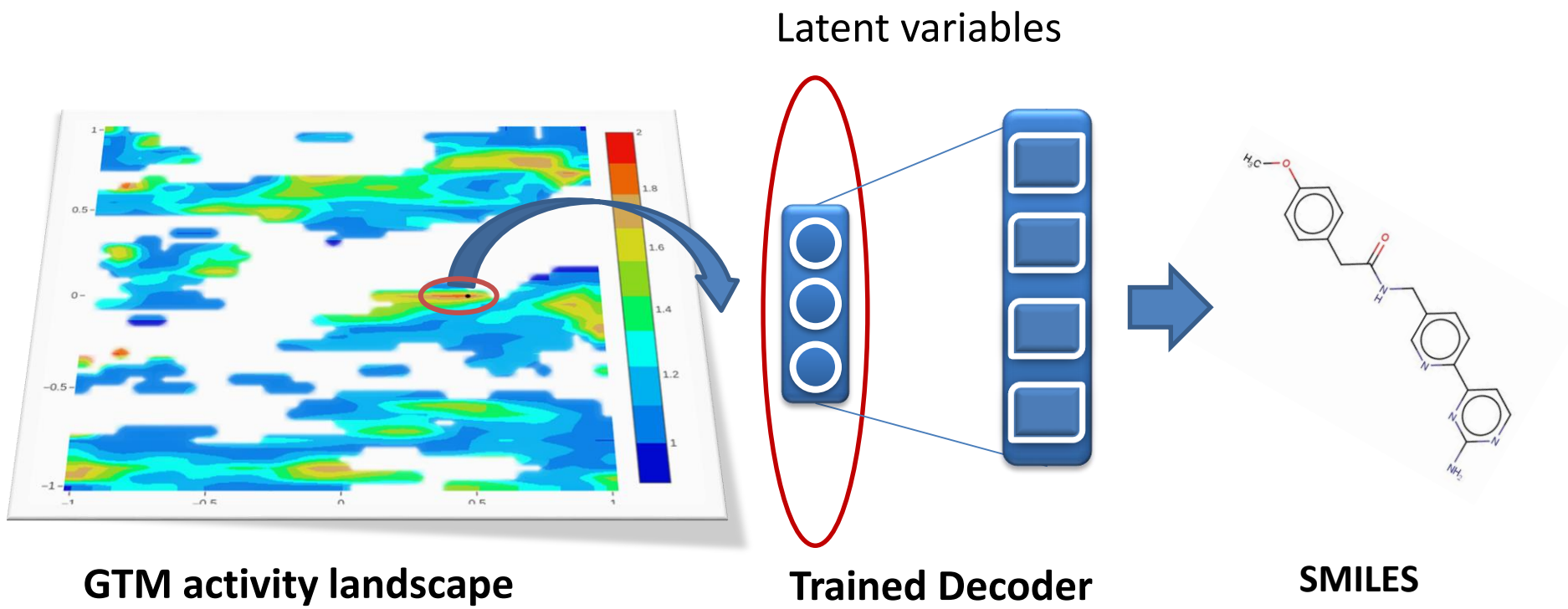


Trained Encoder

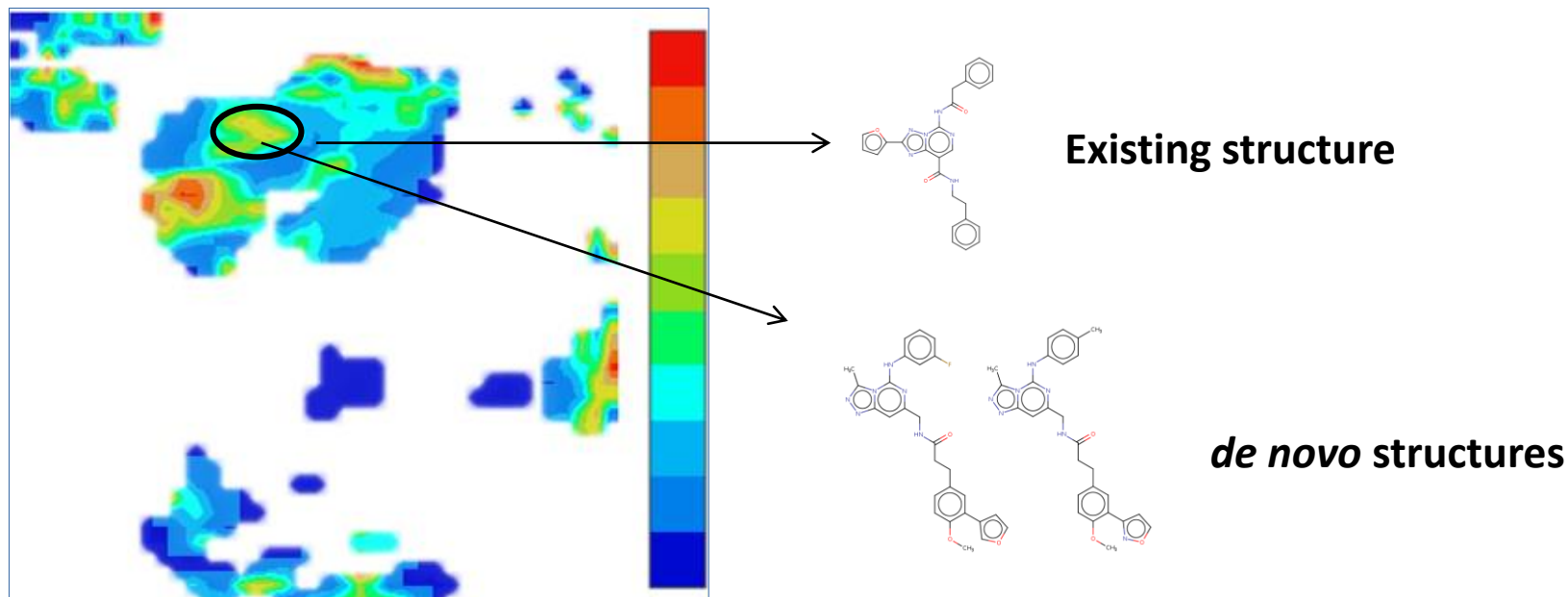


GTM

Generation of novel structures from specific areas of the map

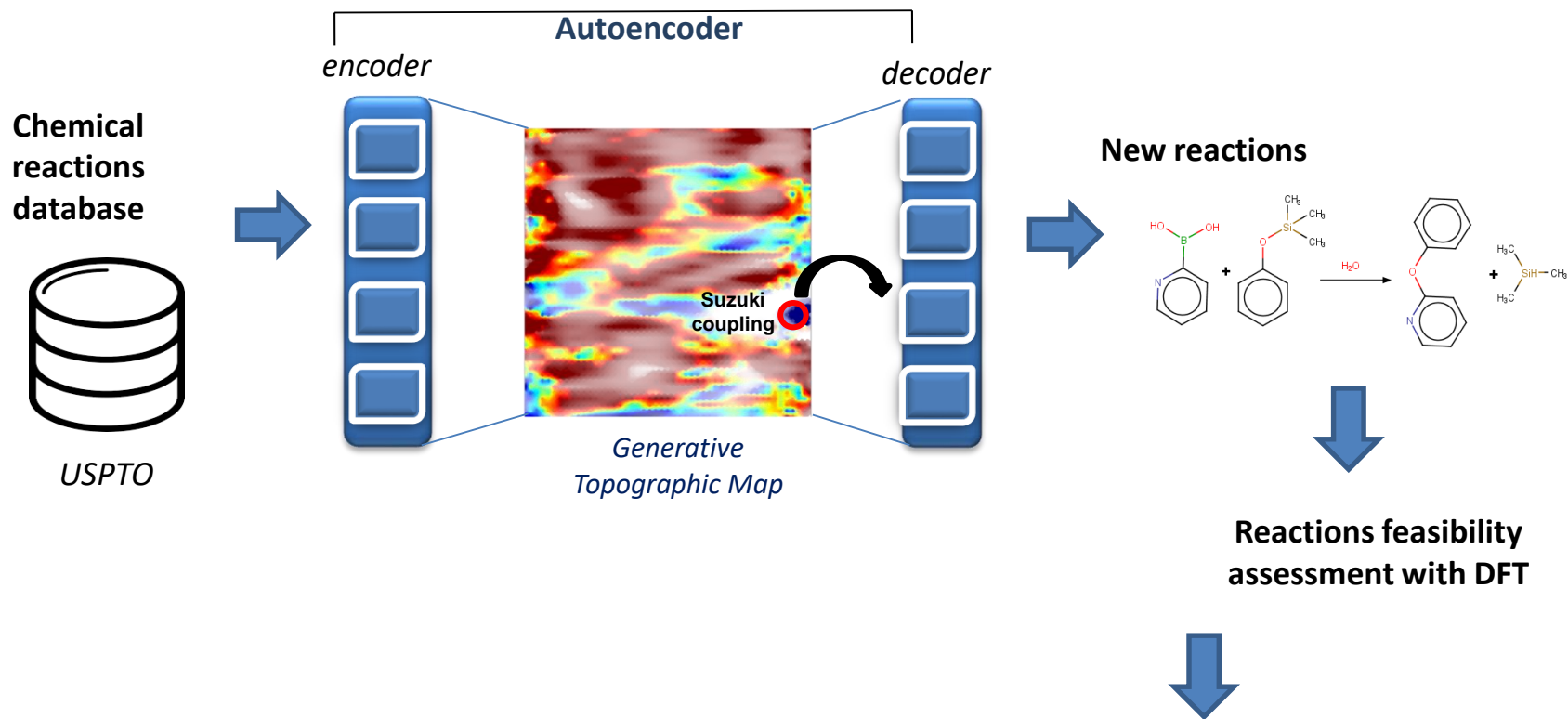


Case study: Generation of inhibitors of A2a receptor

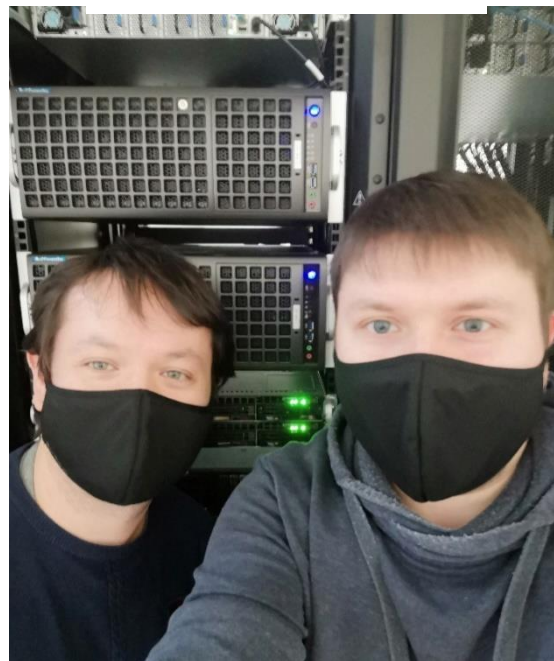


- **Generated structures are enriched with new scaffolds**
- **According to docking experiments they are efficiently able to bind A2a**

AI-driven design of new types of chemical reactions



- **13 new (with respect to the training data) Suzuki-like reactions have been detected**
- **5 of them have been found in recent publications**



Collaboration

- **ITN Marie-Curie BigChem**
- **Federal University of Kazan**
- **Chumakov Research Center RAS**

- **Enamine**
- **eMolecules**
- **Janssen Pharmaceutical**
- **TOTAL**
- **SOLVAY**