

# Chemography concept in chemical space exploration

#### Alexandre Varnek

University of Strasbourg

8ADD workshop Olomouc, 27th January 2025



# **Chemical Space - definition**



# What is « Chemical Space »

- ensemble of molecules ?
- mathematical object ?

Chemical space can be viewed as being analogous to the cosmological universe in its vastness, with chemical compounds populating space instead of stars *C. Lipinski and A. Hopkins, 2004* 

A set of molecules forms a chemical space, for which the relationships between the objects (graphs of descriptor vectors) are established *A.Varnek & I. Baskin, 2011* 

One can formulate chemical space as a mathematical and usually high-dimensional space where distances represent similarities between molecules which can be represented in the form of chemical space maps by applying various dimensionality reduction methods

J. L Reymond, 2025

# **Chemical Space representations**



# **Chemical Space representations**

#### **Graph-based chemical space**

- Scaffolds analysis
- Molecular Matched Pairs
- Activity cliffs
- Chemical space networks (e.g., MMP-based)

#### Vector-based chemical space

- Data visualization and analysis in descriptors space
- Chemical space networks (e.g., similarity-based)

# In silico drug design: Big Data problem



# **Chemography: efficient solution for Big Data handling**



#### Data visualization: dimensionality reduction problem



#### Initial chemical space (N-dimensional)

Latent chemical space (2-dimensional)

# **Dimensionality reduction methods**

Acetylcholinesterase dataset (DUD) : 100 actives and 100 inactives











Multi-Dimensional Scaling

Canonical Correlation Analysis Independent Component

Analysis

**Exploratory Factor** Analysis

Sammon map









Isomap



Locally Linear Embedding





t-SNE

Laplacian Eigenmaps



Autoencoder dimensionality reduction







## **Generative Topographic Mapping : areas of application**



#### **Tunable ISIDA fragment descriptors**



# Generative Topographic Mapping (GTM)



## **GTM Landscapes**





#### Class landscape





C<sub>3</sub>

 $C_4$ 

 $C_5$ 

c<sub>2</sub>

#### **Activity landscape**





<b>p</b> <sub>1</sub>	р <sub>2</sub>	p <sub>3</sub>	p <sub>4</sub>	$\mathbf{p}_{5}$
-----------------------	----------------	----------------	----------------	------------------

300

200

Each landscape can be encoded by a set of special GTM-descriptors characterizing either *structures* only or *structures* & *activity* distributions

#### **GTM Landscapes as predictive models**



## **GTM:** Chemical space analysis



**Density landscape of ChEMBL** (1.8 M cpnds, ISIDA descriptors)

#### Each zone of density landscape can be associated with some chemotypes

A4. Halogenated N-heterocycles

### **GTM:** density and activity landscapes for ChEMBL



- identify the most similar ChEMBL compounds,
- predict its pharmacological profile for > 700 biological activities

#### Pairwise libraries comparison

Task : identification of structural motifs unique for a given library



#### Data analysis: histograms vs chemography



#### Case study 1 : Commercial vs Biologically relevant data



#### Hierarchical GTM navigation of the chemical space



Maximum Common Substructures (MCS)





NH O

#### **Commercial vs Biologically relevant data**





#### Chemspace Atlas: Multiscale Chemography of Ultralarge Libraries for Drug Discovery

Yuliana Zabolotna, Fanny Bonachera, Dragos Horvath, Arkadii Lin, Gilles Marcou, Olga Klimchuk, and Alexandre Varnek\*

Cite This: https://doi.org/10.1021/acs.jcim.2c00509





## Case study 2: Proprietary Library Reshaping



#### Sigma-Aldrich



Boehringer

### Case study 3: Freedom space (5B) vs REAL space (40B)



Delivering Discovery Solutions®



- GTM shows very little overlap between Freedom and REAL spaces
- Freedom space is more enriched with drug-like compounds

## **Multiple libraries analysis**

Task : selection a library most similar to the reference one



#### **Generation and analysis of general-purpose DELs**



**DNA-Encoded Library:** combinatorial collection of small molecules covalently attached to the short DNA tag

#### DEL challenge

#### **Screening libraries**



#### **DNA-encoded libraries**



Parallel screening in separate "wells"

Individual compounds may be cherry-picked

Simultaneous screening in a single tube



Entire library as an object must be considered

#### **DNA-Encoded Libraries (DEL)**



#### How to select an optimal DEL for a particular drug discovery task?

#### GTM-based similarity assessment DEL / reference library (ChEMBL)



# Focused Library design: case study



**79.000** Building blocks from eMolecules





- **2500** DELs designed (size: 1M-1B)
- **2.5 B** compounds generated (1M cmpds per DEL)

**2500** comparative landscapes DEL<sub>i</sub>/reference (ChEMBL)



#### How many DELs cover ChEMBL chemical space ?



# **Chemical Library Space**

- A GTM encodes a chemical library as a vector calculated from the property/activity landscape
- Ensemble of vectors can be used to build a meta-GTM ( $\mu GTM$ ) where each data point represents a library



#### µGTM of the DEL space



R. Pikalyova et al. J Chem Inf Model. 2023 63 (17), 5571-5582

#### µGTM: DEL reaction types



- Most of the DEL space is covered by coupling-based libraries
- Very few DELs that are purely heterocyclization-based
- The coupling-based DELs are more similar to ChEMBL than the heterocyclization ones.

# Cartography of ultra-large combinatorial libraries



# **Pipeline of GTM construction for combinatorial library**



#### **CoLiNN – Neural Network for GTM preparation without structure enumeration**



#### CoLiNN skips the enumeration step, thus accelerating GTM construction

Pikalyova R., T. Akhmetshin et al., et al. ChemRxiv, 2024, DOI: <u>10.26434/chemrxiv-2024-qh3bn</u>
#### **CoLiNN – Neural Network for GTM preparation without structure enumeration**

#### 1. Building Block Embedding

Building blocks:



#### 2. Reaction Embedding

**Reactions:** 



3. Responsibility vector prediction





#### Model performance as a function of training set size

Pikalyova R., T. Akhmetshin et al., et al. ChemRxiv (2024).



# **Drug resistance:** Cartography analysis of proteins and nucleic acids mutations

- Human Immunodeficiency Virus (HIV)
- Staphylococcus aureus bacteria



- 38 million people infected with HIV
- Highly effective therapy is available
- Increasing emergence of drug resistance





HIV infected patients



J.Y.Yeo, G.-R. Goh, C.T.-T. Su, S. K.-E. Gan, Viruses 2020, 12, 297.

HIV infected patients



J.Y.Yeo, G.-R. Goh, C.T.-T. Su, S. K.-E. Gan, Viruses 2020, 12, 297.

### HIV drug resistance: data



**4324** HIV protease (PR) and reverse transcriptase (RT) amino acid sequences and associated resistance profiles to 6 PR and 8 RT inhibitors

Sequence	Drug resistance profile			
Indinavir	Darunavir	Nelfinavir	Lopinavir	•••
P·N·I·W·K·T Resistant	Susceptible	Resistant	Susceptible	
PLITKT Resistant	Susceptible	Resistant	Resistant	
P Q I T K T Susceptible	Susceptible	Resistant	Susceptible	

HIV Drug Resistance Database. https://hivdb.stanford.edu/. Datasets updated on 03.02.2021.

#### Sequence descriptors



### HIV-1 protease drug resistance



- Homodimer
- 99 amino acid residues per monomer
- 8 approved protease inhibitors

NH<sub>2</sub>

## **Resistance landscapes for protease inhibitors**



#### **Resistance-determining mutation patterns in protease**







### Predicting drug resistance for emergent strains





### Antimicrobial resistance: Staphylococcus aureus



- Bacterial antimicrobial resistance (AMR) responsible for 1.27 million global deaths
- AMR could result in US\$ 1 trillion additional healthcare costs by 2050 (the World Bank estimates)
- Increasing emergence of antibiotic resistance

## Tackling "fat data" challenge in genomic data



### Gene-specific landscapes linking genes with drug resistance





# Chemography-guided generation of novel entities using AI tools

- •Molecules
- •Reactions
- •Peptides

#### Autoencoder performing SMILES reconstruction



#### AutoEncoder: sampling using a seed vector



**Goal**: to identify a seed vector from which valid structures possessing a given activity can be generated

### AutoEncoder: example of sampled structures



The numbers correspond to the Tanimoto similarity

#### Chemography-guided molecule generation



GTM identifies a zone in the latent space from which "useful" structures are sampled. Such zone is detected either by the human or by computer algorithm

B. Sattarov et al. J. Chem. Inf. Model., 2019, 59(3), 1182-1196

#### 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 chemical transformations



#### **Reactions: methodological problems:**

- Complexity: reaction equation contains several molecular graphs of two types reactants and products
- Reaction novelty detection and reaction feasibility assessment are not well established

#### Complexity reduction: Condensed Graph of Reaction (CGR)



#### CGR SMILES is almost twice shorter than conventional reaction SMILES

A. Varnek et al., J. Computer-Aided Molecular Design, 2005, 19, 693-703

## AI-driven design of new Suzuki-like reactions



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

#### AI-driven design of reactions: experimental validation



W. Bort et al., Nature Scientific Reports, 2021, 11, 3178

#### In silico design of novel peptides against Methicillin Resistant Staphylococcus Aureus (MRSA)



Diverse range of infections (from mild skin infections to life-threatening conditions, e.g. sepsis)

Duerden BI. Eye (Lond). 2012. 26.



People with MRSA are 64% more likely to die as opposed to its non-resistant variant

World Health Organization (2018). Antimicrobial resistance fact sheet.



Formation of biofilms (more resistant at conventional antibiotics dosages)

Choi, V. et al. Nat Rev Microbiol 21, 555-572 (2023).

### Design of antimicrobial peptides: generation workflow



### **Training set composition**



#### **Inclusion criteria:**

- 10-14 amino acids long
- Only natural amino acids
- Only linear peptides

#### Cartography-based design of novel peptides



Pikalyova K., et al., et al. BioRxiv, 2024, DOI: 10.1101/2024.11.17.622654

## **Experimental results**



THE UNIVERSITY OF BRITISH COLUMBIA

Synthesis & experimental testing of peptides against MRSA biofilms and planktonic cells



#### Screening results

8/8 peptides more active against MRSA biofilms compared to control



 $\circ$ 

5/8 peptides active against MRSA planktonic cells compared to control

Most active peptide is 10-fold more active against MRSA biofilms well-studied control peptide IDR-1018



## AutoEncoder vs Molecular descriptors space



GTM Class landscapes for A2a-receptors binders (1303 actives and 3618 inactives)

69

no structures generation •

- generation of new structures

**Goal**: development of deep-learning architecture able to generate structures with desired activities using *any* descriptor space

#### Sampling from any descriptor space: inverse-QSAR task



Attention-based Conditional Variational Autoencoder (ACoVAE) deep-learning architecture able to generate structures from any descriptor space



#### Neighborhood preservation benchmark: GTM vs t-SNE, UMAP and PCA





**Congeneric series** 103 subsets extracted from ChEMBL

**Descriptor sets** Morgan fingerprints, MACCS keys, ChemDist







**Comparative analysis** of neighborhood preservation according to 18 metrics


## Why GTM ?

GTM considers data probability distribution which makes it suitable for Big Data analysis

GTM manifold can be built on a small representative subset and it can accommodate new data



GTM landscapes can be used in regression or classification tasks

A data collection can be encoded by a vector which enables fast chemical libraries' comparison

Coupling of GTM with the AI technologies facilitates *de novo design* 

## Teams









- MSC-DN Marie-Curie AiChemist
- ITN Marie-Curie BigChem
- ITN Marie Curie TubInTrain
- Institute of Organic Chemistry, Kiev, Ukraine

- Eli Lilly
- SANOFI
- Enamine
- eMolecules
- Novalix
- Janssen Pharmaceutical
- TotalEnergy
- SOLVAY

## Cartography is an efficient way to explore a (chemical) space

