# De novo drug design

Pavel Polishchuk

Institute of Molecular and Translational Medicine
Faculty of Medicine and Dentistry
Palacky University

pavlo.polishchuk@upol.cz
qsar4u.com

real datasets

**SCIFINDER®**
A **CAS** SOLUTION

~ 160 M compounds

**REAXYS®**

~ 105 M compounds

Commercial

**PubChem**

~ 102 M compounds          Free

## ZINC

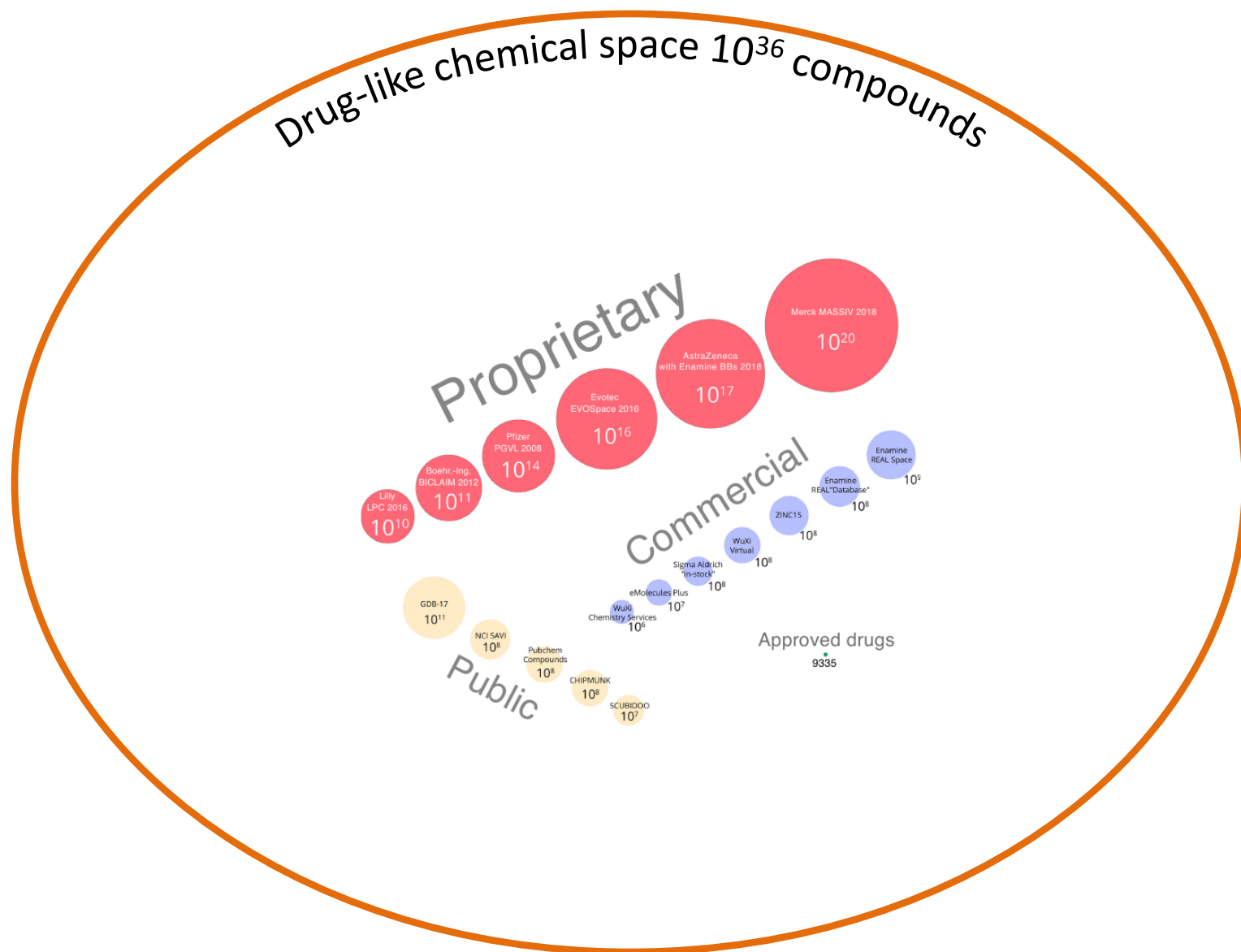up to 1 B commercially available compounds

virtually enumerated dataset

## GDB-17

166 B compounds = $1.66 \times 10^{11}$

Hoffmann, T.; Gastreich, M., The next level in chemical space navigation: going far beyond enumerable compound libraries. *Drug Discovery Today* **2019**, 24, 1148-1156. (https://doi.org/10.1016/j.drudis.2019.02.013)

Drug-like chemical space $10^{36}$ compounds

Proprietary

Commercial

Public

Merck MASSIV 2018
$10^{20}$

AstraZeneca with Enamine BBs 2018
$10^{17}$

Evotec EVOSpace 2016
$10^{16}$

Pfizer PGVL 2008
$10^{14}$

Boehr.-Ing. BICLAIM 2012
$10^{11}$

Lilly LPC 2016
$10^{10}$

Enamine REAL Space
$10^{9}$

Enamine REAL "Database"
$10^{8}$

ZINC15
$10^{8}$

WuXi Virtual
$10^{8}$

Sigma Aldrich "in-stock"
$10^{8}$

eMolecules Plus
$10^{7}$

WuXi Chemistry Services
$10^{6}$

GDB-17
$10^{11}$

NCI SAVI
$10^{8}$

Pubchem Compounds
$10^{8}$

CHIPMUNK
$10^{8}$

SCUBIDOO
$10^{7}$

Approved drugs
9335

# Virtual screening *vs.* de novo design

## Virtual screening


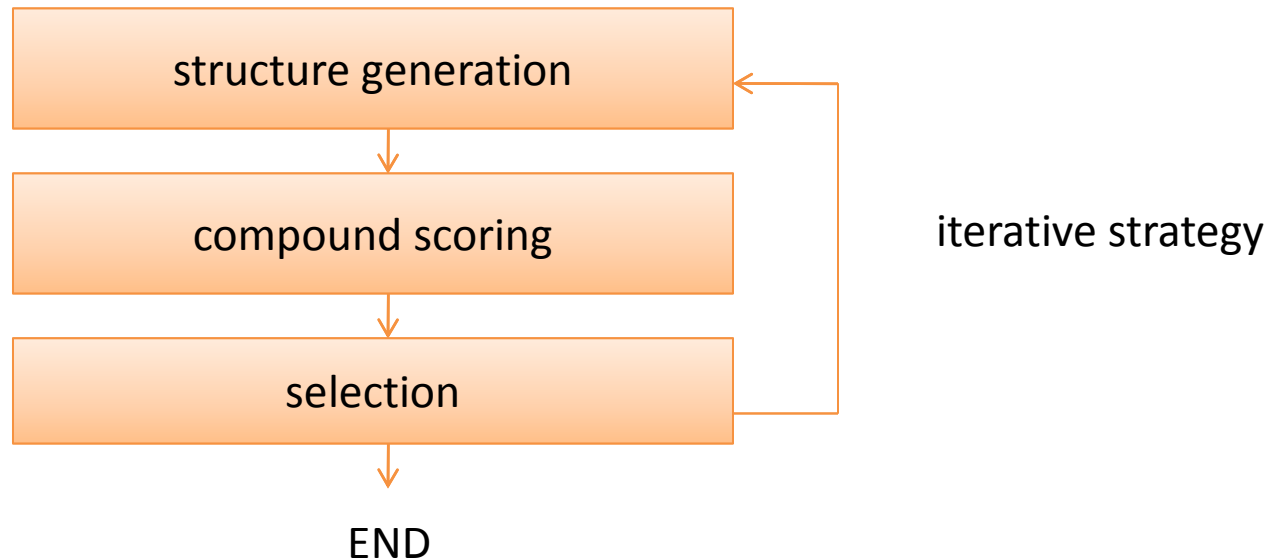
$10^9$-$10^{11}$ compounds

10-100 compounds

## De novo design



~$10^{36}$ drug-like compounds

10-100 compounds
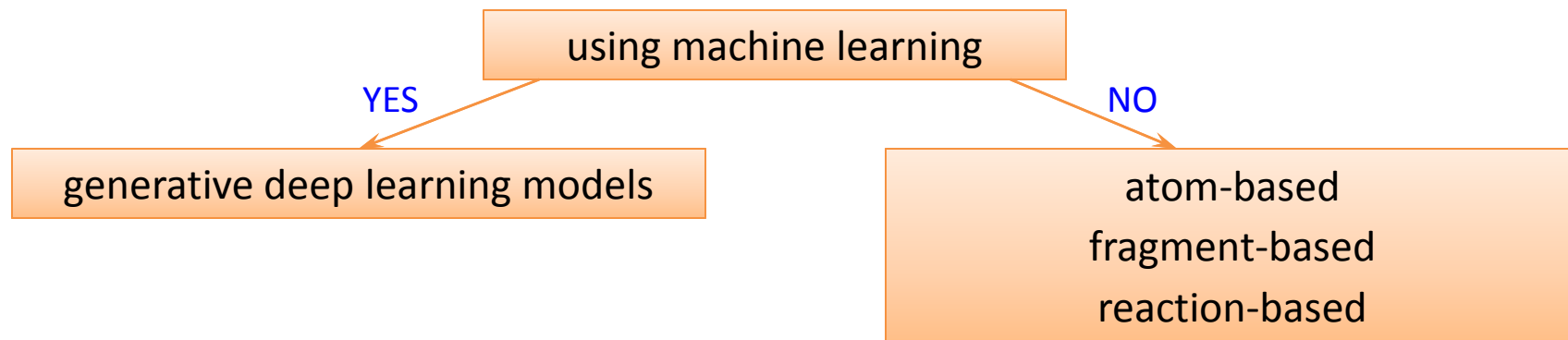
1. **Structure generation** - how to create/assembly new structures

2. **Compound scoring** - how to estimate/predict a property of a compound

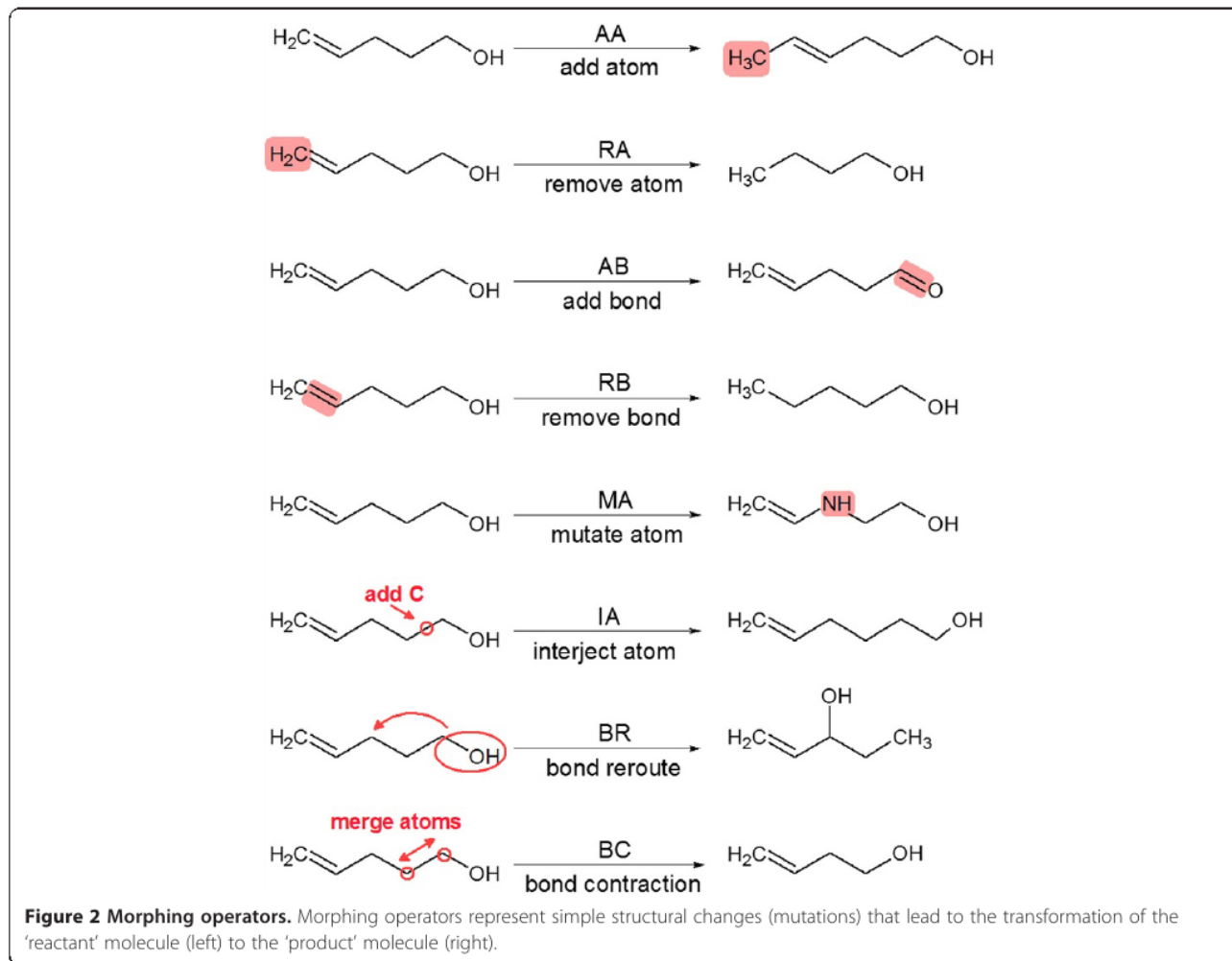3. **Search strategy** - how to find compounds with optimal properties

# De novo structure generation

using machine learning

YES

generative deep learning models

NO

atom-based
fragment-based
reaction-based

- **atom-based** - uses simple rules like add/change/remove atom/bond to perturb structures
- **fragment-based** - uses fragment library to create structures
- **reaction-based** - uses a set of reaction rules and a library of reactants

# Atom-based structure generation

Molpher



**Figure 2 Morphing operators.** Morphing operators represent simple structural changes (mutations) that lead to the transformation of the 'reactant' molecule (left) to the 'product' molecule (right).

**Figure 9 An example of the path between pentamidine (CID 4735) and 2-imino-3-(1H-indol-3-yl)propanoic acid (CID 5599) from dataset D3.** The path was generated using Morgan fingerprint, Tanimoto coefficient, and synthetic accessibility filter turned on. The arrows' labels show the used morphing operators (see Figure 2). The depiction of the path was done by OpenBabel [74].
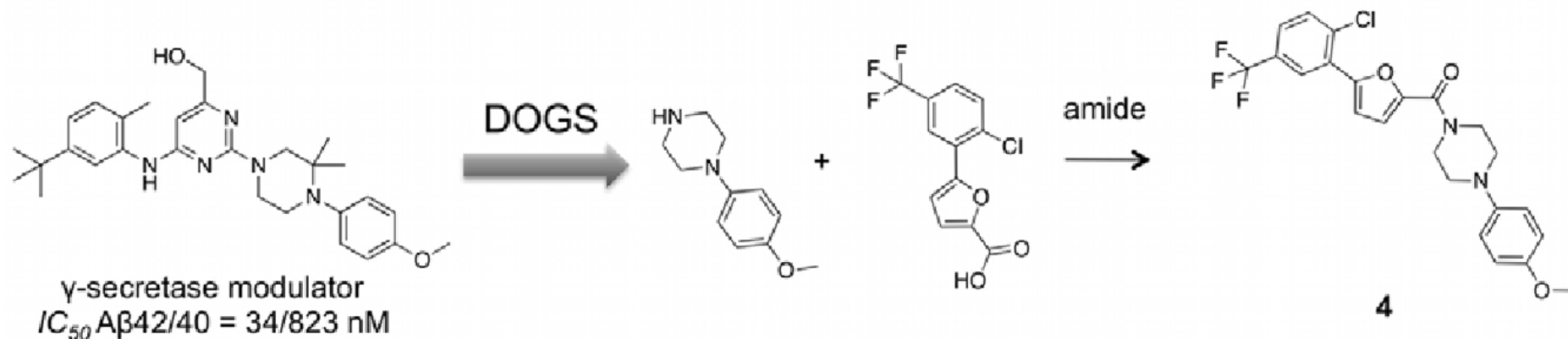
22 steps

# Atom-based structure generation

| parameters | atom-based |
|---|---|
| exhaustiveness  of chemical space search | ++++<br>very small steps;<br>more suitable for systematic exploration of local chemical space |
| structure novelty | +++* |
| structure diversity | +++* |
| chemically valid structures | - |
| synthetically feasible | --- |
| combinatorial explosion  / time consuming | --- |

atom-based ≈ *ab initio*

DOGS

# Reaction-based structure generation

DOGS

γ-secretase modulators



inverse modulator

Hartenfeller, M.; Zettl, H.; Walter, M.; Rupp, M.; Reisen, F.; Proschak, E.; Weggen, S.; Stark, H.; Schneider, G., DOGS: Reaction-Driven de novo Design of Bioactive Compounds. *PLOS Computational Biology* **2012**, 8, e1002380.

# Reaction-based structure generation

## Retinoid X Receptor(RXR) Modulators



Supporting figure 5: Synthesis of de novo mimetics **1a** and **3**-**8**. Reagents and conditions: (a) EtOH, HOAc, *μ*w, 100°C, 3-6 h, 43-78%; (b) montmorillonite K10, *μ*w, 90°C, 30 min, 41-85%.

Merk D., et al. *J. Med. Chem.,* **2018**, 61 (12), pp 5442–5447

| | reaction-based |
|---|---|
| exhaustiveness of chemical space search | +<br>depends on reactant library and reaction rules;<br>only grow molecules |
| structure novelty | ++ |
| structure diversity | ++ |
| chemically valid structures | +++ |
| synthetically feasible | +++ |
| combinatorial explosion / time consuming | +++ |

reaction-based ≈ empirical

# Fragment-based structure generation

GROW

MUTATE

LINK

REDUCE

BREED

Superimpose ligands

Bond match

Split & recombine ligands

BREED: HIV-1 protease inhibitors



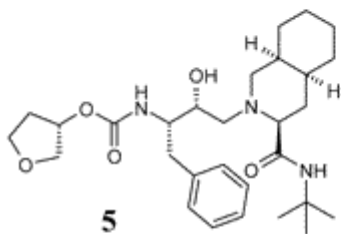$K_i$ = 0.4-0.6 nM

$K_d$ = 1.1 nM

$K_i$ = 1.7 nM

$K_d$ = 0.3 nM

known

designed

$IC_{50}$ = 160 nM

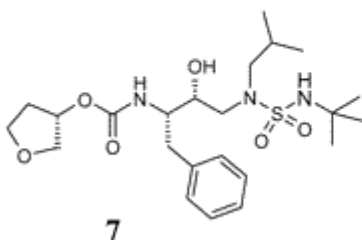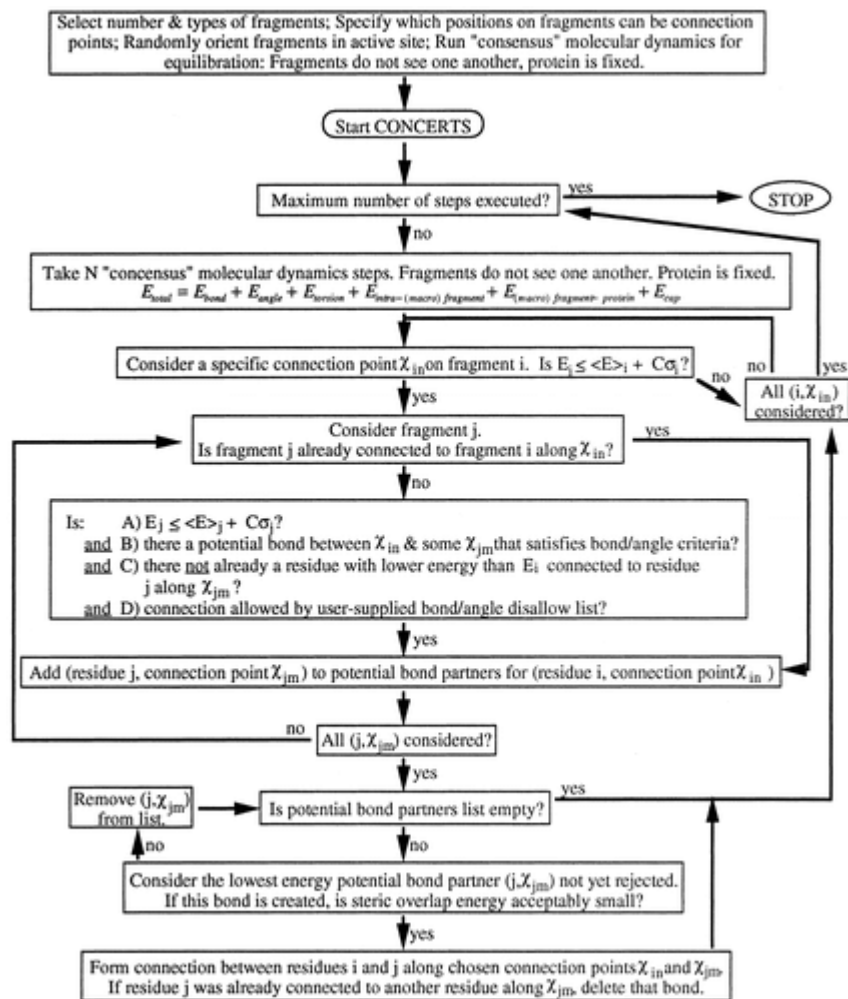$K_i$ = 0.1 nM

$K_i$ = 42 nM

Pierce A.C., Rao G., Bemis G.W. *J. Med. Chem.*, **2004**, 47 (11), pp 2768–2775
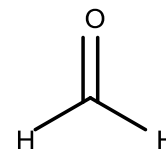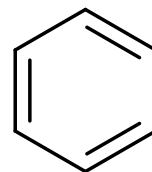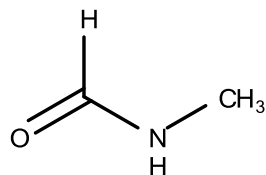
**CONCEPTS**



MD of fragments which are linking or breaking during the simulation in order to create more favorable structures

formation of certain bonds was forbidden: O–O, N–N, N–O, S–O, O–C–O, O–N–O, N–C–N, $C_\alpha$–$C_\alpha$, C–$C_\alpha$–C
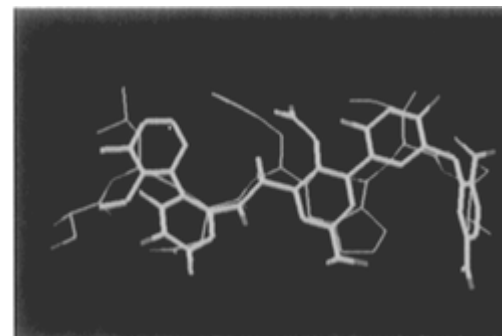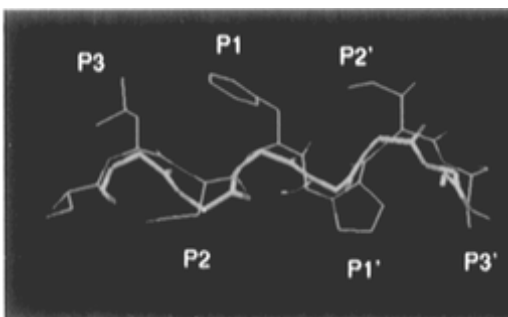
Pearlman D.A., Murcko M.A. *J. Med. Chem.*, **1996**, 39 (8), pp 1651–1663

## CONCEPTS: HIV-1 protease inhibitors



+ 19 side chains

Pearlman D.A., Murcko M.A. *J. Med. Chem.*, **1996**, 39 (8), pp 1651–1663

# Fragment-based structure generation

FOG



Kutchukian, P. S.; Lou, D.; Shakhnovich, E. I., FOG: Fragment Optimized Growth Algorithm for the de Novo Generation of Molecules Occupying Druglike Chemical Space. *Journal of Chemical Information and Modeling* **2009**, 49, 1630-1642.

# Fragment-based structure generation

eMolFrag



Liu, T.; Naderi, M.; Alvin, C.; Mukhopadhyay, S.; Brylinski, M., Break Down in Order To Build Up: Decomposing Small Molecules for Fragment-Based Drug Design with eMolFrag. *J. Chem. Inf. Model.* **2017**, 57, 627-631

# Fragment-based structure generation

CReM: chemically reasonable mutations

exhaustive fragmentation
cutting single bonds

taking context of radius R



DB of replacements

| environment (radius = 3) | fragments |
|---|---|
| 🔶 🔴 | ⬡ 🔺 🟨 ... |
| ... | ... |

mutually exchangeable fragments

# Fragment-based structure generation

CReM: chemically reasonable mutations
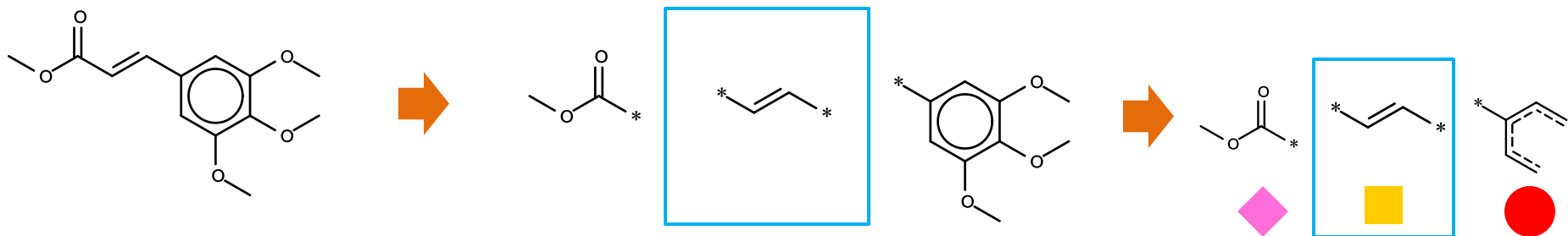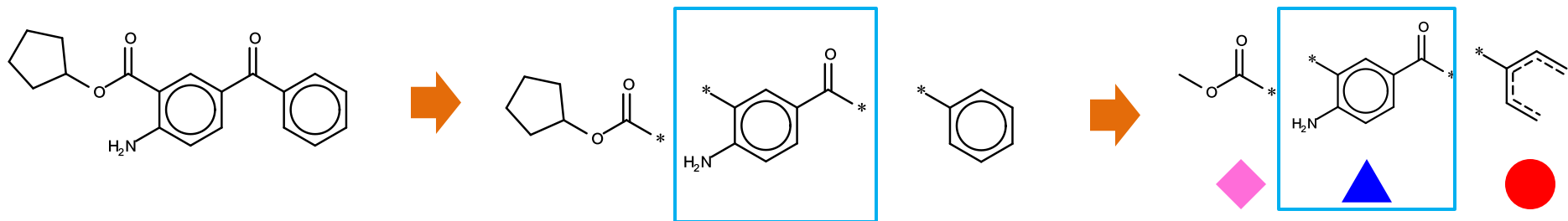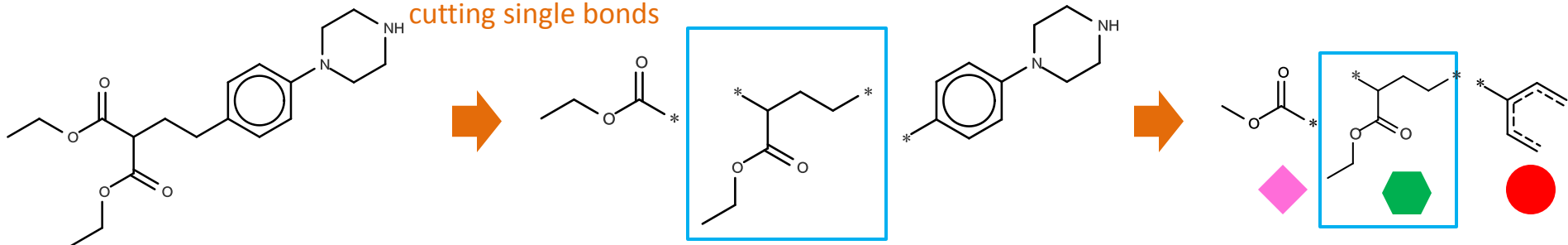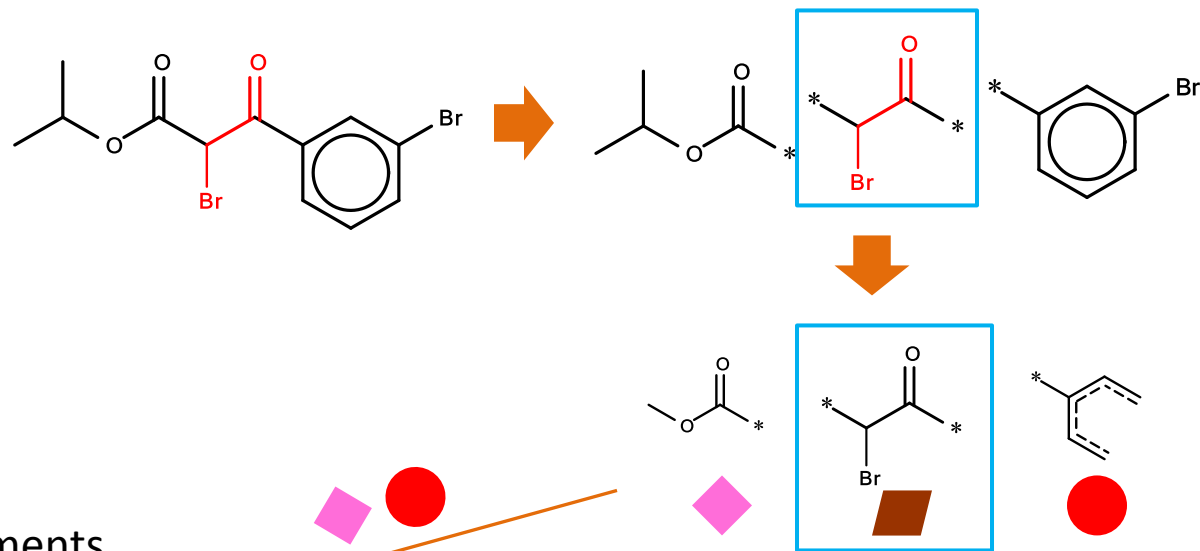


Generated structures are always chemically valid!

# Fragment-based structure generation

| | fragment-based |
|---|---|
| exhaustiveness of chemical space search | ++ <br> variable, controlled by the size of fragments to replace |
| structure novelty | ++ |
| structure diversity | ++ |
| chemically valid structures | (+++) |
| synthetically feasible | (++) |
| combinatorial explosion / time consuming | ++ |

fragment-based ≈ semi-empirical

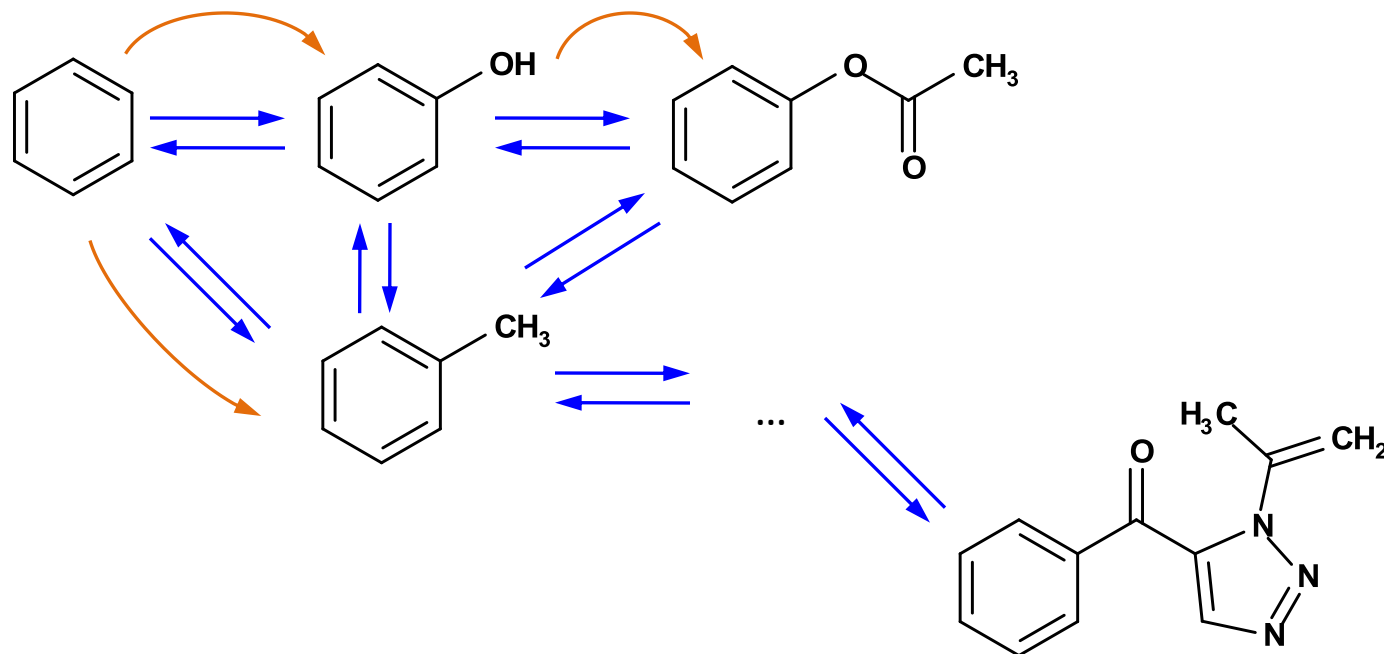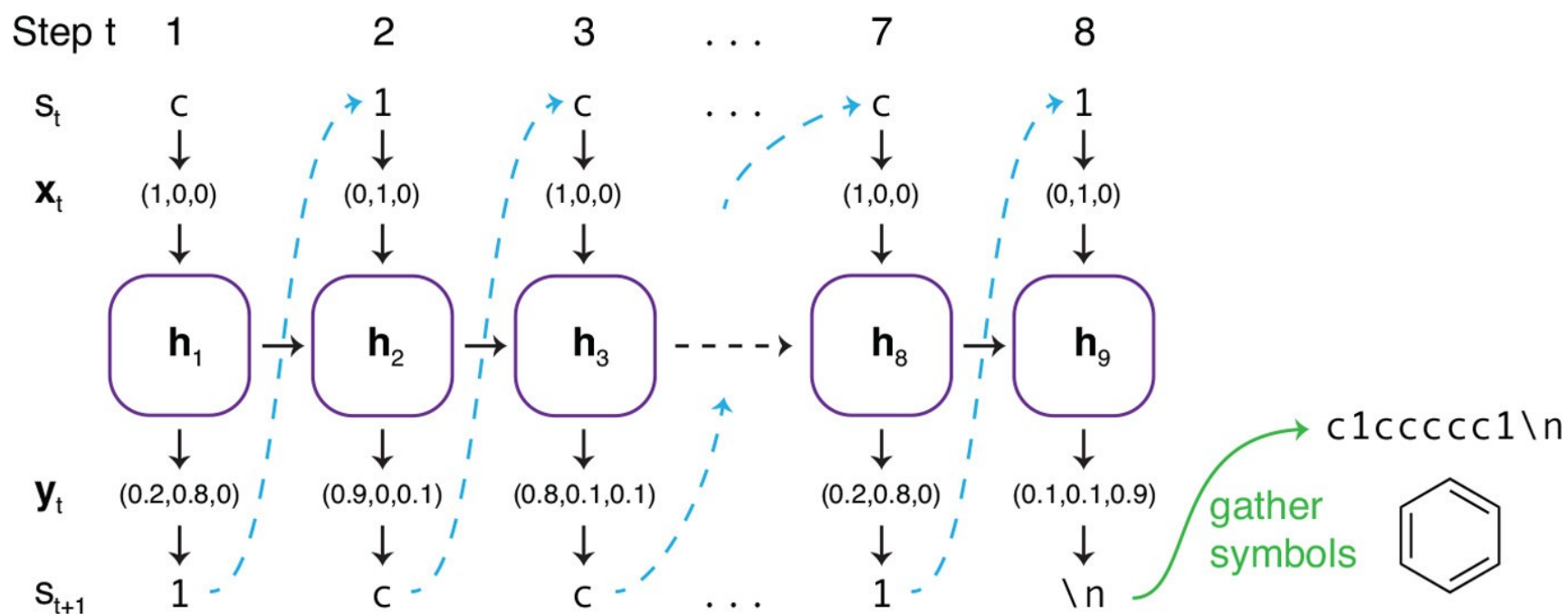| | Reaction-based | Fragment-based |
|---|---|---|
| Prerequisites: | reaction rules set<br>database of building blocks | database of fragments |
| Abilities & issues: | • molecules are more likely to be feasible<br>• not all moves are allowed<br>• usually only increase complexity<br>• some molecules can be unreachable | • do not control synthetic feasibility<br>• many moves are allowed<br>• arbitrary direction of exploration<br>• cover larger chemical space |

# De novo structure generation

Summary

| | atom-based | fragment-based | reaction-based |
|---|---|---|---|
| exhaustiveness of chemical space search | ++++ <br> very small steps; more suitable for systematic exploration of local chemical space | ++ <br> variable, controlled by the size of fragments to replace | + <br> depends on reactant library and reaction rules; only grow molecules |
| structure novelty | +++* | ++ | ++ |
| structure diversity | +++* | ++ | ++ |
| chemically valid structures | - | (+++) | +++ |
| synthetically feasible | --- | (++) | +++ |
| combinatorial explosion / time consuming | --- | ++ | +++ |

## Recurrent neural network (RNN)



Segler, M. H. S.; Kogej, T.; Tyrchan, C.; Waller, M. P., Generating Focused Molecule Libraries for Drug Discovery with Recurrent Neural Networks. *ACS Central Science* **2018**, 4, 120-131.

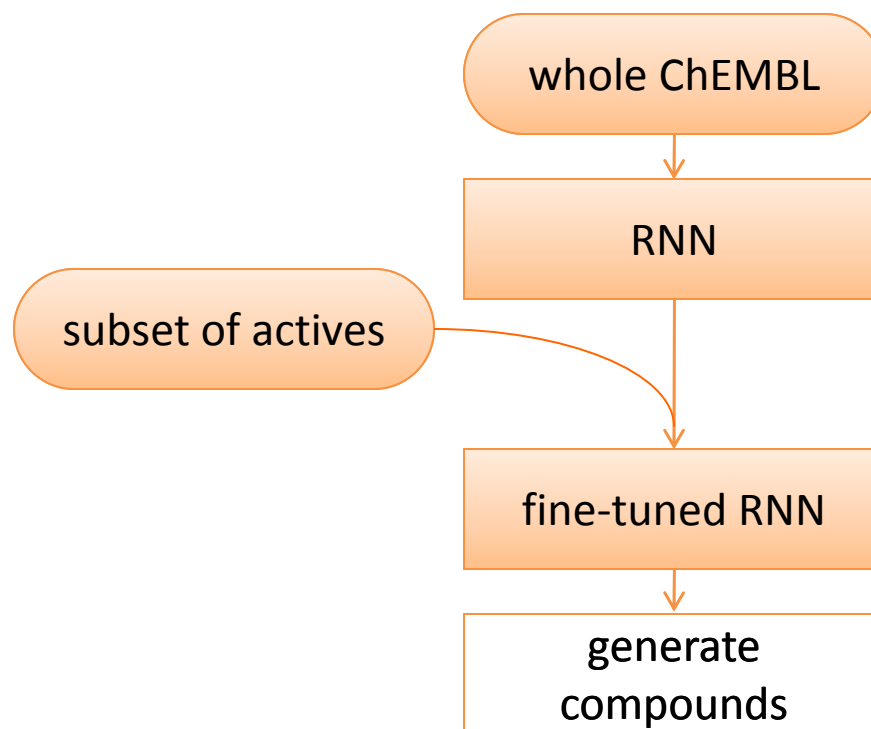# Deep learning models for structure generation

unsupervised generation

transfer learning

Segler, M. H. S.; Kogej, T.; Tyrchan, C.; Waller, M. P., Generating Focused Molecule Libraries for Drug Discovery with Recurrent Neural Networks. *ACS Central Science* **2018**, 4, 120-131.

# Deep learning models for structure generation

unsupervised generation



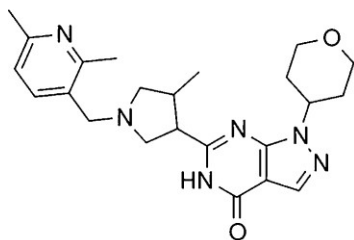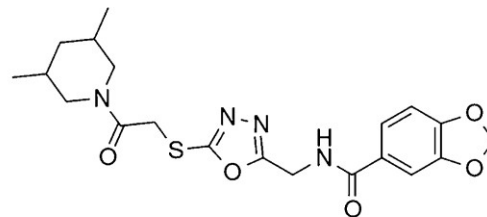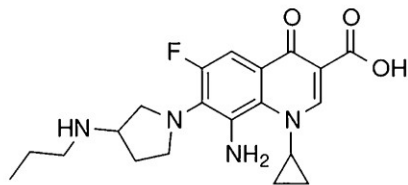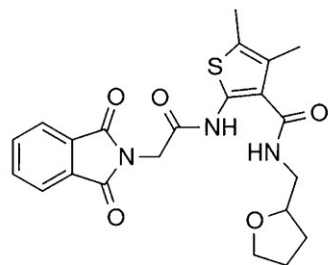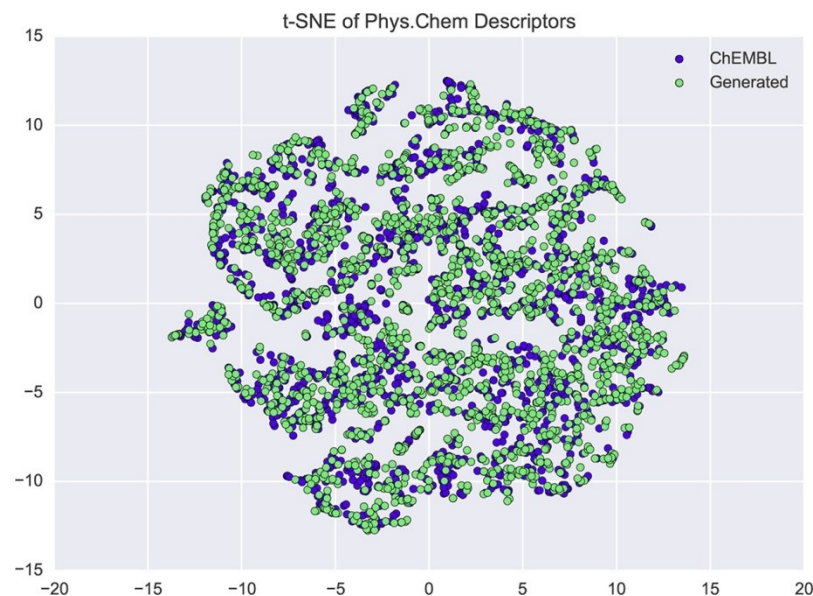976 327 compounds
97.7% chemically valid
11.5% were duplicated with ChEMBL
1.7% of duplicates

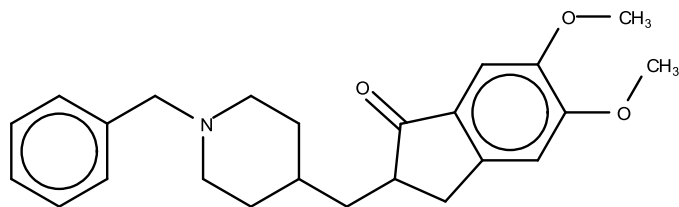75% passed AZ filters (similar to ChEMBL)
12% of scaffolds were common with ChEMBL

Segler, M. H. S.; Kogej, T.; Tyrchan, C.; Waller, M. P., Generating Focused Molecule Libraries for Drug Discovery with Recurrent Neural Networks. *ACS Central Science* **2018**, 4, 120-131.

# Deep learning models for structure generation

| | deep learning |
|---|:---:|
| exhaustiveness of chemical space search | ++ |
| structure novelty | ++ |
| structure diversity | ++ |
| chemically valid structures | ++ |
| synthetically feasible | ? |
| combinatorial explosion / time consuming | +++ |

Issue of SMILES based representation -
the same structure can be represented by different SMILES



COc1cc2CC(CC3CCN(Cc4ccccc4)CC3)C(=O)c2cc1OC
COc1cc2c(cc1OC)C(=O)C(CC1CCN(Cc3ccccc3)CC1)C2
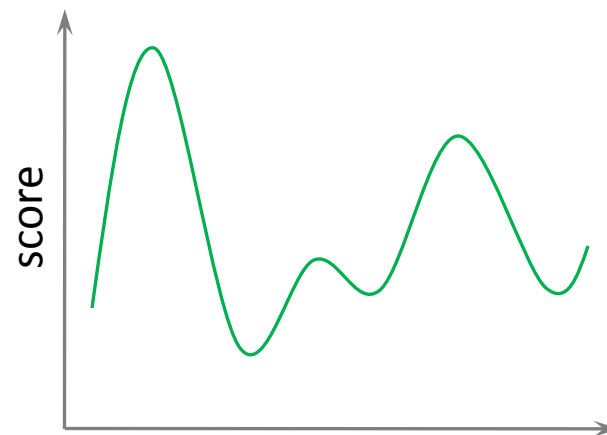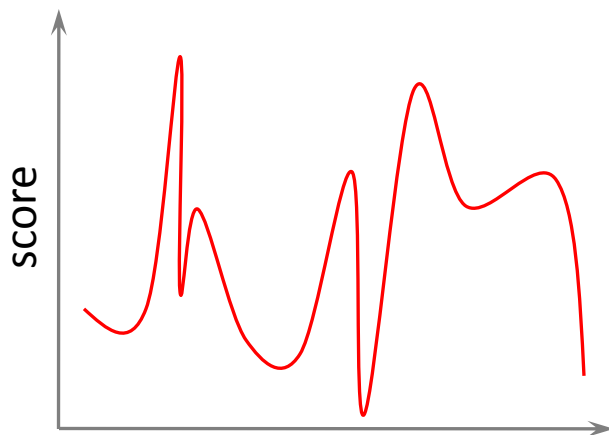
Can be any but preferably smooth to follow the chemical similarity principle:

- physicochemical properties

- similarity measures      } ligand-based scoring functions

- QSAR model prediction

- pharmacophore fit

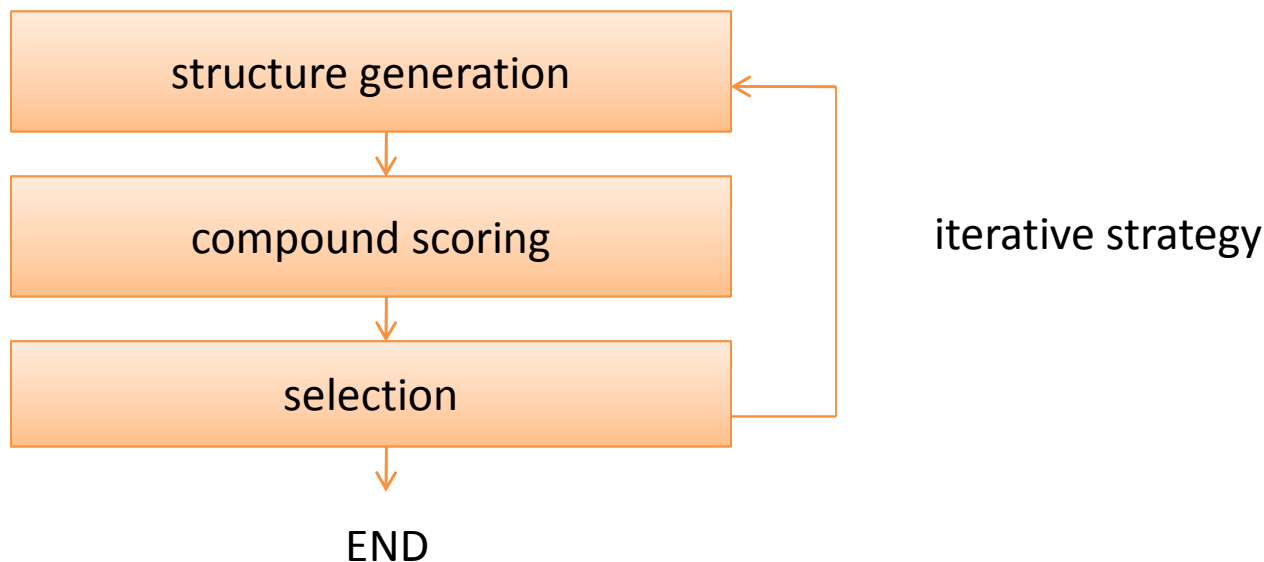- docking scoring      } structure-based scoring functions

- molecular dynamics

…

Can be any , for example:

- greedy search

- Monte Carlo

- evolutionary algorithms, e.g.:

    - genetic algorithm
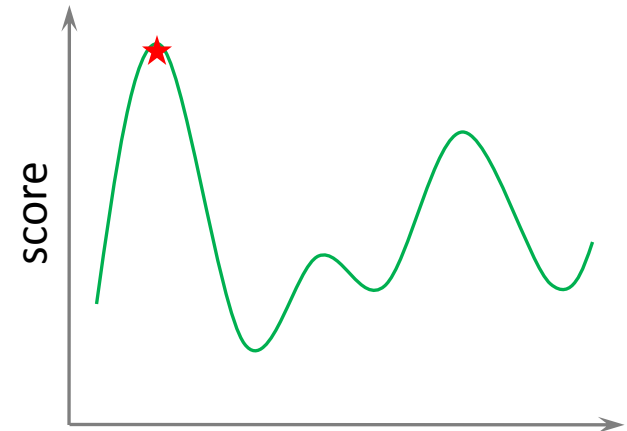
- simulated annealing

…

1. **Structure generation** - how to create/assembly new structures

2. **Compound scoring** - how to estimate/predict a property of a compound

3. **Search strategy** - how to find compounds with optimal properties

| $D_1$ | $D_2$ | $D_3$ | ... | $D_N$ |
|---|---|---|---|---|
| 1 | 0 | 9 | ... | 1 |
| 4 | 0 | 1 | ... | 1 |
| 0 | 2 | 3 | ... | 3 |
| ... | ... | ... | ... | ... |
| 4 | 0 | 0 | ... | 1 |

score

**STRUCTURE ?**

| $D_1$ | $D_2$ | $D_3$ | ... | $D_N$ |
|---|---|---|---|---|
| 11 | 3 | 1 | ... | 15 |

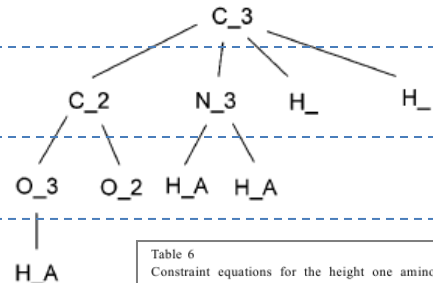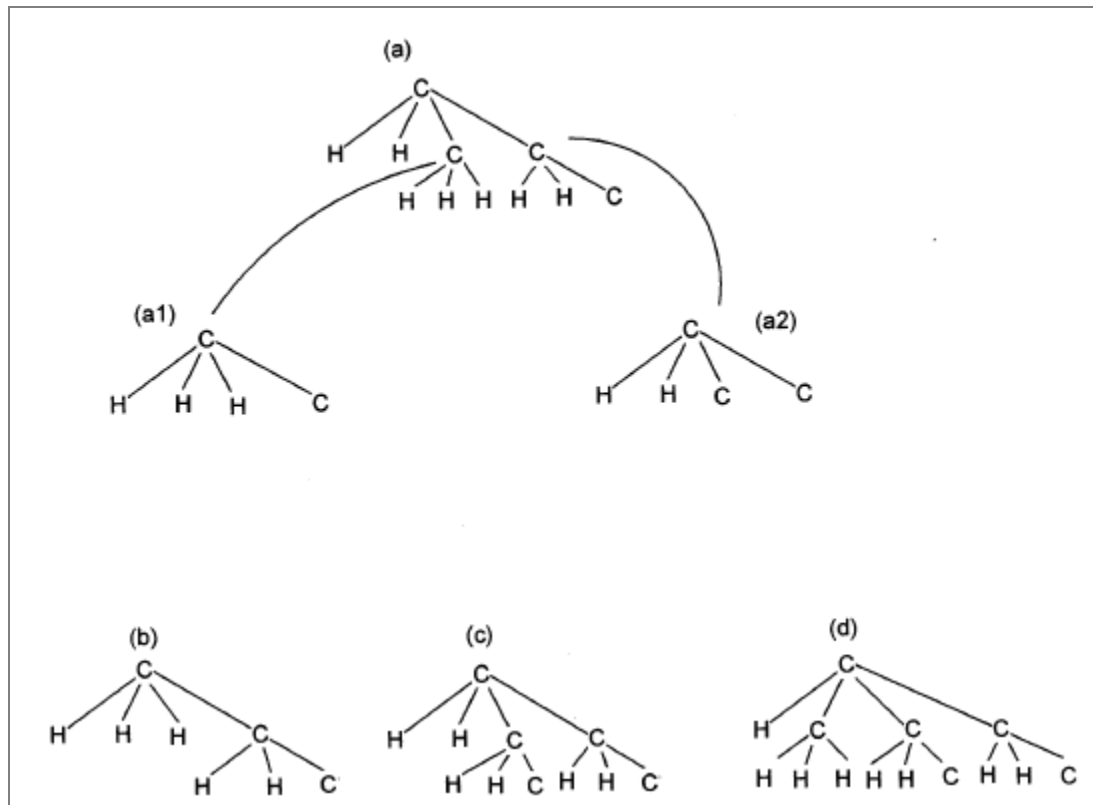## Atom signatures



$\sigma^0$

$\sigma^1$

$\sigma^2$

$\sigma^3$

Table 6
Constraint equations for the height one amino acid signatures in the training set
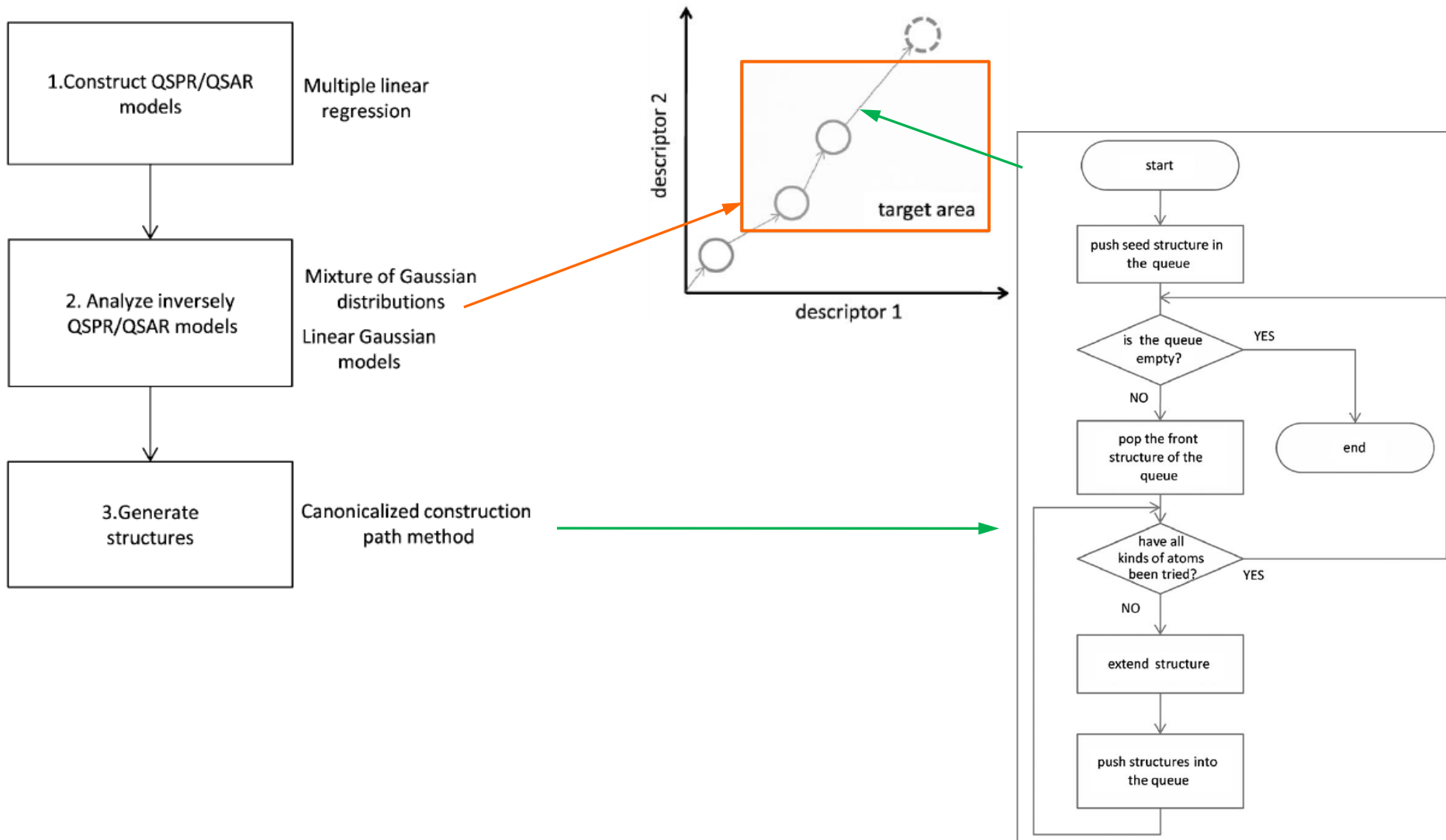
(1) $-x_{44} + x_{46} = 0$
(2) $-x_{38} + x_{47} = 0$
(3) $-x_{22} - x_{27} + x_{45} + x_{47} = 0$
(4) $-x_{10} + x_{45} + x_{46} = 0$
(5) $-x_{34} - x_{37} + x_{41} + x_{42} + x_{43} + x_{44} = 0$
(6) $-x_{21} + x_{43} = 0$
(7) $-x_{16} + x_{40} = 0$
(8) $-x_{13} + x_{39} + x_{42} = 0$
(9) $-x_2 - x_5 + x_{39} + x_{40} + x_{41} = 0$
(10) $-x_{28} - x_{30} - 2x_{31} + x_{33} + x_{35} + x_{36} + x_{37} + x_{38} = 0$
(11) $-x_{18} - x_{24} - x_{26} - x_{27} + x_{32} + x_{36} = 0$
(12) $-x_{14} + x_{35} = 0$
(13) $-x_3 - x_4 - 2x_6 + x_{32} + x_{33} + x_{34} = 0$
(14) $-x_{15} - x_{16} + 2x_{29} + x_{30} = 0$
(15) $-x_5 + x_{28} = 0$
(16) $(x_{20} + x_{25} + x_{26})\%2 = 0$
(17) $-x_{15} + x_{23} + x_{25} = 0$
(18) $-x_{12} - x_{14} + x_{19} + x_{23} + x_{24} = 0$
(19) $-x_9 + x_{17} + x_{19} + x_{20} + x_{21} + x_{22} = 0$
(20) $-x_1 - x_4 + x_{17} + x_{18} = 0$
(21) $-x_8 + x_{11} + x_{12} + x_{13} = 0$
(22) $-x_3 + x_{11} = 0$
(23) $(x_7 + x_8 + x_9 + x_{10})\%2 = 0$
(24) $-x_1 - x_2 + x_7 = 0$

Eqs. (16) and (23) are modulus equations, which can be expressed as homogeneous equations by adding a dummy variable. For example Eq. (16) would read $x_{20} + x_{25} + x_{26} - 2z_1 = 0$. The % sign indicates the modulus is to be used.

Faulon J-L, Churchwell CJ, Visco DP, The Signature Molecular Descriptor. 2. Enumerating Molecules from Their Extended Valence Sequences – J. Chem. Inf. Comput. Sci., **2003**, 43 (3), pp 721–734
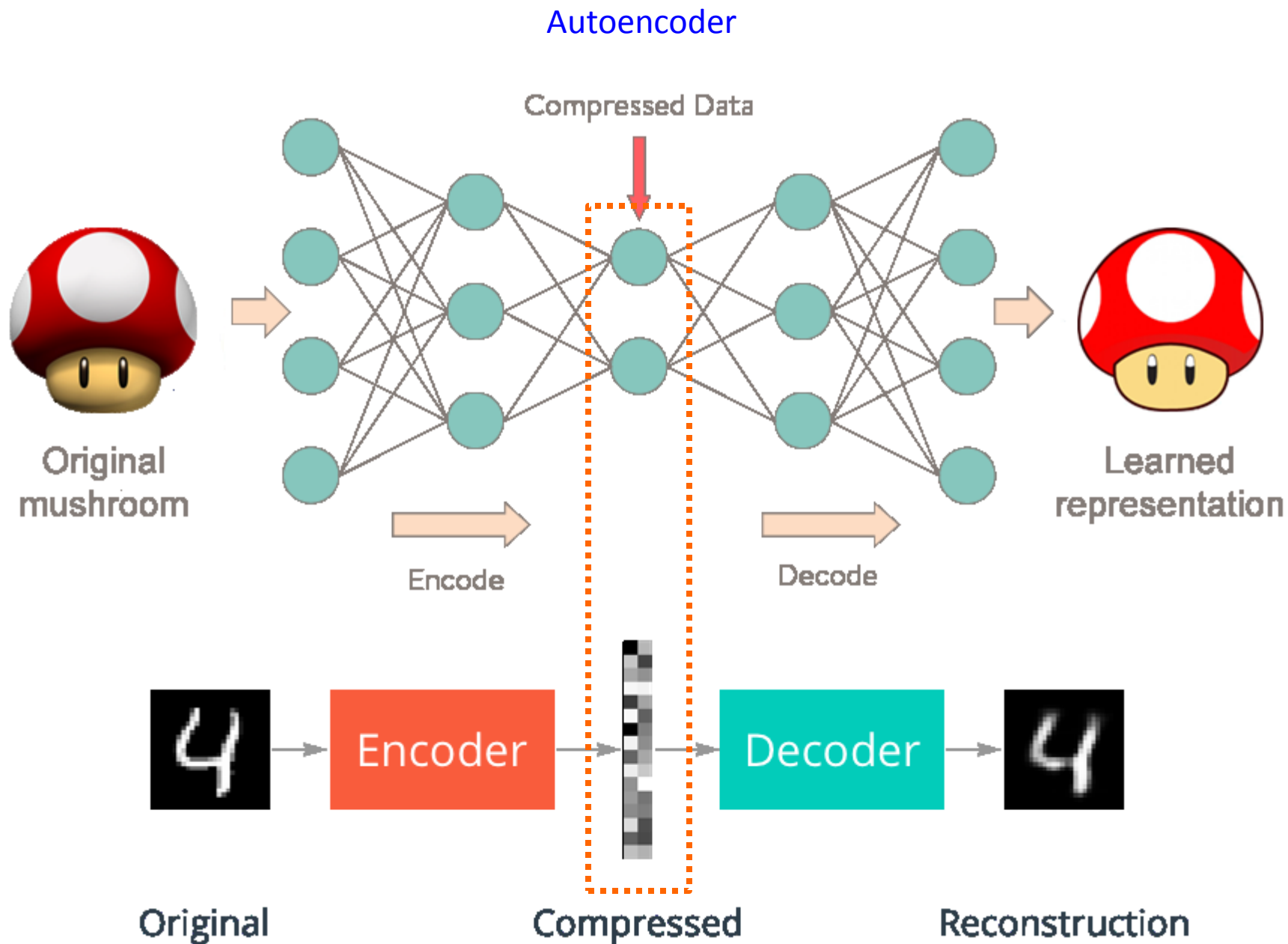
## Inverse QSAR with monotonically changed descriptors



Miyao, T.; Arakawa, M.; Funatsu, K., Exhaustive Structure Generation for Inverse-QSPR/QSAR.
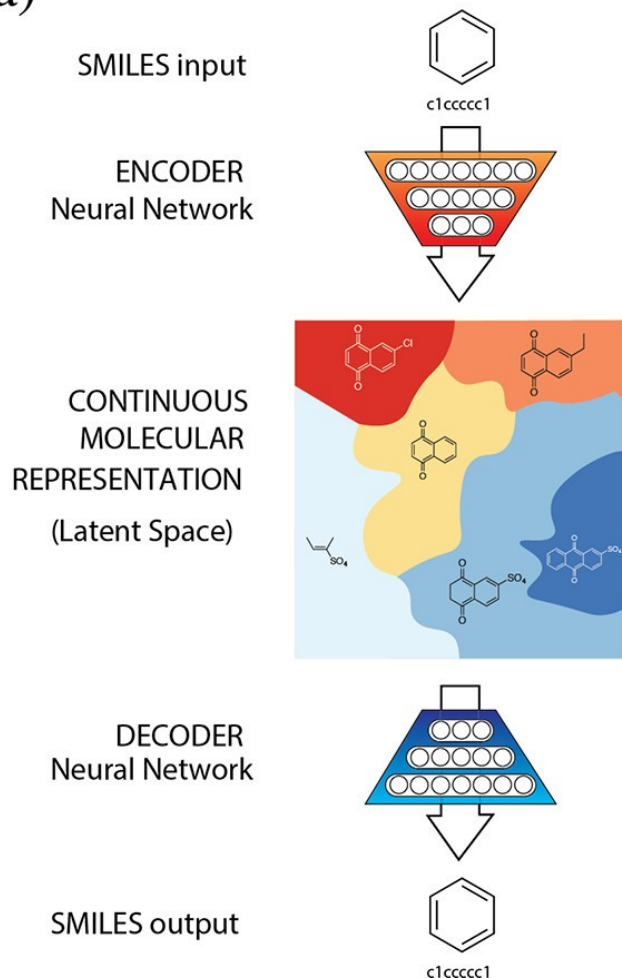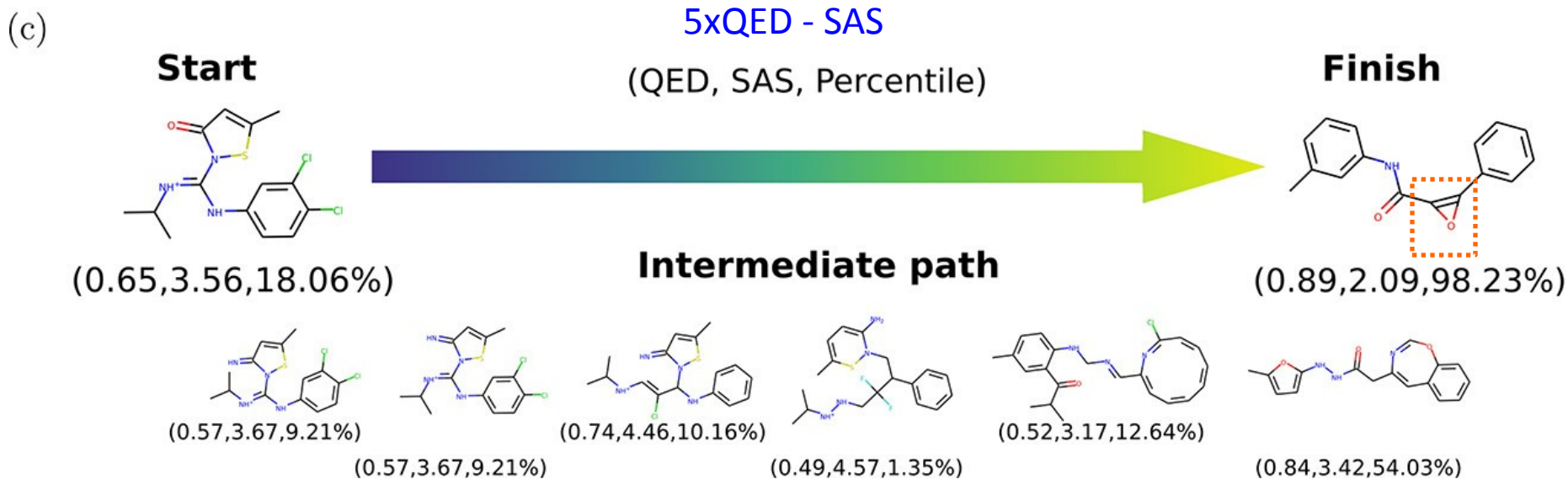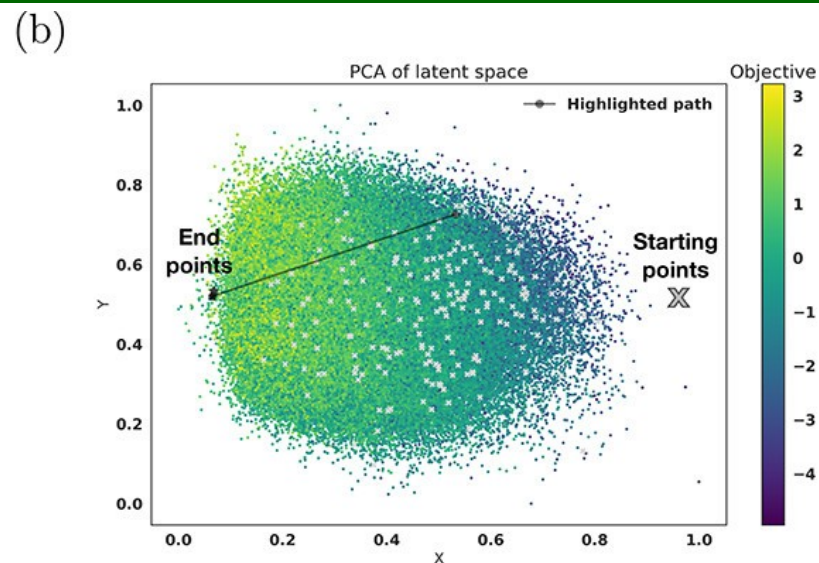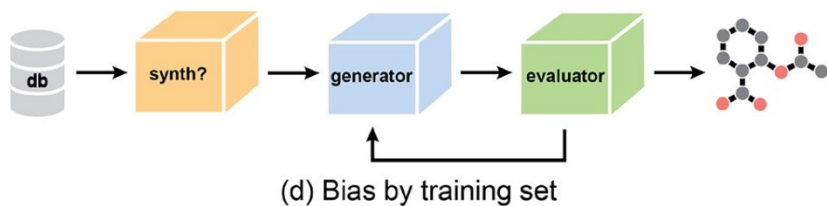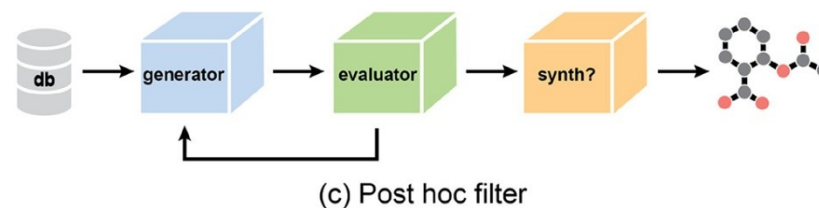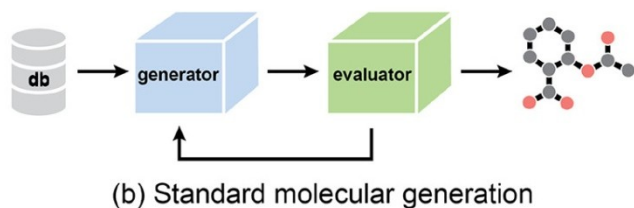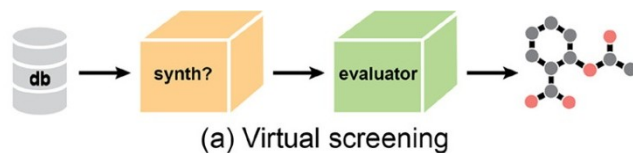*Molecular Informatics* **2010**, 29, 111-125.

## Autoencoder

Gómez-Bombarelli, R.; Wei, J. N.; Duvenaud, D.; Hernández-Lobato, J. M.; Sánchez-Lengeling, B.; Sheberla, D.; Aguilera-Iparraguirre, J.; Hirzel, T. D.; Adams, R. P.; Aspuru-Guzik, A., Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules. *ACS Central Science* **2018**, 4, 268-276.

(a)

(b)

(c) 5xQED - SAS

**Start** (QED, SAS, Percentile) **Finish**

(0.65,3.56,18.06%)

**Intermediate path**

(0.89,2.09,98.23%)

(0.57,3.67,9.21%)     (0.74,4.46,10.16%)     (0.52,3.17,12.64%)

(0.57,3.67,9.21%)     (0.49,4.57,1.35%)     (0.84,3.42,54.03%)

Gómez-Bombarelli, R.; Wei, J. N.; Duvenaud, D.; Hernández-Lobato, J. M.; Sánchez-Lengeling, B.; Sheberla, D.; Aguilera-Iparraguirre, J.; Hirzel, T. D.; Adams, R. P.; Aspuru-Guzik, A., Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules. *ACS Central Science* **2018**, 4, 268-276.

(a) Virtual screening

(b) Standard molecular generation

(c) Post hoc filter

(d) Bias by training set

(e) Bias by heuristics

(g) Explicit constraints

Gao, W.; Coley, C. W., The Synthesizability of Molecules Proposed by Generative Models. *Journal of Chemical Information and Modeling* **2020**

# Assessment of synthetic feasibility

Journal of Cheminformatics

**SOFTWARE**                                    **Open Access**

# AiZynthFinder: a fast, robust and flexible open-source software for retrosynthetic planning

Samuel Genheden[1*], Amol Thakkar[1,2], Veronika Chadimová[1], Jean-Louis Reymond[2], Ola Engkvist[1] and Esben Bjerrum[1*]

## Chemical Science

ROYAL SOCIETY OF CHEMISTRY

**EDGE ARTICLE**                                **View Article Online**
                                                View Journal | View Issue

Check for updates

## Retrosynthetic accessibility score (RAscore) — rapid machine learned synthesizability classification from AI driven retrosynthetic planning†

Amol Thakkar, *[ab] Veronika Chadimová, [a] Esben Jannik Bjerrum, [a] Ola Engkvist [a] and Jean-Louis Reymond [*b]

Journal of Cheminformatics

**RESEARCH ARTICLE**                            **Open Access**

# SYBA: Bayesian estimation of synthetic accessibility of organic compounds

Milan Voršilák[1,2], Michal Kolář[3,4], Ivan Čmelo[1] and Daniel Svozil[1,2*]
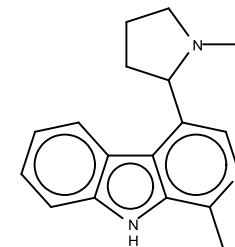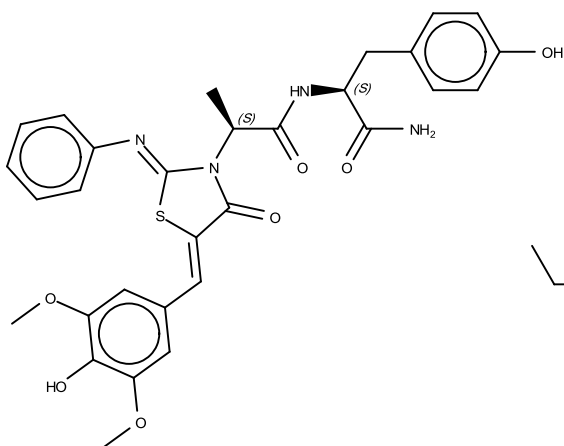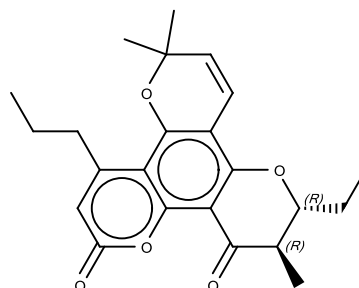
1.2
CHEMBL618

1.5
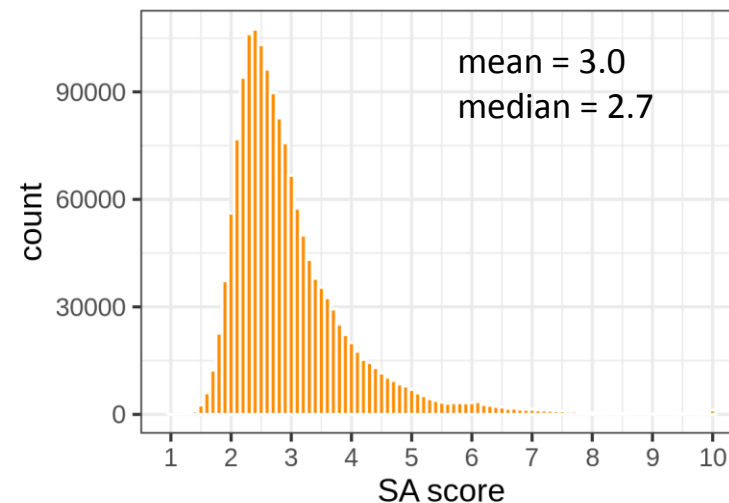CHEMBL3310985

2.0
CHEMBL595820

2.5
CHEMBL503660

3.0
CHEMBL500286

3.5
CHEMBL582554

4.0
CHEMBL7633

mean = 3.0
median = 2.7

Ertl, P.; Schuffenhauer, A., Estimation of synthetic accessibility score of drug-like molecules based on molecular complexity and fragment contributions. *Journal of Cheminformatics* **2009**, 1, 8. (10.1186/1758-2946-1-8)

## Content of fragmented library

all ChEMBL compounds (1 554 160)

compounds with SA score ≤ 2.5 (572 527)

compounds with SA score ≤ 2 (107 806)

## Context radius

1
2
3   less conservative replacements
4
5   more conservative replacements

(d) Bias by training set

(g) Explicit constraints

**CReM DB**
- ● all
- ◇ SA2
- ⬦ SA2.5

**radius**
- ● 1
- ● 2
- ● 3
- ● 4
- ● 5

**other approaches**
- ○ best from ChEMBL
- □ Graph GA
- △ SMILES GA
- ▽ SMILES LSTM

**bias type**
- ● SA bias
- ● no bias

a — SAScore vs total guacamol score for 10 hard tasks

Fig. 1 | V-SYNTHES approach to modular screening of Enamine REAL Space. A general overview of the four-step algorithm (left) and examples for each step (right). Asterisks in step one show the attachment points of synthons; arrows show possible pairing of minimal synthons with real synthons.
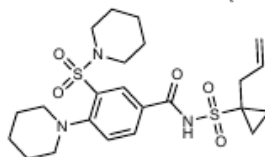
Sadybekov, A. A. et al, Synthon-based ligand discovery in virtual libraries of over 11 billion compounds. *Nature* **2021**. (10.1038/s41586-021-04220-9)

- De novo design can efficiently explore much larger chemical space  than virtual screening

- There are multiple  approaches to generate chemically valid structures, all of them have their pros and cons

- The main issue of de novo design is synthetic feasibility of  generated compounds

- There are several ways how to control synthetic feasibility

# Thank you for your attention