



10th RDKit UGM

14-15 October 2021

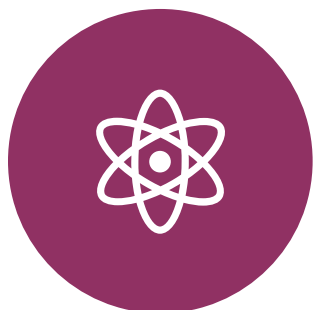
15<sup>th</sup> October 2021

# **BROWSER BASED EXPLORATION OF THE GDB CHEMICAL SPACE USING AI PLANNED SYNTHESIS FACILITATES EXPERIMENTAL ENGAGEMENT**

AMOL THAKKAR

# INTRODUCTION - OVERVIEW

*Currently individual components and we want to combine them – lower barrier to entry*



GDB CHEMICAL  
SPACE



AI PLANNED  
SYNTHESIS



BROWSER BASED  
EXPLORATION



EXPERIMENTAL  
ENGAGEMENT

*Enabling experimentation by linking chemical library visualisation to predicted synthesis in one tool*

# GDB CHEMICAL SPACE – HOW MANY MOLECULES ARE POSSIBLE?

Target  
Identification

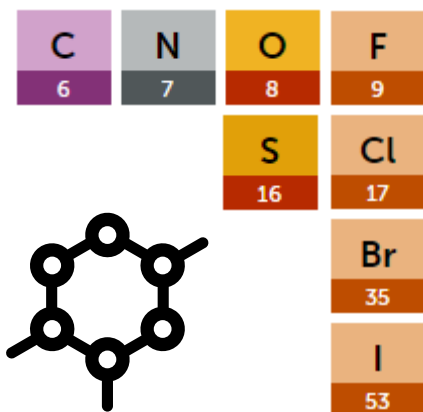
Enumeration

Filtering

Visualisation

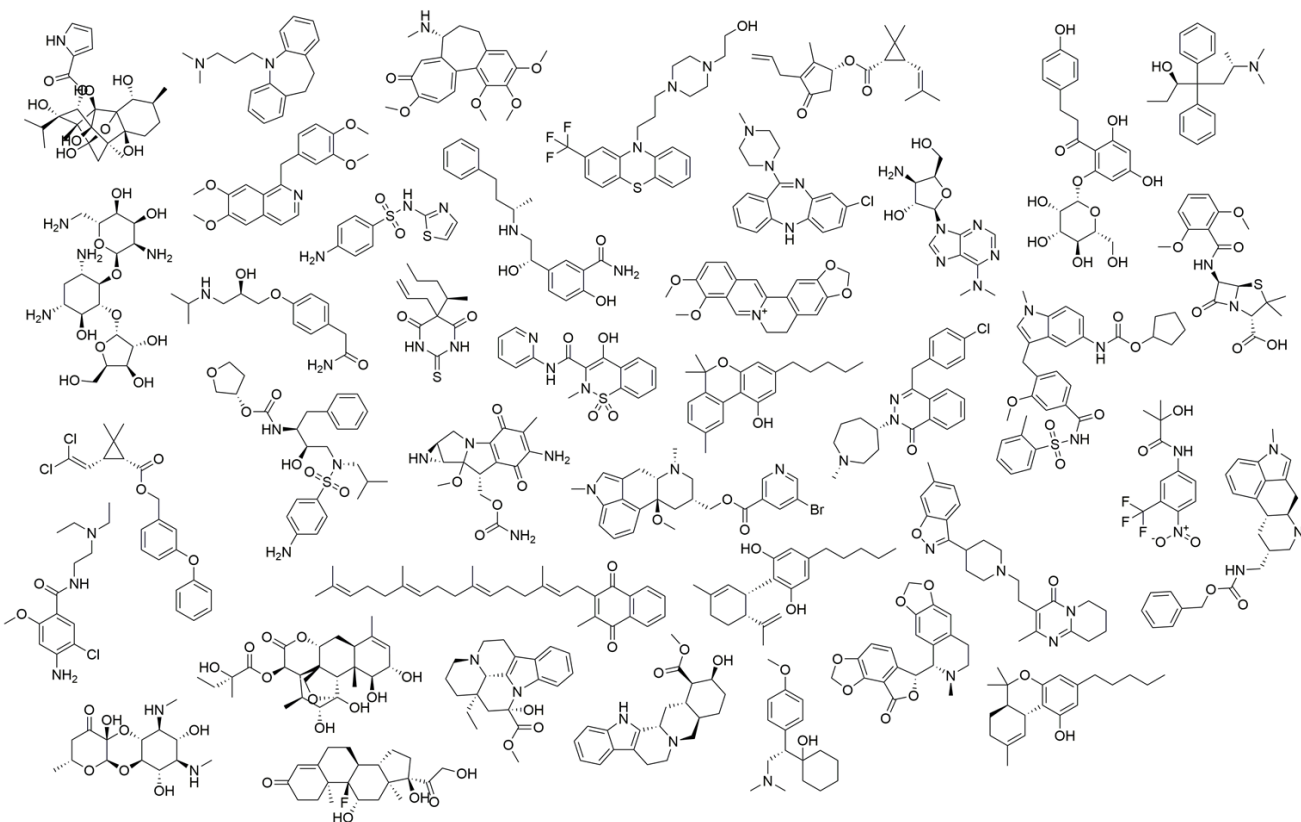
Synthesis  
Planning

Wet Lab  
Candidates



**ca. 166 billion**  
**GDB17**

Ruddigkeit et. al., *J. Chem. Inf. Model.* **2012**, 52, 11, 2864–2875



**ca. 10 million**

**FDB17**

Visini et. al., *J. Chem. Inf. Model.* **2017**, 57, 4, 700–709

**ca. 10 million**

**GDBMedChem**

Awale et. al., *Mol. Inf.* **2019**, 38, 1900031

**ca. 10 million**

**GDBChEMBL**

Bühlmann et. al., *Front. Chem.* **2020**, 8 (46)

# TACKLING COMPUTER AIDED SYNTHESIS COLLABORATIVELY



SYNTHESIS  
PLANNING



WET LAB  
CANDIDATES

How do we **synthesise** the **compounds** that are **generated** or **suggested** by a **chemist** or **computer**?

*u<sup>b</sup>*

b  
UNIVERSITÄT  
BERN



AstraZeneca

# PRIOR ART

- Guide chemists in the synthesis of molecules
- Codify and organise the techniques used in organic synthesis
- Propose a variety of synthetic routes given any molecule

The responsibility for the final evaluation of the merit of the routes lies with the chemist.

***We need a method for getting feedback on lots of routes quickly – improving synthesis plans by iteration***

5



## LHASA—Logic and Heuristics Applied to Synthetic Analysis

DAVID A. PENSACK

Central Research and Develop. Dept., E. I. du Pont de Nemours and Co.,  
Wilmington, Del. 19898

E. J. COREY

Dept. of Chemistry, Harvard University, Cambridge, Mass. 02138

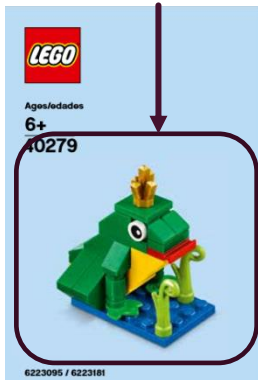
Despite the wealth of knowledge about various chemical reactions, there exists no formal framework of interrelationships to guide the chemist in the synthesis of even moderately complex molecules. The LHASA (Logic and Heuristics Applied to Synthetic Analysis) project is an attempt to codify and organize the techniques used in organic synthesis.

One important aspect of the project has been the writing of a general purpose computer program which will aid the laboratory chemist and will employ both the basic and more complex techniques for synthetic design as elucidated by this study. The program (hereafter also called LHASA) is intended to propose a variety of synthetic routes to whatever molecule it is given. The responsibility for final evaluation of the merit of the routes lies with the chemist. The

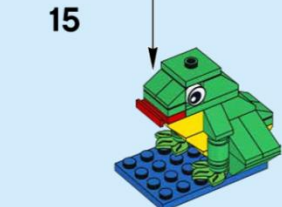
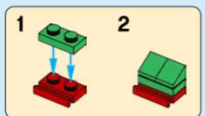
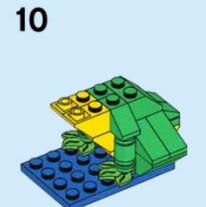
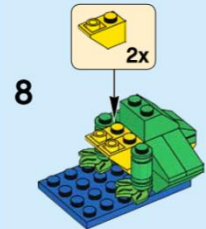
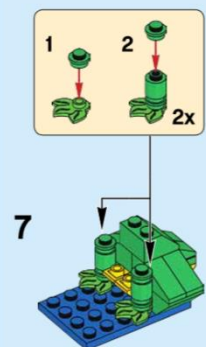
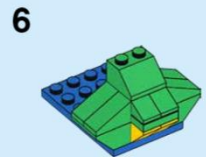
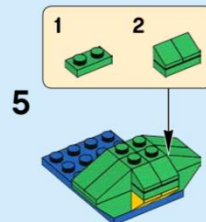
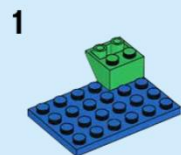
# LEGO – ANALOGY TO SYNTHESIS PLANNING

## Building Lego as an analogy to synthetic route planning

### Target Molecule



Warning! Choking hazard: Small parts.  
Achtung! Erstickungsgefahr: Kleine Teile.  
Attention! Danger d'étouffement: Petits éléments.  
Avvertenza! Rischio di soffocamento: Piccole parti.  
Waarschuwing! Verslikingsgevaar: Kleine onderdelen.  
Advertencia! Peligro de asfixiantamiento: Partes pequeñas.  
Advarsel! Kvælningsfare: Små dele.  
Vigbútur! Kjöfnunahætta. Lítil hlutir.  
Varoitus! Tukehtumisvaara. Pieniä osia.  
Varoitus! Kivertysvaara. Pieniä osia.  
Advarsel! Kvælningsfare. Små delar.  
Atención! Riesgo de asfixia. Pequeñas partes.  
Προειδοποίηση! Κίνδυνος ασφυξίας. Μικρά μέρη.  
警告! 窒息の危険があります。細い部品があります。  
경고! 질식사 위험. 작은 부품 포함.  
Внимание! Риск удушья. Мелкие детали.  
Ostrzeżenie! Niebezpieczeństwo uduszenia się. Małe części.  
UPOZORENIE! Nebezpečí zadušení. Malé části.  
UPOZORENIE! Nebezpečenie duszenia. Małe części.  
Figyelmeztetés! Fulladásveszély. Kis alkatrészek.  
Попередження! Ризик задихатися. Дрібні деталі.  
UPOZORENIE! Opasnost od gušenja. Mali dijelovi.  
УПОЗОРЕНІ! Опасност од гушення. Мали делови.  
Предупреждение! Опасност од задушение. Мали делови.  
Opozorilo! Nevarnost zadušitve zaradi tulke. Majhni deli.  
Avertament! Pericol de sufocare internă. Partii mici.  
Внимание! Опасност от задушване. Мали части.  
Brødningsskilt! Advarsel om risiko. Små detaljer.  
Hoztatal! Figyelmeztetés! Veszélyes rész.  
prejimat! Pavejos upspringti. Smalios detalės.  
Uviti! Bodljeno! Varnostno! (Nevarnost!) Kijelo! Pericolo!  
تحذير! خطر الاختناق! أجزاء صغيرة.  
Peringatan! Bahaya tersedak. Bagian kecil.  
Atvarsel! Bahaya tersedak. Bahagian kecil.



### Known Building Blocks





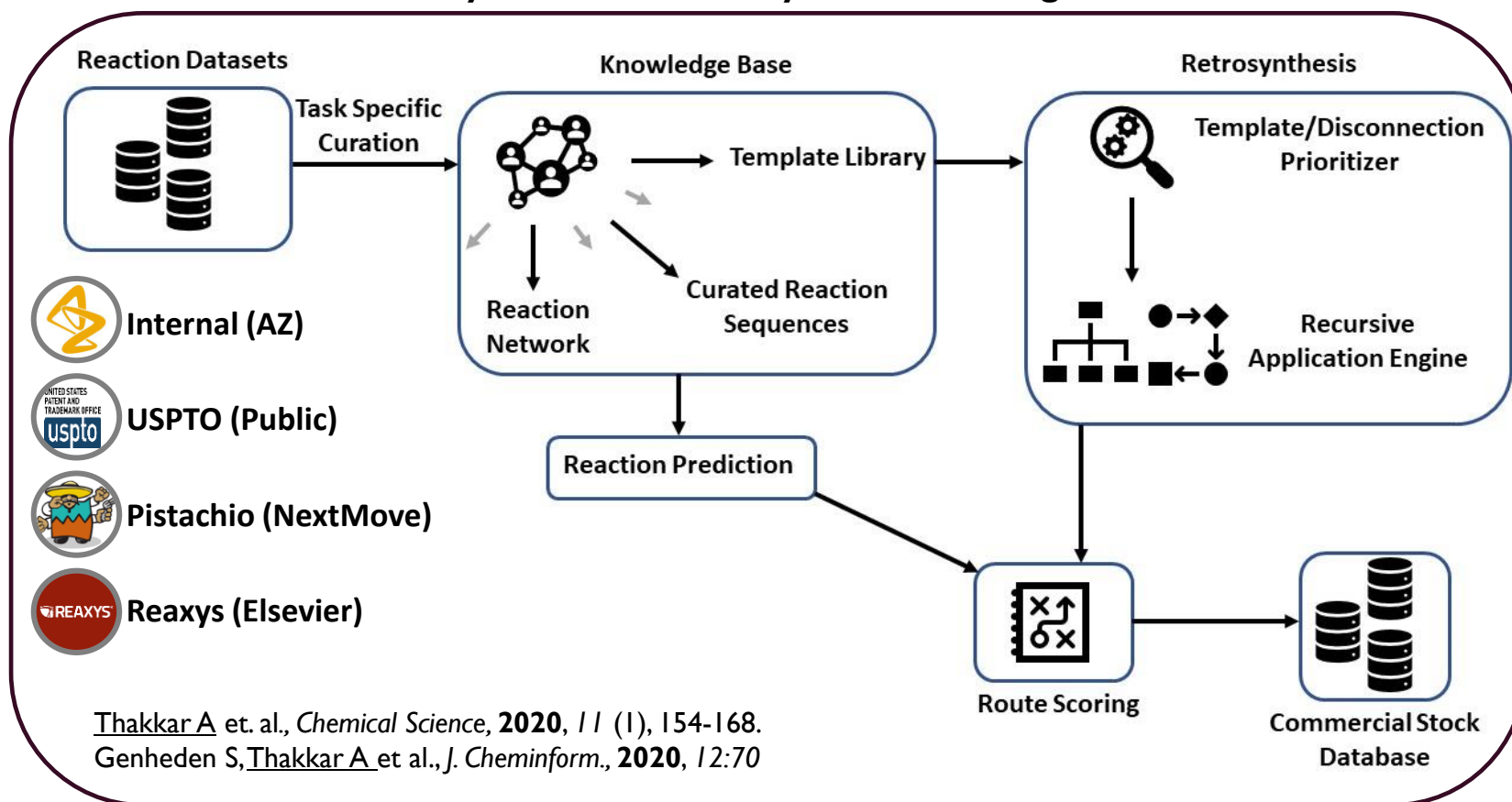
# OPEN-SOURCE APPROACHES FOR SYNTHESIS PLANNING

Code

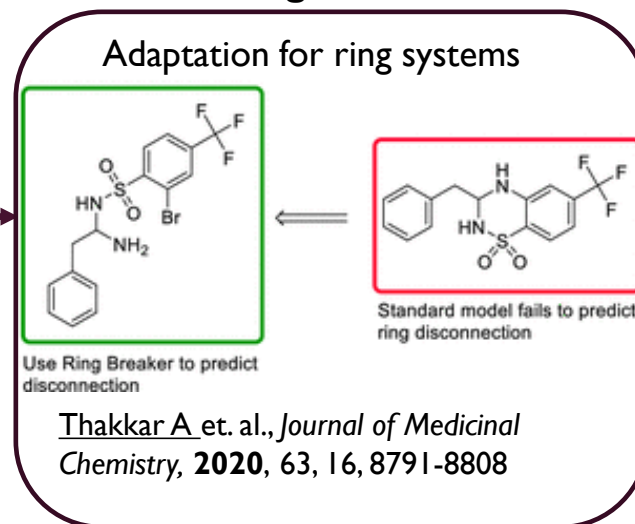
Paper



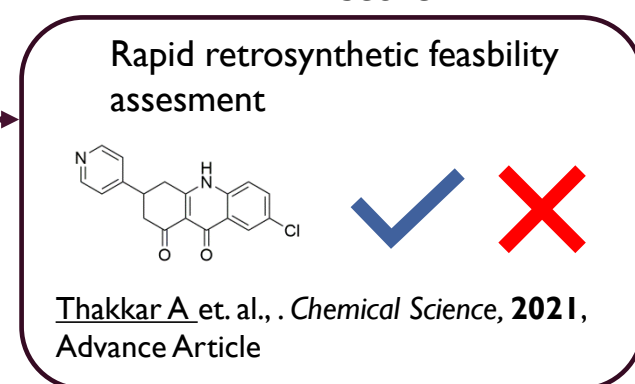
## AiZynthFinder – Retrosynthetic Planning



## RingBreaker



## RA score



# TEMPLATE EXTRACTION

Code

Paper



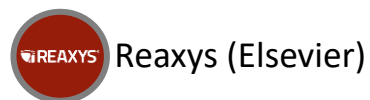
## Dataset



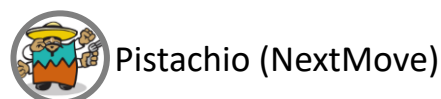
Proprietary



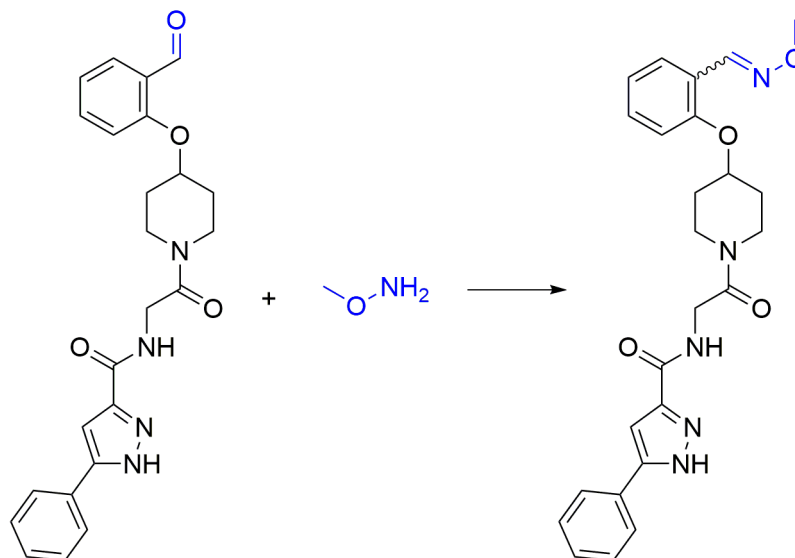
Literature



Patents

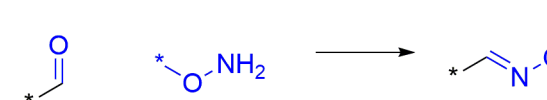


## Reaction in Dataset

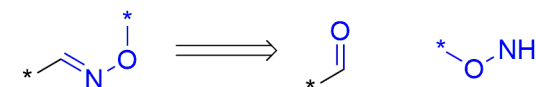


## Extract Rules

Forward reaction rule:



Retro reaction rule:



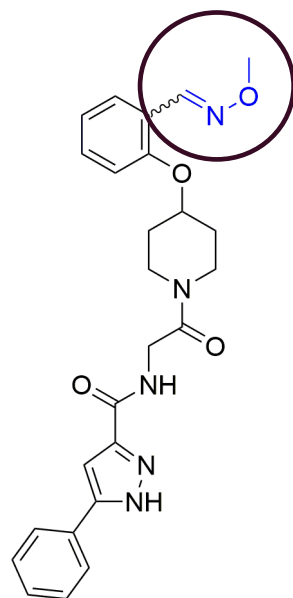
Thakkar et al., *Chemical Science*, **2020**, 11 (1), 154-168.  
Segler et al., *Nature*, **2018**, 555, 604-610



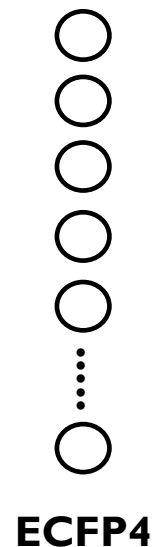
# NEURAL NETWORK TRAINING – TEMPLATE PRIORITISATION

Code

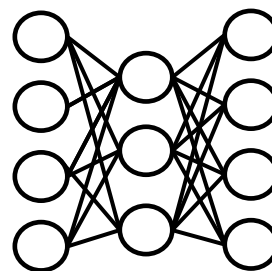
Paper



**Product**



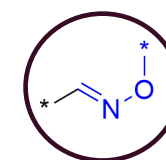
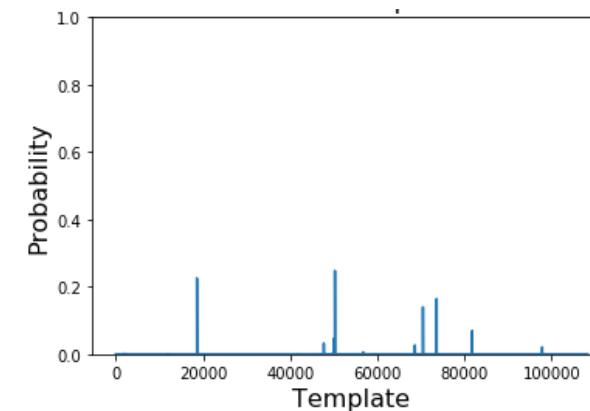
**ECFP4**



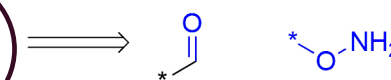
**Prioritisation  
Network**



**Predict  
Retrosynthesis  
Template**



Retro reaction rule:

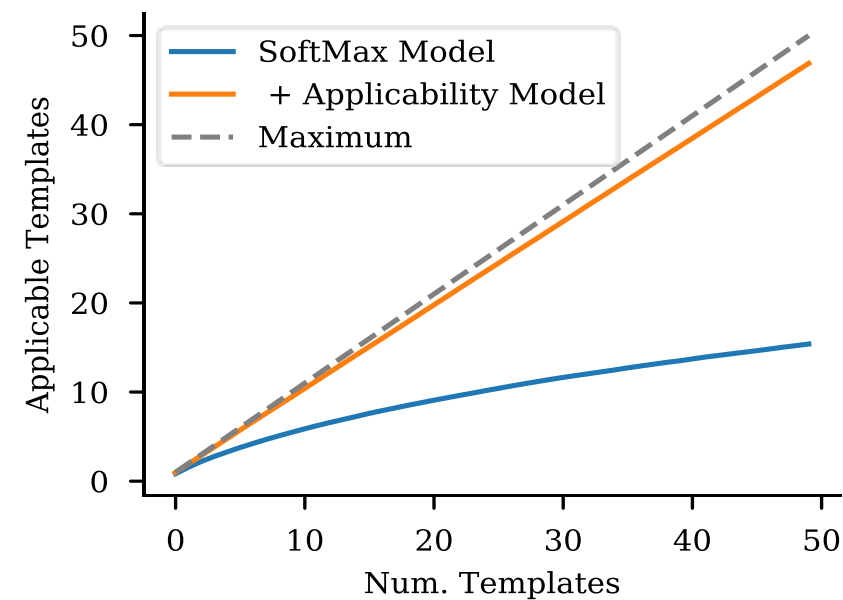
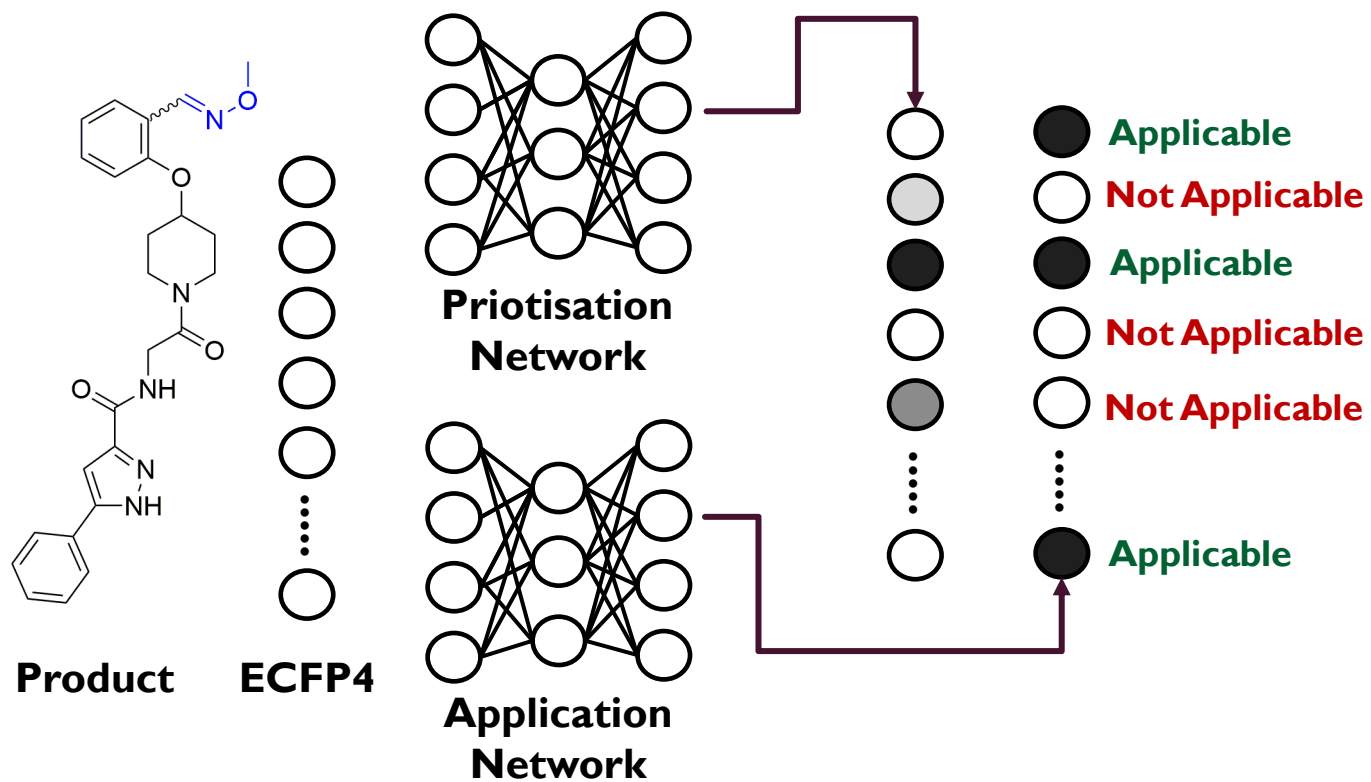
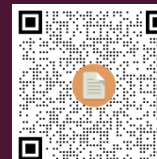


**Is the template  
applicable in silico?**

**Product and template  
must have a  
substructure match**

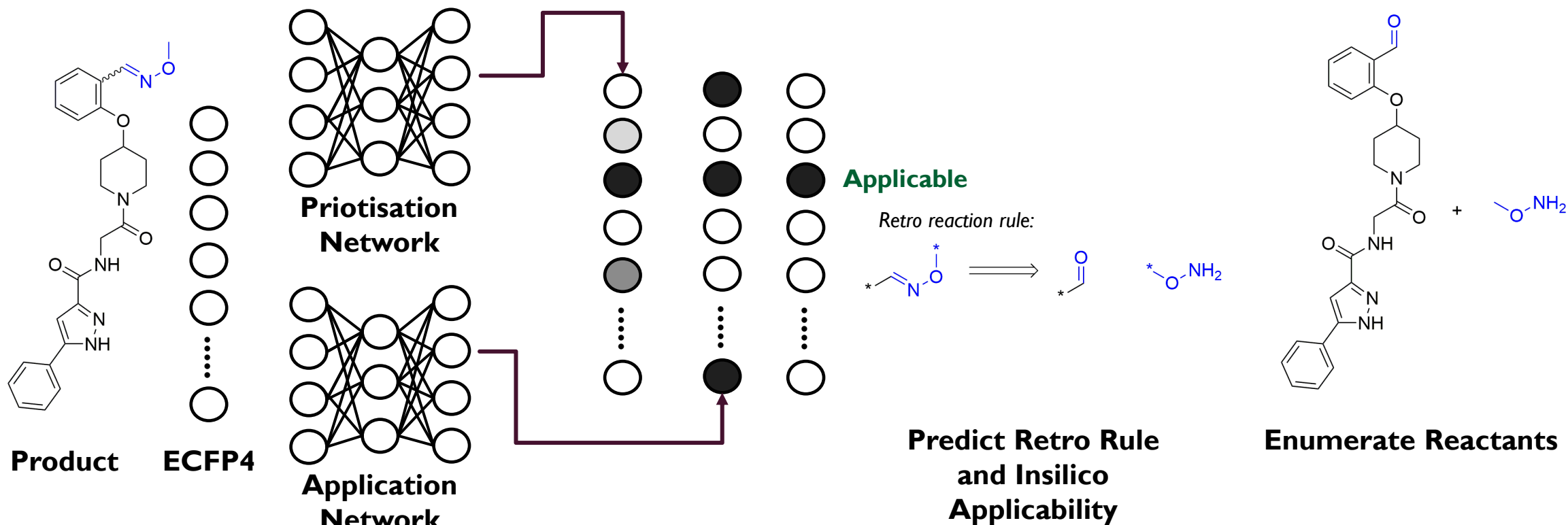
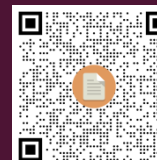
# IMPROVING TEMPLATE PRIORITISATION

Paper

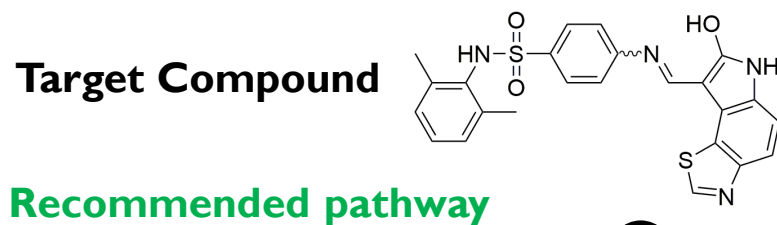


# IMPROVING TEMPLATE PRIORITISATION

Paper



# Paper



**Node expansion according to:**

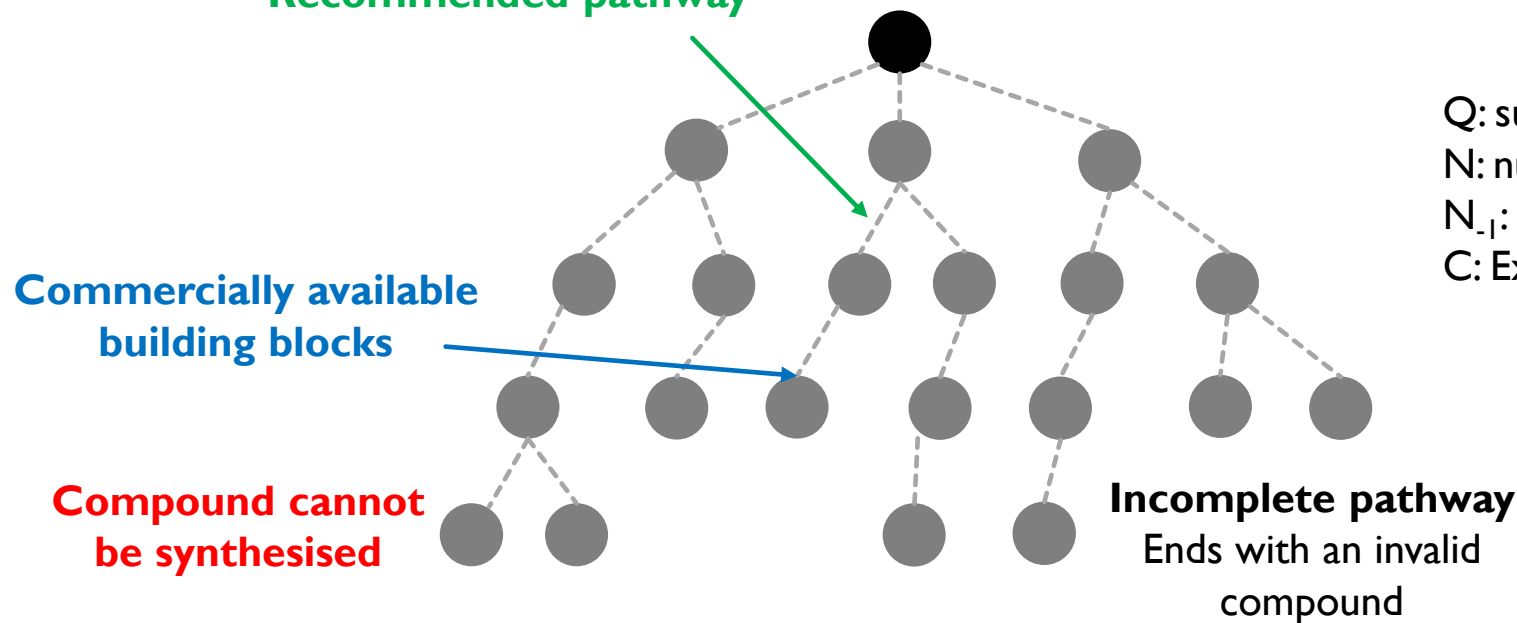
$$UCB = \frac{Q}{N} + C \times \sqrt{2 \times \frac{\ln N_{-1}}{N}}$$

Q: sum of previous rewards

N: number of child node visitations

$N_p$ : Number of parent node visitations

### C: Exploration/Exploitation hyperparameter



# ACCESSIBILITY OPTIONS FOR END USERS

Code

Paper



## GitHub Repo

- For developers
- For end users familiar with scripting and command line

README.md

Unwatch 23 Unstar 170 Fork 39

### AiZynthFinder

license MIT tests passing codecov 83% code style black release v3.0.0 Open in Colab

AiZynthFinder is a tool for retrosynthetic planning. The algorithm is based on a Monte Carlo tree search that recursively breaks down a molecule to purchasable precursors. The tree search is guided by a policy that suggests possible precursors by utilizing a neural network trained on a library of known reaction templates.

### Prerequisites

Before you begin, ensure you have met the following requirements:

- Linux, Windows or macOS platforms are supported - as long as the dependencies are supported on these platforms.
- You have installed [anaconda](#) or [miniconda](#) with python 3.6 - 3.8

The tool has been developed on a Linux platform, but the software has been tested on Windows 10 and macOS Catalina.

### Installation

#### For end-users

First time, execute the following command in a console or an Anaconda prompt

```
conda env create -f https://raw.githubusercontent.com/MolecularAI/aizynthfinder/master/env-users.yml
```

## Graphical User Interface

- Deployed internally in AZ
- Allows interactive exploration
- Functionality not in repo

## IBM RXN Graphical user interface below:

IBM RXN Projects < test on Monday < Retrosynthesis Outcome

test on Monday Sequence 6 Created by: Teodoro Laino

Share Copy & Create New Add Sequence to Collection Start Robo Process

AI model: 2020-07-01 Automatic Mode High Confidence Confidence: 0.998 Optimization score: 1.1 Optimization time: 23 s Predictions: 8 4/F 30

Molecule commercially available on eMolecules.com Not able to find a synthetic path

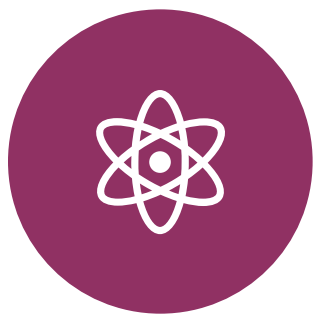
0.998 - Unrecognized

Sequences Generated

- Sequence 0
- Sequence 1
- Sequence 2
- Sequence 3
- Sequence 4
- Sequence 5
- Sequence 6
- Sequence 7
- Sequence 8
- Sequence 9
- Sequence 10
- Sequence 11

© 2019 IBM

# ACCESSIBILITY OPTIONS FOR END USERS



**GDB CHEMICAL  
SPACE**



**AI PLANNED  
SYNTHESIS**



**BROWSER BASED  
EXPLORATION**

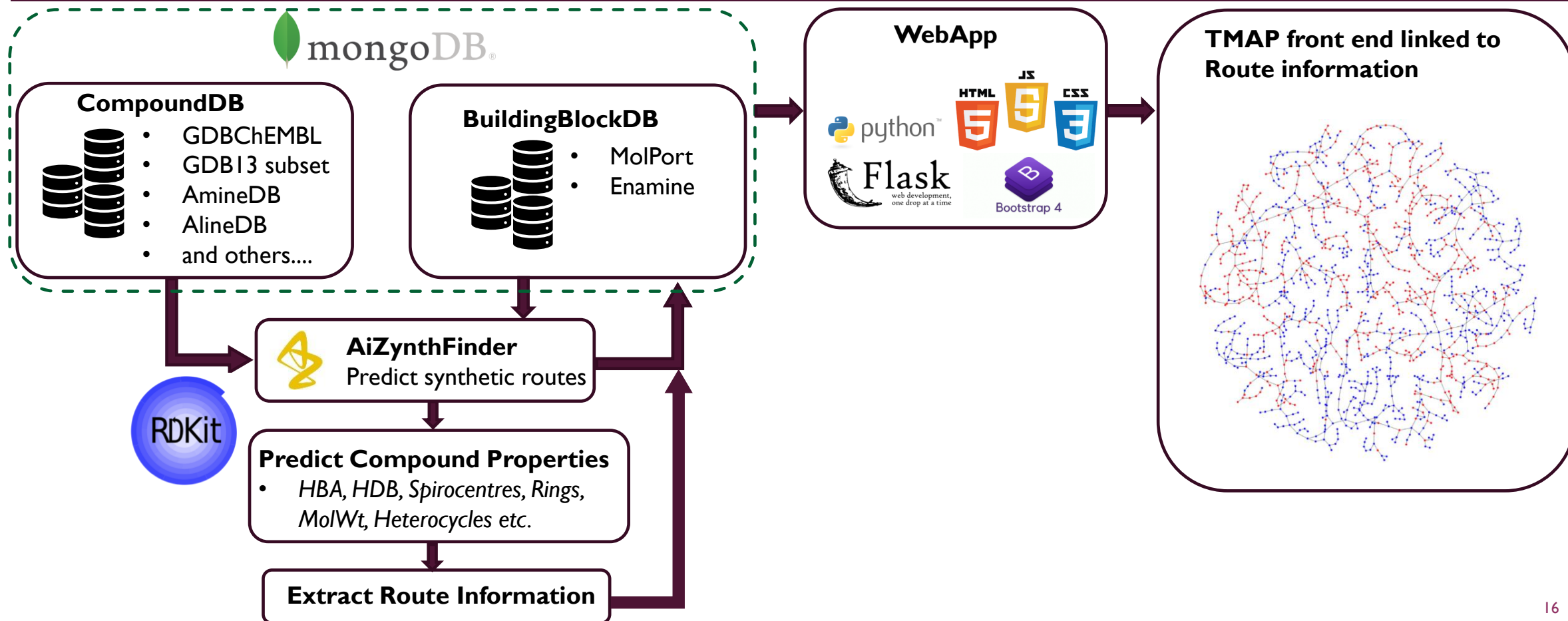


**EXPERIMENTAL  
ENGAGEMENT**

- × Scripting or command line knowledge
- × Step wise compound submission
- × Wait for batch runs
- × Switching between tools for visualisation, calculation of properties, synthesis prediction, prioritisation



# GDB ROUTE BROWSER



# TMAP – VISUALISATION OF CHEMICAL LIBRARIES

- Feature vector generated using minhashing
- Index into LSH forest
- Generation of k-nearest neighbour graph
- Minimum spanning tree algorithm for layout
- Embed into 2d using co-ordinates

<https://tmap.gdb.tools/>

D. Probst, J.-L. Reymond, Visualization of very large high-dimensional data sets as minimum spanning trees. *J Cheminform* **12**, 12 (2020).

tmap

## GETTING STARTED

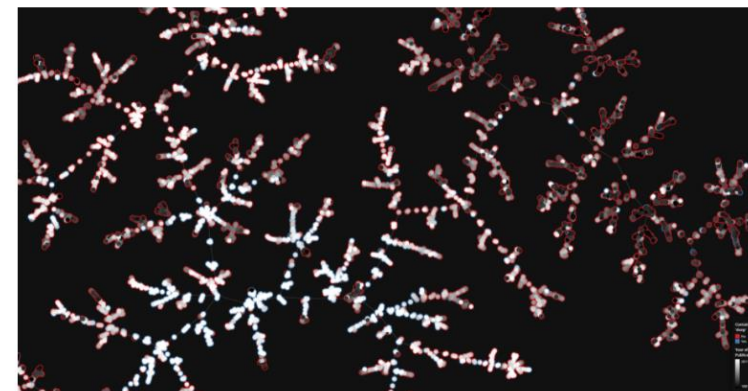
Supported Environments  
Installation  
Laying out a Simple Graph  
MinHash  
LSH Forest

## EXAMPLES

COIL20  
MNIST  
Fashion MNIST  
ChEMBL  
FDB17 and ChEMBL  
Natural Product Atlas  
DSSTox  
Protein Data Bank  
RNA Sequencing  
ProteomeHD  
PubMed Central  
MiniBooNE  
Gutenberg  
NIPS  
Drugbank  
Flowcytometry  
MOLECULENET.AI  
Quantum Mechanics  
Physical Chemistry

## Getting started

GitHub



**tmap** is a very fast visualization library for large, high-dimensional data sets. Currently, tmap is available for Python.

**NEW:** We now provide a web-service that allows for the creation of TMAP visualizations for small chemical data sets.

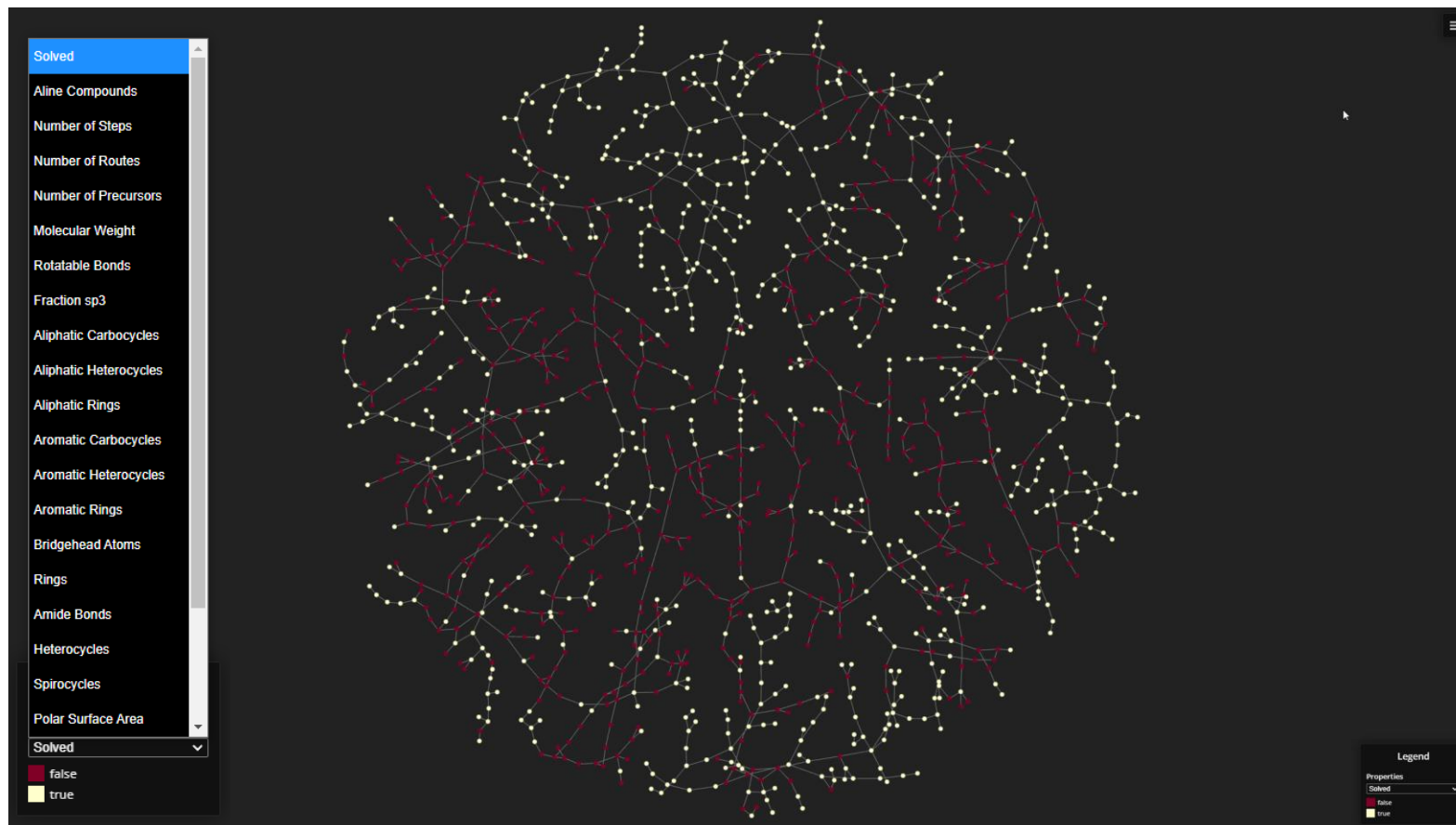
Try TMAP

## Supported Environments

Language	Operating System	Status
Python	Linux	Available
	Windows	Available <sup>1</sup>
	macOS	Available
R		Available <sup>2</sup>

# GDB RETROSYNTHESIS

Dataset	Compounds	% Solved by AiZynthFinder	Number of Routes	Number of Steps	TMAP Created
AmineDB	44,929	17.7	427,493	2,829,675	Y
DiamineDB	1,323	56.2	12,207	59,987	Y
GDB13_ABCDEFGH	994,840	19.7	8,501,323	54,933,389	N
GDBChEMBL_A	1,490,508	33.8	12,826,692	75,029,188	N
<b>Total</b>	<b>2,531,600</b>	<b>28.0</b>	<b>21,767,715</b>	<b>132,852,239</b>	



## COLORATION BY PROPERTIES

### Diamine Database

Number of rings: Max 2

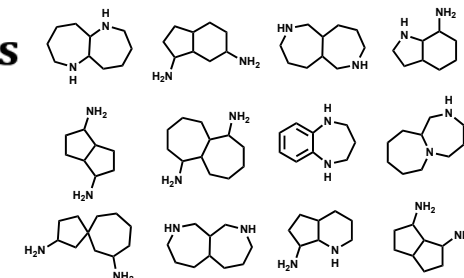
Ring Sizes: 5, 6, and 7

Number of Amines: Max 2

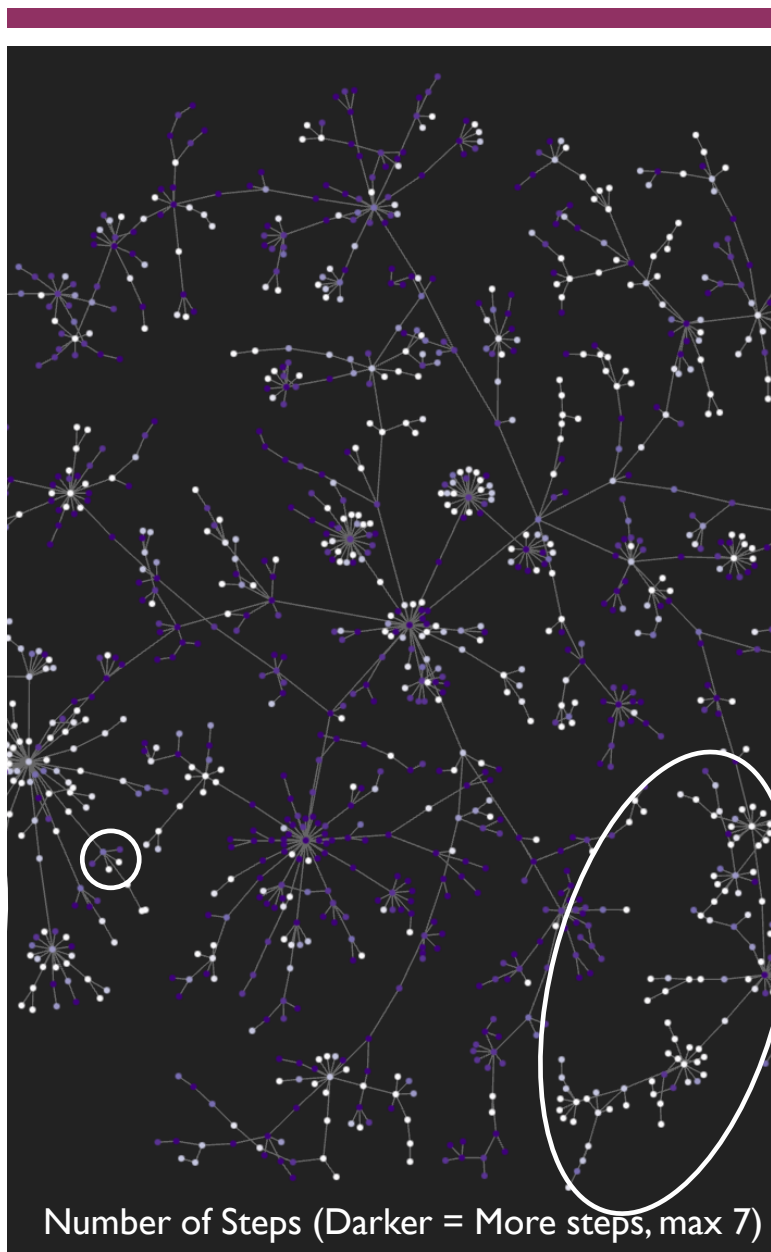
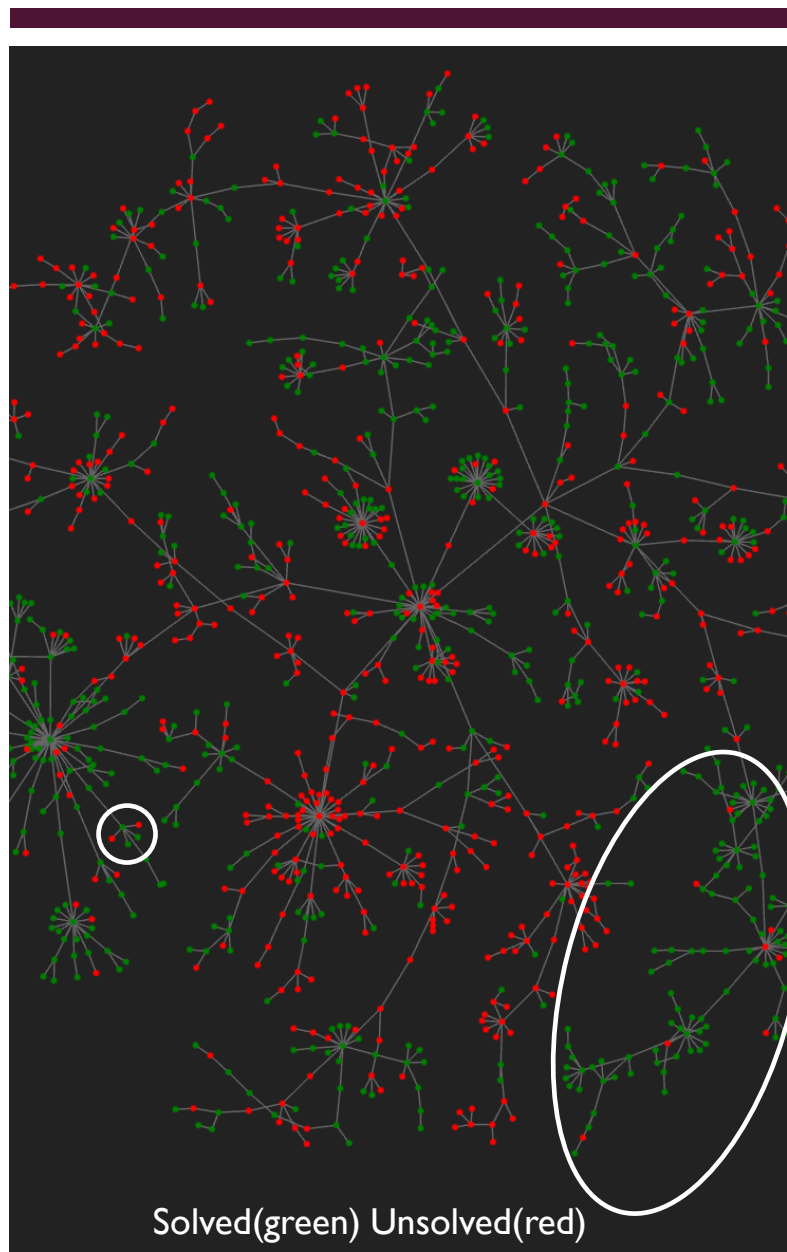
Exo- and Endocyclic

**1323 Molecules**

**57 % Novel**

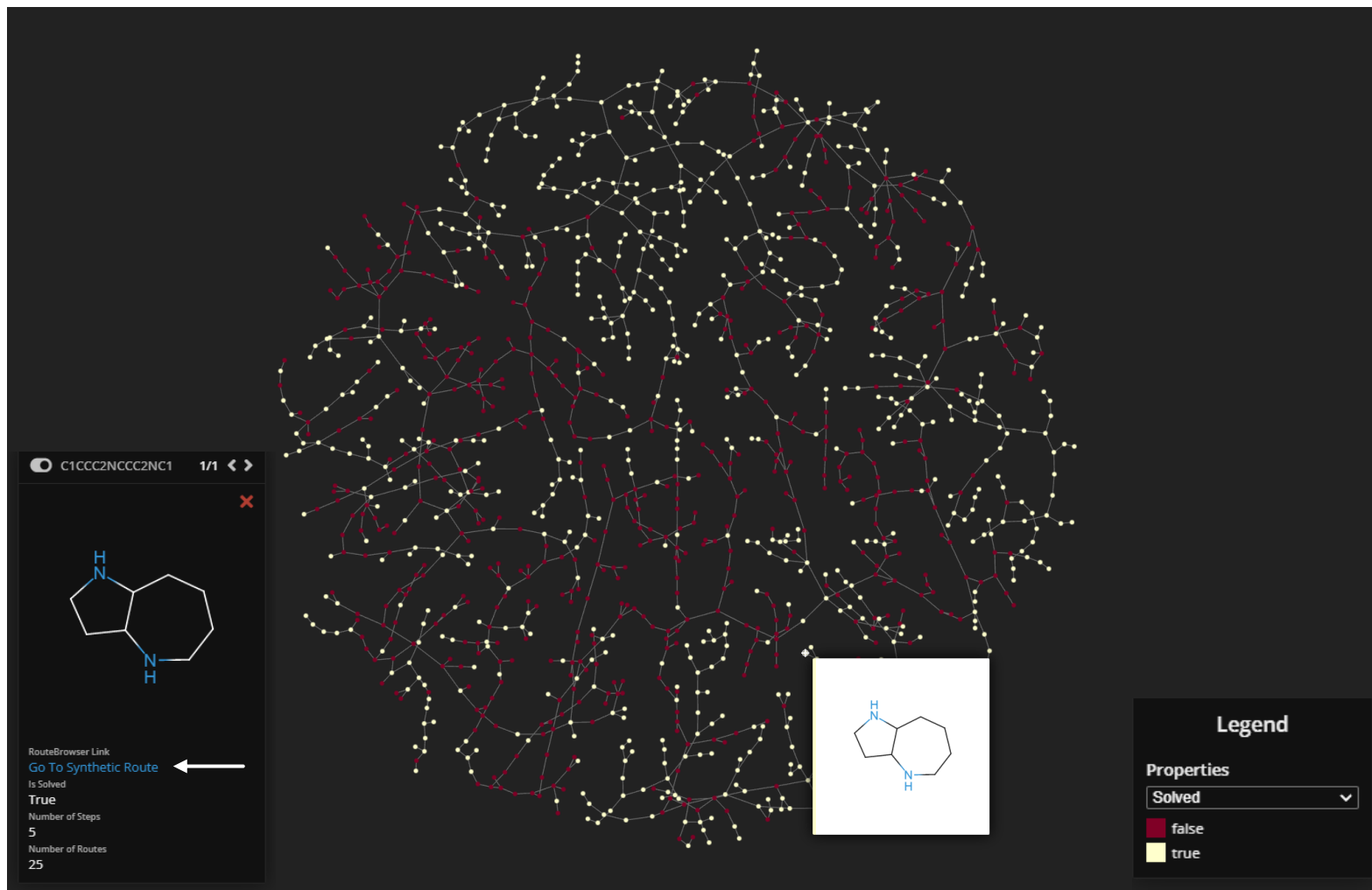


- Helps prioritisation of compounds
- View of related compounds
- Relationships between compounds in the library



## ENABLING A RICHER UNDERSTANDING OF CHEMICAL SPACE

- Short synthetic routes to privileged scaffolds
- Chemists know the reactions for disconnection
- AiZynthFinder knows the reactions and the building blocks available



PRECOMPUTED  
SYNTHETIC ROUTES FOR  
EACH COMPOUND



## Compound Identifiers:

**SMILES:** C1CCC2NCCC2NC1

**InChI:** InChI=1S/C8H16N2/c1-2-5-9-8-4-6-10-7(8)3-1/h7-10H,1-6H2

**InChI key:** SLCIXBLNPIPMOC-UHFFFAOYSA-N

## Route Details:

25 routes predicted in 180.65 seconds

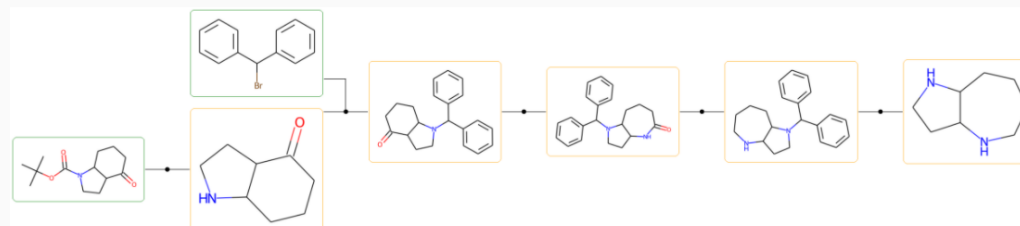
Currently showing: Route 1 of 25 ▸

Number of Steps: 6

Number of Precursors: 3

## Predicted Route:

\*Compounds highlighted in a green box are commercially available



Like

Dislike

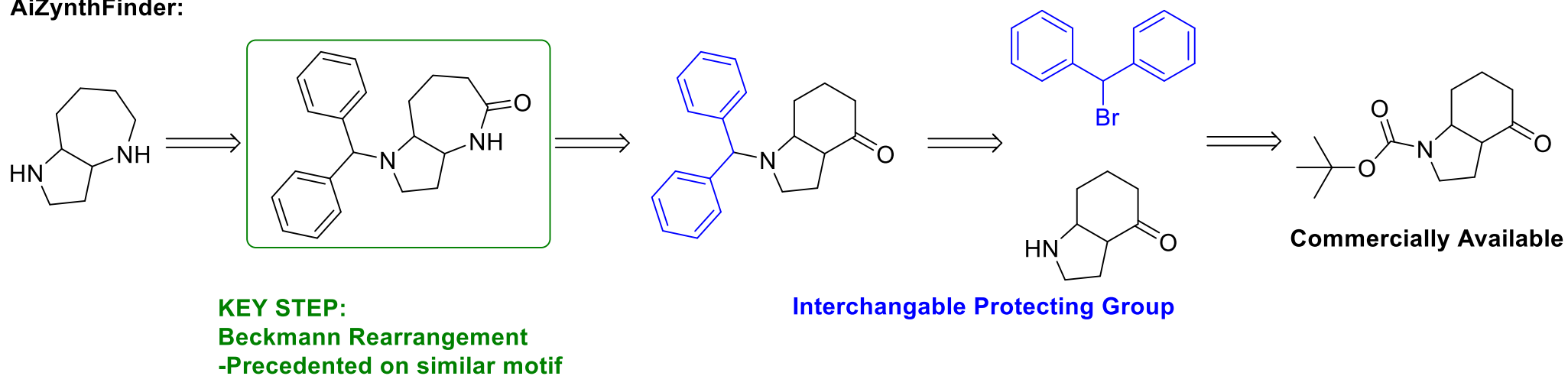
## Precursors Available:

# APPLICATION TO GDB – EXPERIMENTAL VALIDATION

ASKCOS: **No Route Found**

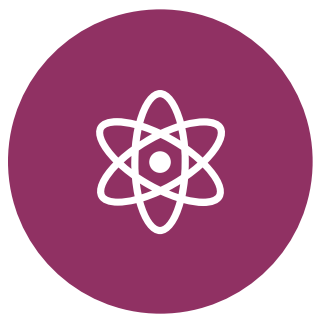
IBM RXN: **No Route Found**

AiZynthFinder:



Successfully synthesized in the lab!

# LINKING LIBRARY VIZUALISATION TO SYNTHESIS



**GDB CHEMICAL  
SPACE**



**AI PLANNED  
SYNTHESIS**



**BROWSER BASED  
EXPLORATION**



**EXPERIMENTAL  
ENGAGEMENT**

- × Scripting or command line knowledge
- × Step wise compound submission
- × Wait for batch runs
- × Switching between tools for visualisation, calculation of properties, synthesis prediction, prioritisation

- ✓ One tool
- ✓ Precomputed routes – fast access
- ✓ Allows whole library to be visualised
- ✓ Enables easier prioritisation using synthetic route information

# SUMMARY

- AiZynthFinder – An Open Access Retrosynthetic Planning Tool
- Established methods for improving AiZynthFinder using data augmentation and specialised models.
- Built specialised tools to target specific needs arising from chemists.
- The tools have demonstrated project impact as shown for the case of GDB.
- Nearly everything is open source and published open access for the community!
- Linking visualisation to chemical synthesis planning will be released in the near future (currently writing and refactoring the code)

# ACKNOWLEDGEMENTS

## Professor Jean-Louis Reymond

Dr Ola Engkvist (AZ)

Dr Esben Jannik Bjerrum (AZ)

Dr Samuel Genheden (AZ)

Reymond Group, University of Bern, Switzerland

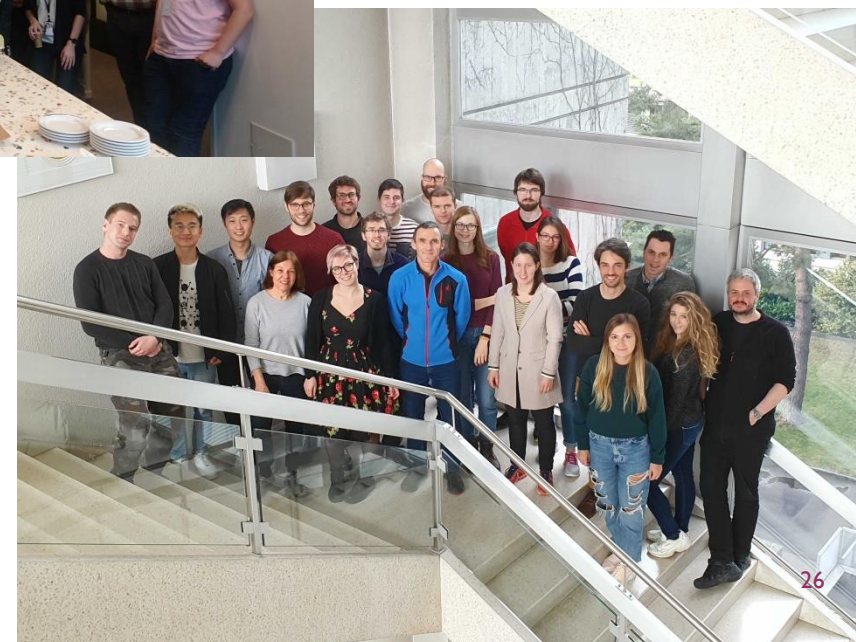
Molecular AI, AstraZeneca, Sweden

Global Chemistry Community, AstraZeneca

Elsevier

NextMove Software

The reviewers for their useful feedback.



**Feedback or Questions?**

**Email:**  
[amol.thakkar@unibe.ch](mailto:amol.thakkar@unibe.ch)