

Paradigm shift in synthesis planning: how to exploit artificial intelligence

Samuel Genheden

MolecularAI, Discovery Sciences, AstraZeneca R&D, Gothenburg



2024-09-16

#### Acknowledgements



Annie Westerlund



Lakshidaa Saigiridharan



Mikhail Kabeshov



Thierry Kogej



Christos Kannas



Ola Engkvist



Varvara Peter Hartog Voinarovska







Supriya Siva Kancharla Manohar Rocío Mercado



Emma Rydholm



# Synthetic Chemistry in Pharma

- Medicinal chemistry:
  - Invent a molecule to treat human disease
  - Design & synthesize new molecules: 1000s per p
  - Small scale: <0.005 g to 50 g
  - Diversity-oriented
  - Speed is key
- Process chemistry:
  - Make large quantities of that molecular for studies and clinic
  - Deliver large quantities of API for studies
  - 50 g to >1,000,000 g scale
  - Purity, sustainability and cost are key





## Computer-aided synthesis planning

- Retrosynthesis analysis
  - One-step retrosynthesis
  - Multi-step retrosynthesis or route prediction



- Database searches
  - Support for predictions
  - Experimental details



## Software landscape

**Free solutions** 

#### **Commercial solutions**

Reaxys





IBM RXN for Chemistry

+ Many GitHub repositories





CAS



## Reaction data landscape

#### **Proprietary sources**

#### Reaxys®



- > 30 M reactions
- 90% from research articles
- > 10 M reactions
- US, European and Asian patents



- Closed company resource
- AZ ELN is > 2 M reactions



- > 100 M reactions
- Rarely used for ML

#### Public-domain sources



ORD

- > 2M reactions
- Open-source extract from US patents
- Questionable quality
- Up to 2016
- > 2M reactions
- Initiative from academia with industry support
- Open data format
- Contains US patents and HTE datasets



# AiZynthFinder and template-based retrosynthesis

#### **Reaction templates**

#### **One-step retrosynthesis**

#### **Monte-Carlo tree search**

Update



1. Databases provide atom-mapped reactions



2. Extract chemical transformation rules: templates







# Template-free model: Chemformer





pubs.acs.org/icim

#### Do Chemformers Dream of Organic Matter? Evaluating a Transformer Model for Multistep Retrosynthesis

Published as part of the Journal of Chemical Information and Modeling virtual special issue "Modeling Reactions from Chemical Theories to Machine Learning."

Annie M. Westerlund,\* Siva Manohar Koki, Supriya Kancharla, Alessandro Tibo, Lakshidaa Saigiridharan, Mikhail Kabeshov, Rocío Mercado, and Samuel Genheden



ABSTRACT: Synthesis planning of new plarmaceutical compounds is a vellknown bottleneck in modern drug design. Template-free methods, such as transformers, have recently been proposed as an alternative to template-based methods for single-step retrosynthetic predictions. Here, we trained and evaluated a transformer model, called the Chemformer, for retrosynthesis predictions within drug discovery. The proprietary data set used for training comprised ~18 M reactions from literature, patents, and electronic lab notebooks. Chemformer was evaluated for the purpose of both single-step and multistep retrosynthesis. We found that the single-step performance of Chemformer was especially good on reaction classes common in drug discovery, with most reaction classes howing a



Article

top-10 round-trip accuracy above 0.97. Moreover, Chemformer reached a higher round-trip accuracy compared to that of a template-based model. By analyzing multistep retrosynthesis experiments, we observed that Chemformer found synthetic routes.

Westerlund et al. J. Chem. Inf. Model. 2024, 64, 3021–3033



# Template-free / Generative models

- Treat retrosynthesis as a language problem
- Translating from product to reactants

- centra nature Moleci COMMUNICATIONS Chemi Philippe S Costas Be Check for u ARTICLE IBM Resea Departmen OPEN https://doi.org/10.1038/s41467-020-19266-y Departmen State-of-the-art augmented NLP transformer models for direct and single-step retrosynthesis
  - Igor V. Tetko<sup>®</sup> <sup>1,2⊠</sup>, Pavel Karpov<sup>1,2</sup>, Ruud Van Deursen<sup>®</sup> <sup>3</sup> & Guillaume Godin<sup>®</sup> <sup>3⊠</sup>
- Model will learn chemical rules and the prioritization of those rules

#### 1. Pretrain a *transformer* model (BART) to learn SMILES



230 M parameters

#### 2. Fine-tune a model to perform retrosynthesis

## On the promises of template-free models

- Template-free retrosynthesis models do not need
  - Reactions with atom-mapping
  - Extracted templates
- Better scaling to larger chemical spaces
- Possibility to extrapolate outside known space



## Data sources and trained models

#### Reaxys

- > 30 M reactions
- 90% from research articles



- > 10 M reactions
- US, European and Asian patents



• AZ ELN is > 2 M reactions



18M unique reactions



## Chemformer vs Template-based: Multi-step performance

- Evaluate route prediction on 5K compounds designed by AZ chemists
- Stock is AstraZeneca internal building block collection (5M)
- Run 100 iteration of retrosynthesis
- Implement novel algorithms for speeding-up Chemformer performance





Solvability case	Percentage
Solved by both	71.6
Only solved by chemformer	23.4
Only solved by template-based	0.74
Solved by none	4.2
Solved by any	95.8

## Do Chemformers dream of organic matter?





#### Are the suggested reactions feasible?

Round-trip accuracy calculated with forward Chemformer, showing how feasible the predicted reactions are



## Are the models exploring different chemical space?

*Extract templates from predicted reaction and look at space of unique templates* 



#### b) Chemformer vs. library

## What is ahead of us?

- Chemformer appears to be a viable option to template-based retrosynthesis
- Hallucination occurs, but is perhaps acceptable
  - Hard to evaluate plausibility on a large scale
- Two major issues
  - Computational speed (engineering)
  - Interpretability and experimental support (science)
- Possible solution: a two-staged approach
  - 1. Predict with template-based retrosynthesis
  - 2. If fails to find synthetic route, try again with Chemformer



# Constrained retrosynthesis

#### Constrained synthesis planning with disconnection-aware transformer and multi-objective search

24 May 2024, Version 1

Working Paper

Annie M. Westerlund 🧔, Lakshidaa Saigiridharan 🧐, Samuel Genheden

Show author details ~

🚯 This content is a preprint and has not undergone peer review at the time of posting.

Download 66 Cite	O Comment
------------------	-----------

#### Abstract



Designing synthesis routes with shared intermediates for a set of target compounds is a common task in drug tep retrosynthesis tools such as AiZynthFinder are frequently used by chemists to generate Although these tools can find solved routes for a majority of target compounds, they may not vhich comply with specific bond constraints. Such bond constraints could be defined in the synthesis routes with common intermediates for the set of compounds. Here, we present a novel proach which aims to generate routes in the feasible region defined by these constraints. The ivided into bonds to break and bonds to freeze. First, we introduce a filter in the search which

## Synthesis of compound series





- In drug discovery you often explore series of related compounds
- Investigate structure-activity relationship for one or more vectors

- Advantageous to synthesis these compounds in similar fashion
- This is not supported by available retrosynthesis tools

Guo et al. Eur. J. Med. Chem. 2019, 178, 767–781. https://doi.org/10.1016/j.ejmech.2019.06.035.

### Constrained synthesis planning



## Novel techniques for enhancing bond-breaking

Selecting and scoring routes with multiple objectives

![](_page_19_Figure_2.jpeg)

Picking node based on Pareto front

Enriching templated-based model with direct bond breaking

![](_page_19_Figure_5.jpeg)

#### Benchmarking on synthetic data

- Extracted synthesis routes from patents and J Med Chem
- Constructing bonds to break or freeze
- Compare five approaches to treating bonds to break
- We can increase the score for bond breaking by steering the search with multi-objective search
- We can increase the number of routes satisfying the constrained with a Chemformer model

![](_page_20_Figure_6.jpeg)

![](_page_20_Figure_7.jpeg)

![](_page_21_Picture_0.jpeg)

# Future outlook

## Challenge 1: route scoring and comparison

- AiZynthFinder has been used for 4 years in production and have predicted synthetic routes for hundred of thousands of compounds
- ... but we don't know if the various improvements made have had a significant effect
- Comparing routes are difficult and mean different things to different people
- What is the best reference set to use for benchmarking?
- Should one compare routes, or should one just score new predictions and monitor an improvement in the score over time?

#### Challenge 2: human-like routes

- The end-goal is that we predict human-like routes with a high confidence of success that can directly taken to the lab
- Currently, we use retrosynthesis as an idea generation
- Currently, route predictions overuse protection/deprotection steps in an unbalanced way, order of steps are non-optimal, steps does not make sense from a forward synthesis perspective, etc.
- To bridge this gap, we need novel algorithms, scores, models etc

#### Acknowledgements

![](_page_25_Picture_1.jpeg)

Annie Westerlund

![](_page_25_Picture_3.jpeg)

Saigiridharan

![](_page_25_Picture_6.jpeg)

Mikhail Thierry Kabeshov

![](_page_25_Picture_8.jpeg)

Christos Kannas

![](_page_25_Picture_10.jpeg)

![](_page_25_Picture_11.jpeg)

Ola Engkvist

![](_page_25_Picture_13.jpeg)

Peter Hartog Varvara Voinarovska

![](_page_25_Picture_15.jpeg)

![](_page_25_Picture_16.jpeg)

![](_page_25_Picture_17.jpeg)

Supriya Siva Kancharla Manohar

Rocío Mercado

![](_page_25_Picture_20.jpeg)

Kogej

Emma Rydholm

![](_page_25_Picture_22.jpeg)

![](_page_25_Picture_24.jpeg)

#### **Confidentiality Notice**

This file is private and may contain confidential and proprietary information. If you have received this file in error, please notify us and remove it from your system and note that you must not copy, distribute or take any action in reliance on it. Any unauthorized use or disclosure of the contents of this file is not permitted and may be unlawful. AstraZeneca PLC, 1 Francis Crick Avenue, Cambridge Biomedical Campus, Cambridge, CB2 0AA, UK, T: +44(0)203 749 5000, www.astrazeneca.com