

Cite this: *Digital Discovery*, 2024, 3, 558

Models Matter: the impact of single-step retrosynthesis on synthesis planning†

Paula Torren-Peraire,^{ID ‡*ab} Alan Kai Hassen,^{ID ‡*cd} Samuel Genheden,^{ID e} Jonas Verhoeven,^{ID b} Djork-Arné Clevert,^{ID d} Mike Preuss^{ID c} and Igor V. Tetko^{ID a}

Retrosynthesis consists of breaking down a chemical compound recursively step-by-step into molecular precursors until a set of commercially available molecules is found with the goal to provide a synthesis route. Its two primary research directions, single-step retrosynthesis prediction, which models the chemical reaction logic, and multi-step synthesis planning, which tries to find the correct sequence of reactions, are inherently intertwined. Still, this connection is not reflected in contemporary research. In this work, we combine these two major research directions by applying multiple single-step retrosynthesis models within multi-step synthesis planning and analyzing their impact using public and proprietary reaction data. We find a disconnection between high single-step performance and potential route-finding success, suggesting that single-step models must be evaluated within synthesis planning in the future. Furthermore, we show that the commonly used single-step retrosynthesis benchmark dataset USPTO-50k is insufficient as this evaluation task does not represent model scalability or performance on larger and more diverse datasets. For multi-step synthesis planning, we show that the choice of the single-step model can improve the overall success rate of synthesis planning by up to +28% compared to the commonly used baseline model. Finally, we show that each single-step model finds unique synthesis routes, and differs in aspects such as route-finding success, the number of found synthesis routes, and chemical validity.

Received 22nd December 2023
Accepted 13th February 2024

DOI: 10.1039/d3dd00252g

rsc.li/digitaldiscovery

1 Introduction

The Design-Make-Test-Analyse (DMTA) cycle is commonly used in small molecule drug discovery to explore novel compounds and indications. Over recent years, it has seen massive changes with the introduction of modern machine-learning approaches.¹ Retrosynthesis, a core task in the Make part of the DMTA cycle of modern drug discovery, is a technique commonly used by organic chemists in synthesis planning. A molecule is successively broken down into smaller subunits until easily synthesizable or purchasable compounds are obtained,^{2,3} where the overall goal is to produce a roadmap for the synthesis of a target compound. With computer-aided

retrosynthesis, researchers in both chemistry and machine learning aim to accelerate the development of chemical synthesis by saving time and resources, addressing more complex molecules or producing more efficient and safe routes. These generated routes can be used by medical chemists to create molecules of interest,⁴ serve as a basis for autonomous chemistry,⁵ or be incorporated into *de novo* drug design to assess synthesizability.⁶

In recent years, retrosynthesis prediction methods have seen an increase in popularity with the alignment of retrosynthesis with modern machine learning approaches⁷ which allow users and developers to consider a larger set of potential synthesis routes. The machine learning field of retrosynthesis prediction is commonly separated into two research fields, referred to as single-step retrosynthesis prediction and multi-step synthesis planning. Where single-step retrosynthesis prediction refers to breaking down a product into a single set of reactants and multi-step synthesis planning refers to the search algorithms used to find synthesis routes leading to purchasable compounds (building blocks).

Specifically, single-step retrosynthesis prediction is a supervised learning task, developed to predict which reactions are relevant to a target molecule, and the corresponding reactants required to produce this reaction. There are two commonly referenced categories of single-step retrosynthesis models,

^aInstitute of Structural Biology, Molecular Targets and Therapeutics Center, Helmholtz Munich – Deutsches Forschungszentrum für Gesundheit und Umwelt (GmbH), Neuherberg, Germany. E-mail: ptorrenp@its.jnj.com

^bIn-Silico Discovery, Janssen Research & Development, Janssen Pharmaceutica N.V., Beerse, Belgium

^cLeiden Institute of Advanced Computer Science, Leiden University, Leiden, The Netherlands. E-mail: a.k.hassen@liacs.leidenuniv.nl

^dMachine Learning Research, Pfizer Research and Development, Berlin, Germany

^eMolecular AI, Discovery Sciences, R&D, AstraZeneca, Gothenburg, Sweden

† Electronic supplementary information (ESI) available. See DOI: <https://doi.org/10.1039/d3dd00252g>

‡ These authors contributed equally to this work.

template-based and template-free.⁷ Template-based methods use reaction templates, an abstraction of the reactions in the data, which summarize the underlying pattern of these reactions. There are different approaches to extracting templates, though in all cases these processes aim to represent the atom and bond structures required to perform a reaction,⁸ where a single template will represent multiple reactions. Template-based methods consider single-step prediction as a classification problem where the task is to predict the appropriate template for the target molecule/product. Examples of template-based methods include NeuralSym,⁹ the first approach in the field which demonstrated the usefulness of using deep neural networks for retrosynthesis prediction, MHNreact,¹⁰ which uses an information retrieval approach to associate products and templates and LocalRetro,¹¹ which uses a graph representation to predict relevant local atom and bond templates for the product.

On the other hand, template-free approaches commonly treat retrosynthesis prediction as a sequence-to-sequence prediction problem,⁸ employing methods seen in natural language processing such as language translation tasks. Instead of extracting and predicting the corresponding templates, the approach aims to learn the underlying reactions to directly predict reactants. Product and reactants are typically introduced as Simplified Molecular-Input Line-Entry System (SMILES), a common text-based representation of chemical entities. Examples of template-free methods include Chemformer,¹² a large pretrained transformer model fine-tuned on the retrosynthesis task, and Augmented Transformer,¹³ a transformer architecture which employs multiple types of augmentation. Other variations of these approaches exist, such as semi-template-based where the molecule is first broken down into subparts then completed to produce chemically viable reactants.^{14–16}

Multi-step synthesis planning focuses on researching novel synthesis route search algorithms using a single-step model to identify retrosynthetic disconnections. Though multiple methods were developed over the decades to address computer-aided synthesis planning,¹⁷ the pioneering machine-learning-based approach in the field uses Monte Carlo Tree Search (MCTS) to plan the traversal of the search tree at runtime guided by a neural network.² Alternative route planning algorithms use an oracle function or heuristics to guide the tree search instead of relying on compute-expensive run time planning. Prominent examples of this are Depth-First Proof-Number (DFPN),¹⁸ which combines classical DFPN with a neural heuristic, Retro*, which combines A* pathfinding with a neural heuristic,¹⁹ or RetroGraph, which applies a holistic graph-based approach.²⁰ Other approaches incorporate reaction feasibility into the tree search²¹ or use synthesizability heuristics in combination with a forward synthesis model.^{22,23} Finally, self-play approaches, motivated by their success in Go,²⁴ learn to guide the tree search by leveraging information gathered from prior runs of synthesis planning.^{25–27}

Single-step retrosynthesis prediction and multi-step synthesis planning are inherently intertwined where the single-step method defines the maximum searchable reaction

network, and the search algorithm tries to efficiently traverse this network by repeatedly applying the chemical information that is stored in the single-step model. However, this connection is not reflected in contemporary research, with only few novel single-step models testing their approaches within a multi-step synthesis planning framework.^{22,28}

Currently, single-step methods are benchmarked by predicting a single retrosynthetic step from a product to reactants. The common benchmark data for these methods, USPTO-50k,^{29,30} consists of around 50k reactions and only has a limited diversity of 10 reaction classes. These methods are typically only tested on reactant prediction and not within multi-step search algorithms, therefore their usability for synthesis planning is not assessed. Similarly, multi-step search algorithms benchmark the route-finding capabilities of their method using a single single-step model, often based on the template-based NeuralSym model,^{2,18–20,27} and evaluate the success rate of finding potential synthesis routes for molecules of interest. However, multi-step approaches do not consider the impact of alternative single-step models, a vital aspect of the search, as the route planning algorithm uses the reaction information stored in the single-step model to find synthesis routes and create alternate reaction pathways within the reaction network.

The current question remains whether contemporary single-step retrosynthesis methods are transferable to the multi-step synthesis planning domain, and their impact on multi-step synthesis planning.^{31,32} In this work, we address the transfer between single-step and multi-step methods by incorporating different contemporary single-step models within a common multi-step search algorithm to analyze the use of these models for multi-step synthesis planning. We explore the effect on performance, analyzing the relationship between contemporary single-step and multi-step performance metrics using both public and proprietary datasets of varying size and diversity. Moreover, we also focus on vital aspects such as model suitability and chemical validity of the predicted routes.

2 Methods

In this work, we develop an evaluation framework to benchmark different single-step models in multi-step synthesis planning (Fig. 1).

2.1 Evaluation scheme

2.1.1 Single-step retrosynthesis. Single-step retrosynthesis methods are evaluated using top-*n* accuracy⁷ (Table 1). The task for single-step retrosynthesis is the correct prediction of (gold-standard) reactants from the product of a known reaction. Here, we measure the percentage of target molecules for which the correct reactants are recovered within top-*n* predictions. Considering that the single-step model defines a possible maximum reaction network for a molecule of interest, published reactions are used to assess the accuracy of the single-step model since they are assumed to be chemically valid. Consequently, the assumption is that if the single-step model can recover a greater number of published reactants, then the



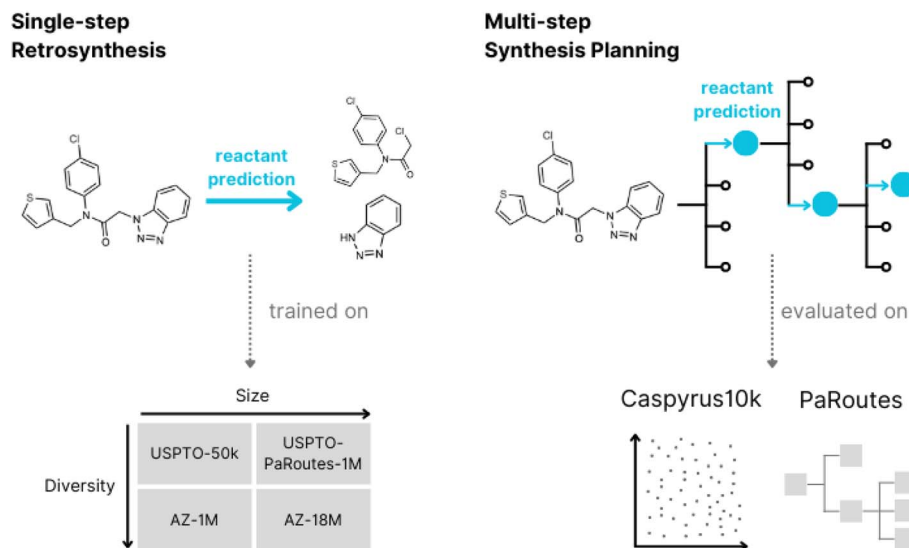


Fig. 1 Evaluation framework for single-step models (AiZynthFinder (AZF), LocalRetro, Chemformer, and MHNreact), trained on different public (USPTO-50k, USPTO-PaRoutes-1M) and proprietary (AZ-1M, AZ-18M) datasets in synthesis planning on Caspyrus10k and PaRoutes.

predictions produced by the model are chemically viable reactions.

2.1.2 Multi-step synthesis planning. On the other hand, for multi-step synthesis planning, the task is the search for likely synthesis routes for a molecule of interest, *i.e.*, a reaction pathway from the target molecule to a set of available building blocks.⁷ For this, we consider multiple aspects for both the search and the predicted routes.

Within success rate, we measure the percentage of molecules for which the route planning algorithm can successfully return at least one solved synthesis route leading from a molecule to building blocks. This condition is required for synthesis routes since a chemist can only consider routes as a suggestion for experimental evaluation if a complete synthesis route is found. Moreover, we analyze the number of solved routes since not only is it interesting to identify if there is a possible synthesis

route for a molecule, but also how many alternatives are produced, given that different synthesis routes have different route properties.

Nevertheless, algorithmic success does not measure if a found synthesis route is chemically valid, but only if a route into building blocks is found. Route accuracy is used to measure the chemical validity of synthesis routes as predicted routes can be compared to published, experimentally tested gold-standard routes.³³ Naturally, a route planning algorithm should be able to recover the gold-standard routes within the set of predicted, solved synthesis routes. This task is inherently more complex than producing solved routes (success rate) since it requires a sequence of multiple reactions and their intermediates to be correctly predicted and in the correct order. Additionally, we calculate whether there is an exact match between the predicted building blocks and the gold-standard building

Table 1 Evaluation metrics for single-step retrosynthesis and multi-step synthesis planning. Solved synthesis route implies that the produced route leads to building blocks

Task	Metric	Description
Retrosynthesis	Top- <i>N</i> accuracy	Percentage of compounds for which the ground-truth reactants are predicted within the top- <i>n</i>
Multi-step synthesis planning	Success rate	Percentage of compounds where at least one solved synthesis route is produced
	Number of solved routes	Average number of unique solved synthesis routes produced per molecule
	Search times	Average search time per molecule
	Single-step model calls	Average number of single-step model calls per molecule
	Route accuracy	Percentage of compounds where the gold-standard route is predicted within the top- <i>n</i> synthesis routes
	Building block accuracy	Percentage of compounds where the gold-standard building blocks are predicted within the top- <i>n</i> synthesis routes



blocks. Building block accuracy differs from route accuracy since the route reactions and intermediates are not considered. In all cases it must be noted that a gold standard route is only one possible way of synthesizing a target molecule.

Lastly, we consider search times and single-step model calls. Ideally, synthesis planning algorithms should produce routes in a timely manner to reduce allocated computational resources. However, different single-step models can have different inference speeds, and the time required for a search can massively diverge.³² Consequently, the average search time for a molecule with a fixed number of single-step model calls, is measured. Additionally, we report the number of single-step model calls since, in some cases, the method may not reach the maximum iteration limit in the maximum search time. Noteworthy, the maximum search time can be exceeded if the last search iteration is started before the time limit is reached.

2.2 Datasets

2.2.1 Single-step retrosynthesis. Within single-step retrosynthesis datasets, each reaction is unique. They are all curated to comprise a single product leading to one or more reactants. One product can have more than one recorded reaction, and a reaction type can occur multiple times. Here we use four different single-step retrosynthesis datasets, USPTO-50k,³⁰ USPTO-PaRoutes-1M,³³ AZ-1M and AZ-18M³⁴ (Table 2). USPTO-50k is the default benchmark dataset for single-step retrosynthesis prediction. It features 50 016 reactions comprising ten reaction classes extracted from the original USPTO dataset,²⁹ which originates from the United States Patent and Trademark Office. USPTO-PaRoutes-1M is a processed version of the original USPTO grant and application data. This single-step dataset is specifically developed to train single-step retrosynthesis models to benchmark multi-step algorithms.³³ The dataset contains single-step reactions and excludes gold-standard synthesis routes and their corresponding reactions for multi-step benchmarking. Here, we use the PaRoutes 2.0 dataset, which contains 1 198 554 single-step reactions.³⁴

Additionally, we use two datasets based on the proprietary AstraZeneca dataset.^{34,36} The first, AZ-18M, is the complete cleaned dataset from AstraZeneca, which includes Reaxys,³⁷ Pistachio (a superset of USPTO-PaRoutes-1M),³⁸ and AstraZeneca Electronic Laboratory Notebooks (ELN) data. This dataset contains 18 697 432 single-step reactions.³⁴ Moreover, to obtain a dataset representative of AZ-18M with a comparable size to USPTO-PaRoutes-1M, we randomly subsample 1M reactions from AZ-18M to produce AZ-1M.

To evaluate single-step models, we split all reaction datasets into random 80% training, 10% validation, and 10% test hold-out splits. In the case of USPTO-PaRoutes-1M, to replicate the original data split size,³³ the hold-out split ratio is 90% training, 5% validation, and 5% test. We defer from using the original hold-out splits since they are based on template stratification. For AZ-18M, we randomly subsample 100k molecules from the complete test set of 1.8 million reactions to avoid excessive evaluation computation.

2.2.2 Multi-step synthesis planning. Multi-step evaluation datasets are collections of compounds that are used to test the route-finding capabilities of multi-step synthesis planning algorithms. To evaluate the synthesis planning capabilities of different single-step models, we create a new dataset called Caspyrus10k that consists of a clustered set of 10 000 molecules from a selection of known bioactive compounds, to ensure a reasonable representation of the chemical space.

In detail, we select the high-quality papyrus³⁵ dataset of 1 238 835 molecules, combining sources such as ChEMBL³⁹ and Escape-DB,⁴⁰ where each molecule has an exact bioactivity value measure and is associated with a single protein. We filter those molecules with the Guacamol cleaning strategy⁴¹ to ensure drug-like molecules, removing molecules which do not fit the criteria in the process. As we are interested in these molecules for synthesis planning, we remove the building blocks present in Zinc,⁴² Enamine,⁴³ MolPort,⁴⁴ and eMolecules.⁴⁵ Finally, we cluster the resulting set of molecules using Butina Clustering⁴⁶ using Morgan Fingerprints with a radius of 2, a fingerprint size of 1024, and a Butina cut-off threshold of 0.6. Out of the resulting 137 963 centroids, we remove centroid molecules in clinical phases⁴⁷ to ensure we focus on molecules most relevant to the early drug discovery stages. In detail, we skip two centroids when selecting the largest 10 000 cluster centroids, which results in an evaluation dataset of 10 000 molecules representing roughly 281 000 molecules.

Additionally, we evaluate the synthesis planning capabilities of all single-step models on PaRoutes,³³ a collection of 10 000 gold-standard retrosynthesis routes. This task differs from the general synthesis planning task with Caspyrus10k in that the goal is to recover specific real-world synthesis routes conducted as part of a patent application process and therefore test the chemical validity of the predicted synthesis routes. The gold-standard routes are obtained from USPTO patent data, where we use the $n-1$ set, which contains a single retrosynthesis route for each patent. As stated in the PaRoutes dataset, we use a specialized set of building blocks containing the leaf nodes of

Table 2 Datasets for training single-step retrosynthesis models and evaluating multi-step synthesis planning

Task	Dataset	Description
Single-step retrosynthesis training	USPTO-50k ³⁰	Default single-step retrosynthesis benchmark dataset
	USPTO-PaRoutes-1M ³³	Largest publicly available single-step retrosynthesis dataset
	AZ-1M ³⁴	1M reaction subsample of internal AstraZeneca reactions
	AZ-18M ³⁴	Dataset based on internal AstraZeneca reactions
Multi-step synthesis planning evaluation	Caspyrus10k	10 000 clustered bio-active molecules from papyrus ³⁵
	PaRoutes ³³	Collection of 10 000 gold-standard synthesis routes extracted from patents



all 10 000 routes. Given the specifics of the PaRoutes dataset, the search algorithm has a maximum route length of 10 as this is the longest extracted route length from patents.

2.3 Selected approaches

2.3.1 Single-step retrosynthesis. We select three contemporary single-step methods to evaluate within multi-step synthesis planning (Table 3). The selection is based on their top-*n* accuracy on the commonly used benchmarking dataset, USPTO-50k, ensuring to select models which employ the main research directions within the field, *i.e.*, graph-based neural networks, sequence-to-sequence, and information retrieval. Where possible, we maintain the original implementation of the methods and only report deviations from this.

LocalRetro¹¹ is a template-based method that uses local atom and bond templates. It applies a graph neural network to create embeddings for both atoms and bonds of a product, which are used in a classification task to predict appropriate templates and reaction centers jointly. Contrary to the original implementation of the method, for AZ-1M and AZ-18M we filter for a minimum template frequency of three to avoid an infeasible number of local atom and bond templates.

Chemformer¹² is a template-free method based on a transformer architecture that uses BART⁴⁸ pre-training on molecular SMILES and is then fine-tuned on the retrosynthesis task. It uses product SMILES as input to predict reactant SMILES using beam-search. We set the beam size to 50.

MHNreact,¹⁰ a template-based information retrieval approach, trains separate product and template encoders and uses modern Hopfield networks⁴⁹ to relate products and template embeddings to find the most applicable reaction template. The original implementation uses all template embeddings simultaneously. However, due to large RAM requirements (>300 GB) of this approach for USPTO-PaRoutes-1M, AZ-1M and AZ-18M, the templates are used in batches to train the model. Moreover, we apply a cut-off of a minimum of three template occurrences for AZ-1M and do not show results for AZ-18M as due to increased reaction diversity leading to a much larger number of templates requiring an unfeasible amount of memory.

Additionally, we include a simple template-based model as a baseline referred to as AZF, adapted from NeuralSym,⁹ which is the default model in the most used public route planning

software implementation AiZynthFinder.⁴² Noteworthy, this model architecture is also commonly used to benchmark novel multi-step search algorithms. Templates are extracted using the standard implementation of RDChiral⁵⁰ with a radius of two. Only templates with at least three occurrences are kept for USPTO-50k, USPTO-PaRoutes-1M, and AZ-1M, for AZ-18M templates with at least ten occurrences were kept, following.³⁴

2.3.2 Multi-step synthesis planning. For multi-step synthesis planning, we select Retro*¹⁹ as the search algorithm used in all experiments. Retro* is a best-first tree search algorithm leveraging A*-like pathfinding guided by a neural network, where each algorithm iteration applies a single model call. We select Retro* as the multi-step algorithm since prior work shows minimal differences across multi-step algorithms,⁵¹ though this is only shown for the common NeuralSym model architecture. Moreover, Retro* performs better than MCTS with contemporary single-step retrosynthesis models, which require longer inference times.³² This performance difference is likely because Retro* does not require online planning for search tree traversal, limiting the number of single-step model calls required. Noteworthy, we defer from using a self-play dependent route planning algorithm, even though they have the highest reported benchmark performance²⁷ since self-play algorithms are not training data and single-step model agnostic, *i.e.*, changes in stock or single-step model change the learned self-play tree traversal policy. This aspect is especially problematic for this work as every single-step model and data combination would require self-play training such that it would become unclear whether the single-step model or the self-play aspect is important for route planning. Furthermore, we use Retro* with no cost function, such that the reactant probability of the single-step model is the guiding probability in the tree search. We defer from using the oracle function because it has shown little impact⁵¹ and is trained on USPTO data, which could cause information leakage. The search goal of Retro* is to find synthesis routes that end in building block molecules, however, that information is not used to shape the reward, as in MCTS,^{2,42} where the percentage of building block leaves is used to guide the tree search. Instead, the sole guidance of the tree search comes from the single-step model to prioritize reactions to explore. For all searches, we use a maximum search time of 8 hours (28 800 seconds) and 200 algorithm iterations. Furthermore, the top 50 reactions from the single-step model are added to the search tree at every iteration, deferring from using

Table 3 Selected single-step retrosynthesis models and multi-step synthesis planning algorithm

Task	Approach	Description
Single-step retrosynthesis	LocalRetro ¹¹	Graph neural network predicting the application of local bond and atom templates
	Chemformer ¹²	Template-free sequence-to-sequence transformer
	MHNreact ¹⁰	Template-based information retrieval method relating products and template embeddings
	AZF ^{9,42} (baseline)	Default template-based method
Multi-step synthesis planning	Retro* ¹⁹	Best-first tree search algorithm leveraging A*-like pathfinding guided by the single-step model



a cumulative probability cut-off. Moreover, unless otherwise stated, we use a maximum synthesis route length of 7 and the Zinc⁴² building block set consisting of 17 422 831 molecules.

2.4 Implementation

All single-step retrosynthesis models are incorporated into the AiZynthFinder⁴² synthesis planning framework using a newly developed common single-step model interface, ModelZoo. We extend AiZynthFinder such that any single-step model can be tested and used interchangeably within all implemented multi-step search algorithms. Where possible, the original single-step model code is used. All code and public data is available on GitHub.

2.5 Computational requirements

All single-step models for this work are trained on GPUs (Tesla V100). However, route planning is conducted on CPUs, given that insufficient GPUs are available for embarrassingly parallel evaluation of 10 000 molecules for each single-step model. In total, more than 1.5 million CPU hours were used to create the reported results.

3 Results

3.1 Single-step retrosynthesis prediction

3.1.1 USPTO-50k. As in the respective single-step retrosynthesis publications, the results on the USPTO-50k dataset, commonly used to benchmark and develop new single-step models,⁸ are reproducible. The best-performing methods are the contemporary template-based methods (LocalRetro, MHNreact), which approach over 93% accuracy by top-50 (Fig. 2). Among those methods, LocalRetro is the best performing, closely followed by MHNreact. Chemformer, a template-free method, has the highest top-1 accuracy but stagnates as its performance does not increase with rising top-*n*. AZF is the worst-performing model until the top-10, where it outperforms Chemformer. However, AZF and Chemformer only reach a maximum of 77% by top-50, an almost 19% performance drop-off compared to LocalRetro and MHNreact.

3.1.2 USPTO-PaRoutes-1M. All models perform practically identically on the USPTO-PaRoutes-1M single-step dataset, with a maximum difference of $\pm 4.6\%$ accuracy across all top-*n* (Fig. 2), despite each approach employing different model architectures. At top-1, most models perform similarly, with LocalRetro outperforming the other models by 1%. Within the top-3 accuracy, all contemporary models (LocalRetro, Chemformer, MHNreact) maintain similar performance, whereas AZF performs slightly worse. By top-50, some slight differences are present, where LocalRetro is the best performing model, followed by MHNreact and the slightly worse performing AZF and Chemformer.

3.1.3 AZ-1M. In contrast to the comparably sized USPTO-PaRoutes-1M dataset, for AZ-1M the overall performance drops across all models (Fig. 2). All three contemporary models (LocalRetro, Chemformer, MHNreact) outperform AZF on all top-*n* accuracy levels. Both contemporary template-based

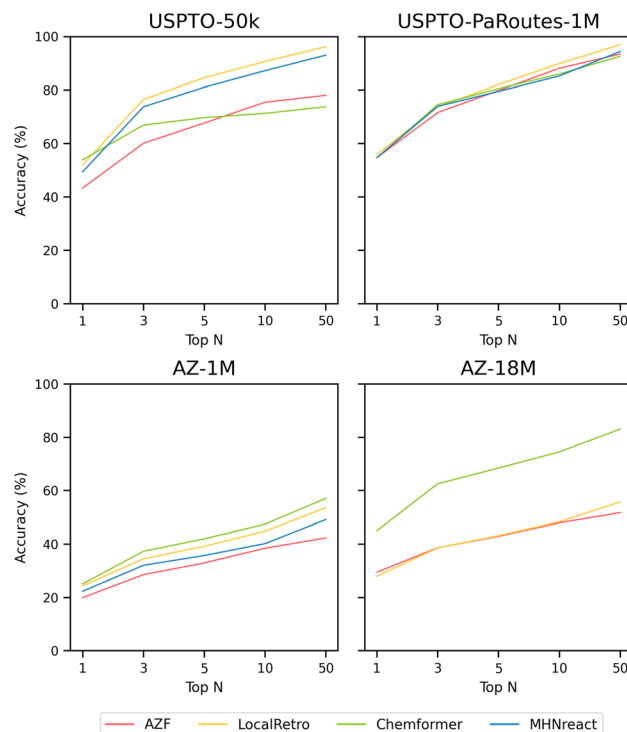


Fig. 2 Single-step retrosynthesis prediction Performance in terms of top-*n* accuracy for AZF, LocalRetro, Chemformer, and MHNreact on different datasets (USPTO-50k, USPTO-PaRoutes-1M, AZ-1M, AZ-18M) (see ESI Table S1†).

models perform similarly, where LocalRetro surpasses MHNreact as top-*n* increases. The template-free model, Chemformer, is the best-performing model throughout, though the difference is initially minimal, it becomes more pronounced across larger top-*n*. At top-50, Chemformer continues as the best-performing model, however it is closely followed by LocalRetro across all top-*n*.

3.1.4 AZ-18M. On the AZ-18M dataset, with an $18\times$ increase of data compared to AZ-1M, Chemformer clearly outperforms the other models (Fig. 2). At top-1, Chemformer already reaches an accuracy of 45.0%, improving upon the other models by at least a +15.5% margin. At top-50, Chemformer reaches 83.1%, outperforming the next best model (LocalRetro) by +27.3%. Noteworthy, both template-based methods (LocalRetro, AZF) perform similarly until top-10. Importantly, it was not possible to obtain results for MHNreact on AZ-18M due to the memory requirements of the method.

3.2 Multi-step synthesis planning

3.2.1 Caspyrus10k. Multi-step metrics of single-step models in synthesis planning are evaluated on Caspyrus10k, specifically route-finding success rate, average number of solved routes per molecule, average number of single-step model calls per molecule, and the average search time per molecule (see Methods). This establishes an overview of the capabilities of different models, trained on different datasets, across a large synthesizable chemical space.



Table 4 Multi-step synthesis planning performance on Caspyrus10k for different single-step models when trained on a diverse set of datasets. Measured by the success rate, indicating the number of molecules where a full synthesis route is found, the average number of solved routes, indicating the ability to produce synthesis route candidates, search times in seconds, and the average number of single-step model calls (see ESI Fig. S1 for distributions)

Training dataset	Model	Overall	Average per molecule		
		Success rate (%)	Solved routes	Search time (s)	Model calls
USPTO-50k	AZF	41.1	36.1	159	199
	LocalRetro	74.1	124	161	200
	Chemformer	62.4	7.37	19 051	177
	MHNreact	51.0	38.0	28 958	99
USPTO-PaRoutes-1M	AZF	66.3	83.5	163	200
	LocalRetro	86.0	324	1218	200
	Chemformer	94.1	463	28 809	147
	MHNreact	64.6	215	28 839	169
AZ-1M	AZF	73.5	124	168	200
	LocalRetro	88.1	321	465	200
	Chemformer	94.5	358	29 109	108
	MHNreact	56.0	77.0	29 116	65
AZ-18M	AZF	76.2	154	154	199
	LocalRetro	87.3	350	2736	200
	Chemformer	90.9	381	30 212	75

3.2.1.1 USPTO-50k. For models trained on the USPTO-50k dataset, LocalRetro is the best-performing model with the highest success rate and average number of solved routes. Regarding success rate, a large disparity of $\pm 32.0\%$ between the best-performing and worst-performing models is present. LocalRetro performs best, with a success rate of 74.1%, followed by Chemformer, MHNreact, and AZF, with each model decreasing in performance by around 10% from the previous one. The average number of solved routes per molecule also differs largely between the different single-step models, with the best-performing model producing almost 17 \times more solved routes than the worst-performing model. Again, LocalRetro performs best with 124 solved routes, followed by MHNreact, AZF, and Chemformer. In terms of single-step model calls, AZF, LocalRetro, and Chemformer approach the 200 model-call limit, yet there is a large disparity in search time. LocalRetro and AZF require only around 160 seconds per molecule, whereas Chemformer reaches an average search time of 5.3 hours (19 051 seconds). Lastly, despite reaching the search time limit, MHNreact has a considerably lower number of model calls.

3.2.1.2 USPTO-PaRoutes-1M. Models trained on the USPTO-PaRoutes-1M dataset have considerable performance differences in synthesis planning, even though they perform similarly on the single-step test data (Fig. 2). With the increased data volume, compared to USPTO-50k, all models solve a much larger portion of Caspyrus10k. The best-performing model in terms of success rate is Chemformer with 94.1%, followed by LocalRetro, AZF, and finally MHNreact. Overall, the average number of solved routes is high for contemporary single-step models. Chemformer finds, on average, 463 solved synthesis routes, followed by LocalRetro and MHNreact with 324 and 215, respectively. In comparison, the baseline AZF model finds only

83.5 solved routes per molecule. Concerning search time, Chemformer and MHNreact both exhaust the maximum search time, where neither reaches the maximum number of model calls. AZF is by far the fastest method, reaching 200 model calls in an average of 163 seconds. LocalRetro reaches the iteration limit within 1218 seconds on average, 7.5 \times slower than AZF but considerably faster than other contemporary models.

3.2.1.3 AZ-1M. For AZ-1M, no clear performance improvement pattern is present in comparison to USPTO-PaRoutes-1M. In terms of success rate, AZF has a +7% gain compared to USPTO-PaRoutes-1M, whereas Chemformer and LocalRetro maintain a very similar success rate. MHNreact, however, drops in route-finding success, reaching only 56.0%. The average number of solved routes slightly increases for AZF compared to USPTO-PaRoutes-1M, whereas the performance decreases by 105 routes for Chemformer and more than halves for MHNreact. LocalRetro performs comparably with a minimal decrease of 3 solved routes. Regarding search time, both Chemformer and MHNreact exhaust the maximum search times, again not reaching the maximum number of single-step model calls. In fact, both models have a particularly low number of model calls, on average carrying out 108 model calls for Chemformer and 65 model calls for MHNreact. Both LocalRetro and AZF reach the maximum iteration limit, but LocalRetro is 2.77 \times slower.

3.2.1.4 AZ-18M. Finally, the success rate of models trained on the considerably larger AZ-18M dataset is comparable to the performance on AZ-1M with no changes beyond $\pm 3.6\%$, even though the single-step performance can differ massively between both single-step datasets (Fig. 2). Compared to AZ-1M, all models produce more solved routes. Chemformer solves the most routes per molecule, followed by LocalRetro and AZF. As for the search times, Chemformer once again reaches the time limit of 8 hours, whereas LocalRetro is considerably faster on



average, beaten only by AZF. AZF and LocalRetro each reach the maximum iteration limit, whereas Chemformer only has 75 single-step model calls on average. Even though Chemformer success rate decreases, it can still produce the highest number of solved routes and the best success rate on AZ-18M.

3.2.2 PaRoutes. Instead of evaluating the general route-finding abilities of single-step retrosynthesis models, PaRoutes focuses on the ability to recover gold-standard routes given a set of molecules and their predefined target building blocks. In terms of multi-step metrics, using the same evaluation as for Caspyrus10k, all models achieve an extremely high success rate of at least 91% (Table 5). In particular, AZF, LocalRetro, and Chemformer find solutions for practically all PaRoutes compounds. The three template-based methods (AZF, MHNreact, LocalRetro) produce a similar number of solved routes per molecule ranging between 159 and 173, whereas Chemformer surpasses these with an average of 524 solved routes per molecule (Table 5 and ESI Fig. S5†). As already seen with Caspyrus10k, Chemformer and MHNreact reach the maximum search time of 8 hours without maxing out the single-step model calls. LocalRetro and AZF perform considerably faster, with AZF taking just 153 seconds on average to reach the maximum of 200 iterations.

The route accuracy of the single-step model in synthesis planning measures how often the gold-standard synthesis route is recovered for a target molecule, where the selected $n-1$ set³³ features only one retrosynthetic route per target-molecule. AZF has by far the best route accuracy overall, recovering 61.8% of gold-standard routes within its top-50 predicted synthesis routes (23.7% at top-1) (Fig. 3). Noteworthy, the performance plateaus after top-10 (at 60.7%) and with little improvement at higher top- n . Both contemporary template-based methods perform similarly across all top- n , but underperform compared to AZF by around -20% (MHNreact: 39.7%, LocalRetro: 36.1%). The template-free Chemformer model is worst-performing across all top- n , reaching only 11.9% by top-50. Noteworthy, the performance for all contemporary models improves until the top-1000 (ESI Fig. S5†), but never reaches the performance of AZF.

Considering the building block accuracy, which measures if the correct building blocks of the reference route are predicted while not considering the route reactions or intermediate molecules, considerable improvements for all models are

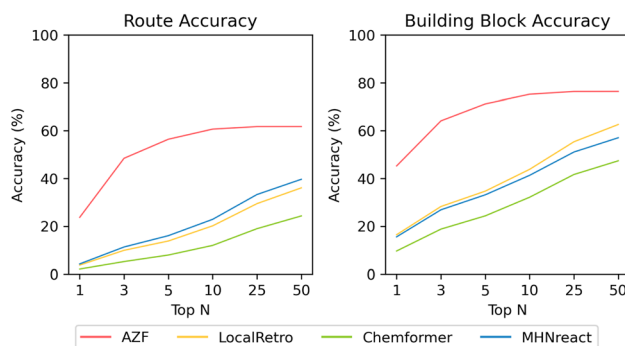


Fig. 3 Multi-step synthesis planning accuracy on PaRoutes gold-standard synthesis routes with different single-step models trained on USPTO-PaRoutes-1M. Route accuracy measures the ability to recover the correct synthesis route within top- n , whereas building block accuracy measures the ability to recover the correct building blocks while not considering reactions and intermediates (see ESI Table S7†).

present compared to the route accuracy. Within the top-50 synthesis route predictions, AZF correctly predicts the building blocks for 76.4% of the gold-standard synthesis routes, a +14.6% increase over its route accuracy. This improvement pattern is also present for the contemporary models within the top-50 predicted synthesis routes, where all three contemporary methods see a considerable improvement with at least a +17% improvement between route and building block accuracy.

4 Discussion

The task of retrosynthesis prediction is commonly treated as two separate machine learning research fields. In this work, single-step retrosynthesis and multi-step synthesis planning are joined to analyze the impact of the single-step model on multi-step synthesis planning (Fig. 1). In particular, the focus is on vital aspects of synthesis planning, the single-step model, the multi-step search algorithm, and their domain-specific applicability.

4.1 Impact on single-step retrosynthesis prediction

Considering the lack of performance transfer across single-step retrosynthesis accuracies (Fig. 2 and ESI Table S1†), it can be stated that the default single-step retrosynthesis benchmark

Table 5 Multi-step synthesis planning performance on PaRoutes for different single-step models when trained on USPTO-PaRoutes-1M. Measured by the success rate, indicating the number of molecules where a full synthesis route is found, the average number of solved routes, indicating the ability to produce synthesis route candidates, search times in seconds, and the average number of single-step model calls (see ESI Fig. S5 for distributions)

Training dataset	Model	Overall	Average per molecule		
		Success rate (%)	Solved routes	Search time (s)	Model calls
USPTO-PaRoutes-1M	AZF	97.1	159	153	200
	LocalRetro	98.9	161	1067	200
	Chemformer	99.7	524	28 538	157
	MHNreact	91.1	173	28 802	156



dataset, USPTO-50k, is problematic both in terms of performance and model scalability. A novel benchmark is required for the single-step retrosynthesis research field, as methods developed for 50 000 data points are not easily transferable to real-world-sized datasets with millions of data points. Naturally, new methods should be developed using larger datasets that better encompass the size and diversity shown in real-world data since development for USPTO-50k limits their transferability (Fig. 2). In terms of dataset size, all models require at least minor refactoring to run on larger datasets or do not scale beyond 1 million data points (MHNreact). Similarly, some USPTO-50k developed models do not conceptually consider the increase in reaction diversity in larger (real-world) datasets. For example, template-based models produce more templates with higher data diversity, requiring more template prediction classes in their classification tasks. Inherently, the number of classes a method can represent limits the number of different templates a method can predict. The solution to the diversity problem for those template-based methods is to remove templates occurring below a threshold and subsequently remove potential valid reaction predictions (see Methods). The natural exception are template-free methods as they are not constrained to reaction templates and show better scalability to more diverse data (Fig. 2).

Additionally, a model performing well on the smaller 50k reaction dataset does not necessarily perform well on larger, more diverse datasets, as the ranking of the best-performing single-step model changes for every dataset. Generally, model performance increases, or stays comparable, with more data available. For instance, for USPTO-PaRoutes-1M, a superset of USPTO-50k with a larger number of reaction classes, the performance increases (AZF and Chemformer) or stays comparable (LocalRetro, MHNreact). This pattern is also present when comparing AZ-1M to its superset AZ-18M, where more data improves the performance slightly (LocalRetro) or substantially (Chemformer, AZF). For AZ-18M, the model with the highest jump in performance is the template-free Chemformer, reaching a top-50 accuracy of 83.1% and substantially outperforming all other template-based methods by +27.3%. Here it seems that the template-based nature of the other two models (AZF, LocalRetro) limits their ability to perform on the largest, most diverse dataset. This indicates that template-based methods may have reached a performance plateau due to not being able to extrapolate beyond known templates, a limitation which is not present for the template-free Chemformer. Interestingly, for USPTO-50k, the template-free method is outperformed by all template-based methods at top-10 accuracy. Looking at the performance of AZF on AZ-18M, it is generally worse than shown in ref. 34. The previous work uses a template-based stratified split for the hold-out split, leading to an even distribution of templates across the different splits and ensuring that every template is present in every split, which can benefit a template-based approach. However, in this work, we address the hold-out split by a strict random split on the reaction level, given the nature of the different single-step methods used. With increased data diversity, single-step performance diminishes for all models comparing the equally sized USPTO-PaRoutes-1M

and AZ-1M (Fig. 2). Data diversity is measured by the number of extracted unique reaction templates from the training splits of both datasets (USPTO-PaRoutes-1M: 314 959, AZ-1M: 439 618), representing different reaction ideas present in the respective datasets. This pattern is especially problematic, as USPTO-50k only includes ten reaction classes (USPTO-50k: 10 196 unique reaction templates).

Noteworthy, USPTO-PaRoutes-1M,³⁴ with its higher number of reactions and reaction diversity, is also not a perfect single-step model benchmark dataset since all single-step models perform comparably on it. Compared to the alternative public dataset USPTO-Full,⁵² the performance of all single-step models is much higher on USPTO-PaRoutes-1M, where LocalRetro has a more than +25% top-50 accuracy improvement.⁸ The difference in single-step performance between USPTO-PaRoutes-1M and USPTO-Full and the equal performance on USPTO-PaRoutes-1M might be explainable by the underlying data sources and their respective preprocessing. USPTO-PaRoutes-1M is a superset of USPTO-Full, where the first contains USPTO grants and applications (3 748 191 total reactions) and the latter only USPTO grants (1 808 938 total reactions).³⁶ In terms of preprocessing, USPTO-Full is noisier compared to USPTO-PaRoutes-1M as the latter applies extensive data cleaning and recreates and standardizes the atom-mapping between reactions with RXNMapper.⁵³ Naturally, given that all tested single-step models perform comparably on the most cleaned, standardized, publicly available dataset, the question remains whether a saturation point in single-step performance is reached on public data.

Directly inferring multi-step synthesis planning results from single-step retrosynthesis results is not possible since single-step model performance metrics do not directly transfer to multi-step route planning success. In fact, it is necessary to evaluate the performance of respective single-step models in a multi-step framework to evaluate their synthesis planning performance. In this study, single-step models performing equally well on the USPTO-PaRoutes-1M single-step task are performing vastly differently in multi-step synthesis planning. For example, Chemformer, compared to MHNreact, has considerable differences in multi-step performance with a nearly $\pm 30\%$ higher success rate and finding double the average number of solved routes per molecule (Table 4). Moreover, LocalRetro has a roughly +20% higher success rate than AZF and finds $3.9\times$ the number of solved routes. Looking at the disparities between USPTO-50k and other datasets, LocalRetro has the highest route-finding success of single-step models trained on the USPTO-50k dataset but is not the best-performing model when trained on larger datasets. Additionally, low single-step model performance on AZ-1M still leads to high multi-step performance. Here, the high diversity of reactions in AZ-1M, compared to the equally sized USPTO-PaRoutes-1M, might be the factor for the low single-step model performance. It seems that with fewer correctly predicted reactions, it is still possible to reach high multi-step performance. This aligns with prior works showing that most molecules can be addressed with relatively few reaction templates.³⁶



4.2 Impact on multi-step synthesis planning

Though multi-step synthesis planning approaches typically do not compare different single-step models, an important finding for multi-step synthesis planning is the sheer increase in performance that can be achieved by merely switching out the single-step model, introducing novel reaction pathways to traverse the underlying reaction network (Table 4). In particular, huge success rate disparities are present within datasets, where the performance difference in finding a synthesis route between the best and worst models can be up as high as $\pm 38.5\%$ (USPTO-50k: $\pm 33.0\%$, USPTO-PaRoutes-1M: $\pm 29.5\%$, AZ-1M: $\pm 38.5\%$, AZ-18: $\pm 14.7\%$). This performance disparity pattern between the best and worst performing models trained on the same dataset is also present for the average number of solved routes per molecule, where the difference in solved routes ranges in the hundreds (USPTO-50k: ± 117 , USPTO-PaRoutes-1M: ± 380 , AZ-1M: ± 281 , AZ-18M: ± 227). The availability of more reaction data can improve the success rate of route planning up to a certain level, where the largest jump is present between USPTO-50k and the considerably larger USPTO-PaRoutes-1M. Noteworthy, public data is on par with private data in terms of multi-step success rate for Chemformer and LocalRetro which have comparable performance when trained on USPTO-PaRoutes-1M or AZ-1M. However, for AZF, public datasets perform much worse as more reaction templates are extractable from private data.³⁴ For MHNreact, private data even decreases the performance as the added complexity highly increases inference times, and only 65 single-step model calls are conducted in a generous 8 hours search window. The availability of more diverse reaction data can increase the average number of solved synthesis routes produced. Generally, we see that as reaction diversity of the single-step data increases so does the number of solved synthesis routes though eventually this performance stagnates or even worsens due to model architecture limitations. All models have either longer run times, if they reach the iteration limit, or a reduced number of single-step model calls, if they reach the time limit, reducing their potential to explore additional synthesis pathways. In the case of

LocalRetro, where the minimum reaction template occurrence is increased from USPTO-PaRoutes-1M to AZ-1M from one to three due to an infeasible number of reaction template classes in the more diverse dataset, the search times massively decrease while even improving the success rate likely due to the decreased number of reaction templates. Finally, template-based models produce their respective most solved synthesis routes using the 18 million reaction dataset, AZ-18M. Chemformer, however, achieves less solved routes compared to USPTO-PaRoutes-1M as the number of single-step model calls is halved for the largest dataset, suggesting that the inference becomes slower with more diverse data.

Even though single-step retrosynthesis models improve the performance of route planning, they are generally not tailored to multi-step search algorithms. Single-step models have slow inference times that can deny high multi-step success rates, as few single-step model calls are possible within a set time limit and can also impede ad hoc synthesis route generation. Attached to the inference problems of single-step models are the algorithmic properties of most multi-step algorithms. Though multi-step algorithms require single-step retrosynthesis models, they are generally developed to address a single molecule as a sequential next-disconnection prediction problem, with few exceptions.²⁰ Single-step models, however, are not optimized for this as they predict reactants for multiple different products simultaneously, typically in a joined GPU batch. Consequently, the combination of single-step and multi-step methods, though both thought for the task of retrosynthesis prediction, are currently not developed to be complementary to each other. Moreover, novel search algorithms, such as implementing asynchronous route planning, could have a substantial impact in this area.

4.3 Impact on domain-specific applications

Retrosynthesis prediction can be viewed as a domain-specific problem where the true objective of synthesis planning is to produce routes that can be used and tested experimentally. Given that there are multiple ways of synthesizing a molecule,

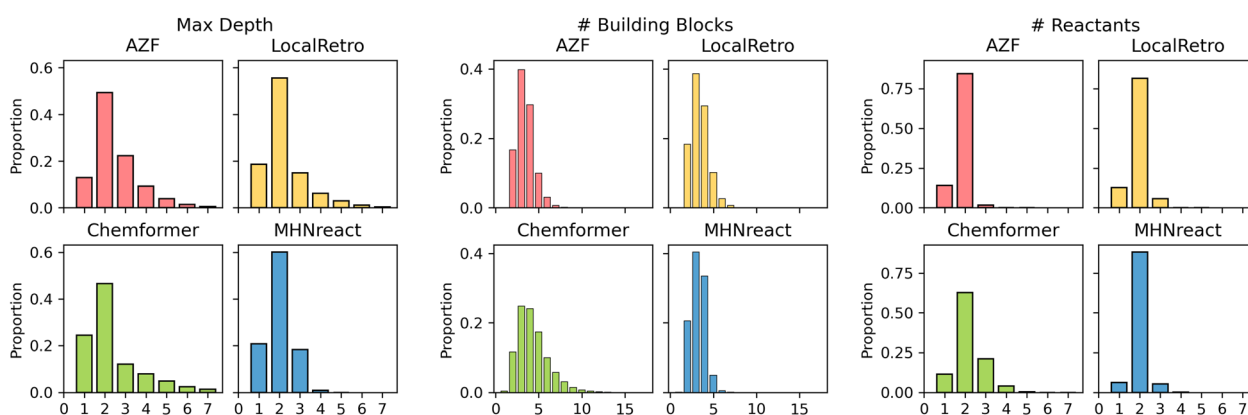


Fig. 4 Caspyrus10k route statistics of top-5 found synthesis routes by different single-step retrosynthesis models trained on USPTO-PaRoutes-1M. Shown are the maximum depth, referring to the longest linear path within the route, the number of building blocks within the route, and the number of reactants per route reaction.



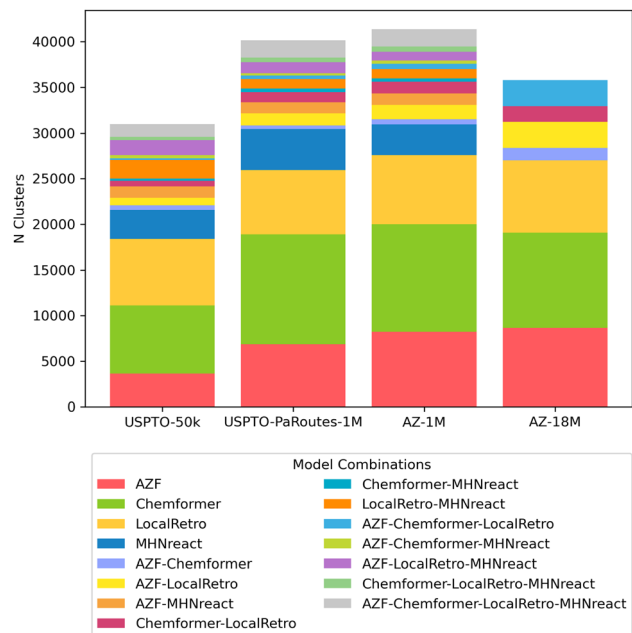


Fig. 5 Distribution and overlap of route clusters per single-step model and dataset when clustering with route-distance package.^{54,55} Clusters were calculated on a per molecule basis, N clusters shows the number of clusters which contained the stated combination of models.

the solution selected will often depend on the reaction preferences of the chemist and the desired route properties. As such, apart from the success rate and the number of solved routes, the route properties and their chemical validity are vital for the usefulness of the produced routes.

Generally, different models produce different route characteristics on Caspyrus10k (Fig. 4), where the template-free method has noticeably different maximum route length, number of building blocks and number of reactants compared to the template-based methods. AZF and LocalRetro generally have very similar distributions across all characteristics, particularly in maximum route length where MHNreact has markedly shorter routes. Since MHNreact carries out a low number of single-step model calls within the maximum search time, it is likely that it is only able to address and solve short routes. Yet, Chemformer generally has a higher proportion of routes with a maximum depth of one, essentially directly predicting building blocks. Additionally, Chemformer predicts a higher number of building blocks per route compared to all template-based methods, yet this effect is reduced with increased training data (ESI Fig. S2†). Within the template-based methods we observe that the majority of reactions are bimolecular, producing two reactants, this is particularly true for MHNreact. Chemformer on the other hand predicts reactions which at times lead to four or more reactants.

Apart from looking at general route statistics of Caspyrus10k route planning results, we cluster the resulting synthesis routes to understand the relationship between different solved routes produced by distinct models within a reaction dataset. In detail, the approximated pairwise edit distance between solved

synthesis routes of the top-5 predictions for each molecule is used to cluster with the route-distance package.^{54,55} Here, different single-step models produce unique route clusters when looking at the same training data, where routes produced by each model are generally unique to that model (Fig. 5). Noteworthy, routes produced by methods that rely on reaction templates (AZF, MHNreact, LocalRetro) tend to cluster together more frequently. Furthermore, models trained on AZ-18M tend to converge more regarding shared routes between models than models trained on USPTO-PaRoutes-1M. Nevertheless, the bulk of routes remains in unique clusters. Noteworthy, we check that the clustering patterns are also present when removing MHNreact (Fig. S3†) to ensure that the missing MHNreact results for AZ-18M are not the sole reason for the difference between AZ-18M and the other datasets.

The availability of solved synthesis routes does not imply that those routes are also chemically valid. Validity can be assessed by comparing the produced routes of a single-step model to gold-standard routes as found in USPTO patents³³ to indicate how valid the produced routes are. Generally, different single-step models are distinctive in their ability to reproduce gold-standard chemistry routes, *i.e.*, route accuracy (Fig. 3). Surprisingly, there is no relationship between the multi-step success rate and the route accuracy of a single-step model. All models achieve at least 91% success rates on PaRoutes target molecules (Table 5) but differ considerably between route accuracies. AZF is the best-performing model regarding route accuracy, recovering 23.7% of routes as the top-1 predicted synthesis route and 61.8% within the top-50 predicted routes. In comparison, contemporary models produce lower route accuracy, even if they produce high success rates. Within those contemporary models, template-based models (LocalRetro and MHNreact) have a considerably higher route accuracy than the template-free approach Chemformer, yet still have a considerable gap in performance compared to the route accuracy of AZF.

Instead of predicting the correct gold-standard synthesis route, an easier task is to predict the right building blocks of the gold-standard route. This means that though the gold-standard route may not be entirely correctly predicted the building blocks are correctly predicted in the synthesis route, *i.e.*, the order of the reactions may be incorrect or intermediate molecules are missing. For the easier task of predicting the correct building blocks, all models improve their performance compared to their respective route accuracy. However, the improvement between route accuracy and building block accuracy is much greater, compared to AZF, for contemporary models that operate on local reaction templates (LocalRetro) or no templates at all (Chemformer), potentially meaning that they are more likely to skip vital aspects of the gold-standard synthesis routes in their route predictions rather than producing a distinct retrosynthesis route than the gold-standard route. Overall, the template-based AZF method performs best regarding building block accuracy.

The performance difference on PaRoutes across different methods might be explainable by the allowed degree of chemical freedom of their respective model architectures. Template-based methods are more constrained by the reaction templates



they apply, which are extracted from training reactions. With this constraint they are made to follow reaction pathways which are more chemically sound since their templates by definition, must be based on previous reactions. In comparison, the template-free Chemformer performs worst across both route and building block accuracy, potentially explained by the non-existent template guidance of the method allowing it to predict non-chemically sound reactions. Interestingly, this is in line with the divergence of Chemformer from general route statistics on Caspyrus10k, as the model predicts a much higher number of building blocks, multi-molecular reactions and routes that only consist of a singular reaction (Fig. 4).

Generally, contemporary approaches can provide a much larger set of route alternatives (ESI Fig. S5†). This is also reflected in the PaRoutes route and building block accuracy, where AZF plateaus by top-10 accuracy, whereas contemporary methods continue to increase their accuracy into very high top- n (ESI Fig. S4†). Given that contemporary models produce more route alternatives, a future research direction, might be the best ranking of synthesis routes, as it can be assumed that desired routes are present within a large set of found synthesis routes.

An underlying assumption for single-step and multi-step synthesis planning is that the single-step model prior indicates the predicted chemical viability of a reaction for a molecule. We assess this assumed relationship by extracting the predicted reaction probabilities and their respective rank for reactions from the top-10 solved routes of the PaRoutes benchmark dataset by analyzing a random subset of 100 000 reactions for each model. Interestingly, models with a smoother

progression between probabilities of higher and lower ranked disconnections (AZF, LocalRetro, MHNreact) tend to perform better at recovering gold-standard routes (Fig. 6 and 3). In contrast, a more skewed, overconfident, distribution towards top-1 predictions tends to perform worse (Chemformer). In all cases, reactions of solved synthesis routes contain both low probability reactions and low prediction rank (ESI Fig. S6†).

Though the routes found within the top-10 predicted routes use reactions with very low reaction probabilities, gold-standard routes are generally only found within the top-5 predicted reactions (ESI Fig. S4†). This suggests that routes with reactions ranked outside the top-5 predicted reactions, though leading to building blocks, produce non-viable route reactions (Fig. 6). The presence of these low-probability reactions can be explained by the search algorithm ranking possible synthesis routes by their ability to reach building blocks and their overall route length. In the tree search itself, the search algorithm prioritizes short and solved routes, which might also include reactions with low probabilities as the overall search goal is to find a synthesis route ending in purchasable building blocks. The effect of low-probability reactions is enforced by adding 50 reactions to the search tree at every time step, even if those disconnections have low probabilities. Noteworthy, it is likely that the tree search algorithm explores those low-probability reactions when the high-probability disconnections are already explored. However, given the overall distribution of reaction priors (Fig. 6), this approach might not be desired for future synthesis planning search algorithms. Furthermore, in future work, it could be interesting to analyze how the synthesis planning results differ when applying only the top-5 predicted reactions, consequently limiting the breadth of the search tree. Given that gold-standard routes are only found within the top-5 predicted routes (ESI Fig. S6†), it opens the question if the resulting synthesis routes are closer to human-desired routes.

When discussing gold-standard synthesis routes, it is important to point out that a gold-standard route is only one way of synthesizing a desired molecule and other valid synthesis routes might also be possible. However, a good synthesis planning application should be able to prioritize real-world routes from a set of all potential routes, even if the favored chemical reactions change over time. Not finding the real-world routes entirely, yet identifying the correct building blocks, indicates that the produced synthesis routes are invalid or potentially missing vital parts of the synthesis route to be directly useful in an experimental setting. Naturally, there is a clear connection between the ability to recover gold standard routes and the ability to predict solved routes at all. High success rates produce route candidates that might be potential real-world synthesis routes but need to consider chemical validity. Because of this lack of validity, candidates are currently treated as initial retrosynthetic ideas. For a real improvement in the field of retrosynthesis, one of the essential questions, beyond improving the generation of possible solved route candidates, is how to evaluate and improve the chemical validity of generated synthesis routes. For this, it is vital to introduce reagents, conditions and yields into synthesis planning in the future and address the chemical feasibility of the generated

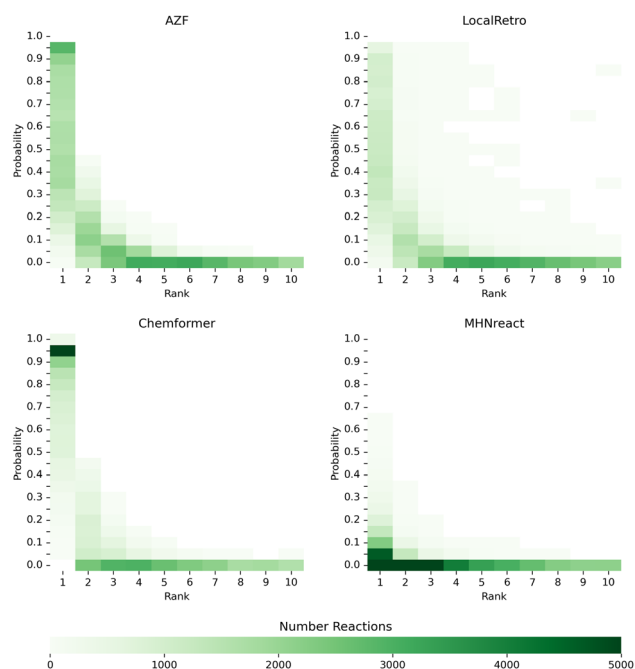


Fig. 6 Single-step model prior and rank distributions of reactions from the predicted and solved PaRoutes synthesis routes. A random sample of 100 000 reactions is extracted from the top-10 predicted routes (see Fig. 3) for each single-step retrosynthesis model trained on USPTO-PaRoutes-1M.



routes. Though there is currently a lack of *in silico* synthesis feasibility evaluation, as methods like round-trip accuracy²² only measure if the product is recoverable from the reactants and do not consider full chemical validity, given that retrosynthesis methods do not produce the relevant reagents and conditions required. Newer works have attempted to address this problem by predicting all required components.²³ Chemical validity, however, could potentially be addressed with new advancements in the field, such as molecular dynamics or quantum chemistry prediction.

Finally, when selecting the single-step retrosynthesis model for route planning, there are trade-offs between different desired search properties, as no approach outperforms all others if one uses a large enough dataset like USPTO-PaRoutes-1M. Clearly, there is a single-step performance advantage of template-free single-step models on large, heterogenous reaction data. However, this advantage comes at the cost of inference speed at multi-step synthesis planning, where template-based models are generally preferred as they can perform over 200-fold faster than template-free. If the overall goal of synthesis planning is a high success rate with a high average number of produced solved routes while accommodating long search times and a high divergence from reference routes, then the template-free approach, Chemformer, may be relevant. With a slightly lower success rate and average number of solved routes but much shorter runtimes and medium divergence from reference routes the successful contemporary template-based model, LocalRetro, is of interest. For very short run times and low divergence from reference routes yet lower success rate and an average number of solved routes, the default single-step retrosynthesis model, AZF, will be of use. Future developed models can aim to address a combination of these goals.

One of the underlying problems in the field is that benchmarking different single-step retrosynthesis models within synthesis planning is time- and resource-intensive. In order to facilitate such benchmarking in the future, we analyze the variance of different subsample sizes of the Caspyrus10k multi-step synthesis dataset such that an approximation of the results can be carried out *in lieu* of running the full datasets for faster benchmarking/prototyping (see ESI Tables S2–S5†). In detail, we repeatedly randomly subsample a subset of molecules (100, 500, 1000, 5000 molecules) and measure the mean and standard deviation across 1000 subsamples (sampling without replacement). Given that the standard deviation is reasonably small for a sample size of 1000 molecules (see ESI Table S4†), we provide a selected set of 1000 molecules if a full evaluation is not feasible (see ESI Table S6†).

Noteworthy, this work only explores three contemporary and a common baseline single-step retrosynthesis models, and even though representative of the common research directions, gives us only a snapshot of possible single-step and multi-step retrosynthesis combinations. In the future, it might be interesting to increase the evaluation framework beyond a well-chosen singular search algorithm as search algorithms are an active research field and newer approaches could improve the performance further.

5 Conclusion

In this work, we create the first in-depth study combining contemporary single-step retrosynthesis with multi-step synthesis planning, analyzing the gains and pitfalls when combining the two research fields. We find that there is no direct relationship between high single-step performance and successfully finding synthesis routes, both for publicly available and large-scale proprietary datasets. We show that, well-performing single-step models also require vital aspects for a successful search, such as fast inference times, scalability to large datasets and prioritization of reactions leading to potential synthesis routes, emphasizing the need to develop and evaluate single-step retrosynthesis models in a multi-step synthesis planning framework. Moreover, we show that the default single-step retrosynthesis benchmark dataset, USPTO-50k, is insufficient as methods developed for this small, homogenous dataset are not transferable to real-world, larger, and more diverse datasets, both in terms of single-step performance and model scalability.

For multi-step synthesis planning, we show that the single-step model is an essential but thus far underexplored aspect of the search algorithm. By merely changing the single-step retrosynthesis model it is possible to improve route-finding success by up to +28%, reaching success rates above 90% compared to the commonly used baseline model, when trained on the same reaction datasets, showing the benefit of further exploring single-step models within future multi-step frameworks. Furthermore, we show that every single-step model produces unique synthesis routes when used in multi-step synthesis planning with varying route characteristics, and each single-step model also differs in important aspects such as route-finding success, the average number of found synthesis routes, search times, and chemical validity. To summarize, we show that the combination of single-step retrosynthesis prediction and multi-step synthesis planning is a crucial aspect when developing future methods.

Data availability

All code is available on the Models Matter (<https://github.com/AlanHassen/modelsmatter>) GitHub repository, and all public data with their respective trained models are available in the Models Matter (<https://figshare.com/s/2eab4132b32229c1efc>) Figshare.

Author contributions

Paula Torren-Peraire: conceptualization, data curation, formal analysis, investigation, methodology, software, visualization, writing – original draft, writing – review & editing. Alan Kai Hassen: conceptualization, data curation, formal analysis, investigation, methodology, software, visualization, writing – original draft, writing – review & editing. Samuel Genheden: data curation, investigation, resources, supervision, writing – review & editing. Jonas Verhoeven: supervision, writing – review & editing. Djork-Arné Clevert: funding acquisition, resources,



supervision, writing – review & editing. Mike Preuss: funding acquisition, resources, supervision, writing – review & editing. Igor Tetko: conceptualization, funding acquisition, resources, supervision, writing – review & editing.

Conflicts of interest

There are no conflicts to declare.

Acknowledgements

This study was partially funded by the European Union's Horizon 2020 research and innovation program under the Marie Skłodowska-Curie Innovative Training Network European Industrial Doctorate grant agreement No. 956832 "Advanced machine learning for Innovative Drug Discovery". Parts of this work were performed using the ALICE compute resources provided by Leiden University.

Notes and references

- 1 R. S. K. Vijayan, J. Kihlberg, J. B. Cross and V. Poongavanam, *Drug Discov. Today*, 2022, **27**, 967–984.
- 2 M. H. S. Segler, M. Preuss and M. P. Waller, *Nature*, 2018, **555**, 604–610.
- 3 E. J. Corey and X.-M. Cheng, *The logic of chemical synthesis*, John Wiley & Sons, Ltd, New York, 1989.
- 4 C. W. Coley, W. H. Green and K. F. Jensen, *Acc. Chem. Res.*, 2018, **51**, 1281–1289.
- 5 C. W. Coley, D. A. Thomas, J. A. Lummiss, J. N. Jaworski, C. P. Breen, V. Schultz, T. Hart, J. S. Fishman, L. Rogers, H. Gao, R. W. Hicklin, P. P. Plehiers, J. Byington, J. S. Piotti, W. H. Green, A. John Hart, T. F. Jamison and K. F. Jensen, *Science*, 2019, **365**, eaax1566.
- 6 F. Miljković, R. Rodríguez-Pérez and J. Bajorath, *ACS Omega*, 2021, **6**, 33293–33299.
- 7 P. Schwaller, A. C. Vaucher, R. Laplaza, C. Bunne, A. Krause, C. Corminboeuf and T. Laino, *Wiley Interdiscip. Rev. Comput. Mol. Sci.*, 2022, **12**, e1604.
- 8 Z. Zhong, J. Song, Z. Feng, T. Liu, L. Jia, S. Yao, T. Hou and M. Song, Recent advances in artificial intelligence for retrosynthesis, *arXiv*, 2023, preprint, arXiv:2301.05864, DOI: [10.48550/arXiv.2301.05864](https://doi.org/10.48550/arXiv.2301.05864).
- 9 M. H. Segler and M. P. Waller, *Chem.—Eur. J.*, 2017, **23**, 5966–5971.
- 10 P. Seidl, P. Renz, N. Dyubankova, P. Neves, J. Verhoeven, J. K. Wegner, M. Segler, S. Hochreiter and G. Klambauer, *J. Chem. Inf. Model.*, 2022, **62**, 2111–2120.
- 11 S. Chen and Y. Jung, *JACS Au*, 2021, **1**, 1612–1620.
- 12 R. Irwin, S. Dimitriadis, J. He and E. J. Bjerrum, *Mach. Learn.: Sci. Technol.*, 2022, **3**, 015022.
- 13 I. V. Tetko, P. Karpov, R. Van Deursen and G. Godin, *Nat. Commun.*, 2020, **11**, 5575.
- 14 V. R. Somnath, C. Bunne, C. W. Coley, A. Krause and R. Barzilay, *Advances in Neural Information Processing Systems*, 2021.
- 15 C. Shi, M. Xu, H. Guo, M. Zhang and J. Tang, *Proceedings of the 37th International Conference on Machine Learning*, 2020, pp. 8818–8827.
- 16 X. Wang, Y. Li, J. Qiu, G. Chen, H. Liu, B. Liao, C.-Y. Hsieh and X. Yao, *Chem. Eng. J.*, 2021, **420**, 129845.
- 17 S. Szymkuć, E. P. Gajewska, T. Klucznik, K. Molga, P. Dittwald, M. Startek, M. Bajczyk and B. A. Grzybowski, *Angew. Chem., Int. Ed.*, 2016, **55**, 5904–5937.
- 18 A. Kishimoto, B. Buesser, B. Chen and A. Botea, *Advances in Neural Information Processing Systems*, 2019.
- 19 B. Chen, C. Li, H. Dai and L. Song, *Proceedings of the 37th International Conference on Machine Learning*, 2020, pp. 1608–1616.
- 20 S. Xie, R. Yan, P. Han, Y. Xia, L. Wu, C. Guo, B. Yang and T. Qin, *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, Washington DC, USA, 2022, pp. 2120–2129.
- 21 K. Lin, Y. Xu, J. Pei and L. Lai, *Chem. Sci.*, 2020, **11**, 3355–3364.
- 22 P. Schwaller, R. Petraglia, V. Zullo, V. H. Nair, R. A. Haeuselmann, R. Pisoni, C. Bekas, A. Iuliano and T. Laino, *Chem. Sci.*, 2020, **11**, 3316–3325.
- 23 D. Kreutter and J.-L. Reymond, Multistep Retrosynthesis Combining a Disconnection Aware Triple Transformer Loop with a Route Penalty Score Guided Tree Search, *chemRxiv*, 2023, preprint, DOI: [10.26434/chemrxiv-2022-8kth-v2](https://doi.org/10.26434/chemrxiv-2022-8kth-v2).
- 24 D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, Y. Chen, T. Lillicrap, F. Hui, L. Sifre, G. Van Den Driessche, T. Graepel and D. Hassabis, *Nature*, 2017, **550**, 354–359.
- 25 J. S. Schreck, C. W. Coley and K. J. M. Bishop, *ACS Cent. Sci.*, 2019, **5**, 970–981.
- 26 Y. Yu, Y. Wei, K. Kuang, Z. Huang, H. Yao and F. Wu, *Adv. Neural Inf. Process.*, 2022, 10257–10268.
- 27 G. Liu, D. Xue, S. Xie, Y. Xia, A. Tripp, K. Maziarz, M. Segler, T. Qin, Z. Zhang and T.-Y. Liu, Retrosynthetic Planning with Dual Value Networks, *arXiv*, 2023, preprint, arXiv:2301.13755, DOI: [10.48550/arXiv.2301.13755](https://doi.org/10.48550/arXiv.2301.13755).
- 28 S. Zheng, T. Zeng, C. Li, B. Chen, C. W. Coley, Y. Yang and R. Wu, *Nat. Commun.*, 2022, **13**, 3342.
- 29 D. M. Lowe, Extraction of Chemical Structures and Reactions from the Literature, PhD thesis, University of Cambridge, 2012.
- 30 N. Schneider, N. Stiefl and G. A. Landrum, *J. Chem. Inf. Model.*, 2016, **56**, 2336–2346.
- 31 H. Tu, S. Shorewala, T. Ma and V. Thost, *NeurIPS 2022 AI for Science: Progress and Promises*, 2022.
- 32 A. K. Hassen, P. Torren-Peraire, S. Genheden, J. Verhoeven, M. Preuss and I. V. Tetko, *NeurIPS 2022 AI for Science: Progress and Promises*, 2022.
- 33 S. Genheden and E. Bjerrum, *Digital Discovery*, 2022, **1**, 527–539.
- 34 S. Genheden, P.-O. Norrby and O. Engkvist, *J. Chem. Inf. Model.*, 2023, **63**, 1841–1846.



- 35 O. J. M. Béquignon, B. J. Bongers, W. Jespers, A. P. IJzerman, B. van der Water and G. J. P. van Westen, *J. Cheminf.*, 2023, **15**, 3.
- 36 A. Thakkar, T. Kogej, J.-L. Reymond, O. Engkvist and E. J. Bjerrum, *Chem. Sci.*, 2020, **11**, 154–168.
- 37 Elsevier Limited, Reaxys, 2023, <https://www.reaxys.com/>.
- 38 NextMove Software, Pistachio, 2023, <https://www.nextmovesoftware.com/pistachio.html>.
- 39 A. P. Bento, A. Gaulton, A. Hersey, L. J. Bellis, J. Chambers, M. Davies, F. A. Krüger, Y. Light, L. Mak, S. McGlinchey, M. Nowotka, G. Papadatos, R. Santos and J. P. Overington, *Nucleic Acids Res.*, 2014, **42**, D1083–D1090.
- 40 J. Sun, N. Jeliaskova, V. Chupakhin, J.-F. Golib-Dzib, O. Engkvist, L. Carlsson, J. Wegner, H. Ceulemans, I. Georgiev, V. Jeliaskov, N. Kochev, T. J. Ashby and H. Chen, *J. Cheminf.*, 2017, **9**, 17.
- 41 N. Brown, M. Fiscato, M. H. Segler and A. C. Vaucher, *J. Chem. Inf. Model.*, 2019, **59**, 1096–1108.
- 42 S. Genheden, A. Thakkar, V. Chadimová, J.-L. Reymond, O. Engkvist and E. Bjerrum, *J. Cheminf.*, 2020, **12**, 70.
- 43 Enamine Ltd, Enamine Building Blocks Catalog, 2023, <https://enamine.net/building-blocks/building-blocks-catalog>.
- 44 Molport SIA, Molport Compound Sourcing, Selling and Purchasing Platform, 2023, <https://www.molport.com/shop/index>.
- 45 eMolecules, Inc., eMolecules Chemical Building Blocks, 2023, <https://www.emolecules.com/products/building-blocks>.
- 46 D. Butina, *J. Chem. Inf. Comput. Sci.*, 1999, **39**, 747–750.
- 47 S. M. Corsello, J. A. Bittker, Z. Liu, J. Gould, P. McCarren, J. E. Hirschman, S. E. Johnston, A. Vrcic, B. Wong, M. Khan, J. Asiedu, R. Narayan, C. C. Mader, A. Subramanian and T. R. Golub, *Nat. Med.*, 2017, **23**, 405–408.
- 48 M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov and L. Zettlemoyer, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 7871–7880.
- 49 H. Ramsauer, B. Schäfl, J. Lehner, P. Seidl, M. Widrich, L. Gruber, M. Holzleitner, T. Adler, D. Kreil, M. Kopp, G. Klambauer, J. Brandstetter and S. Hochreiter, *International Conference on Learning Representations*, 2021.
- 50 C. W. Coley, W. H. Green and K. F. Jensen, *J. Chem. Inf. Model.*, 2019, **59**, 2529–2537.
- 51 A. Tripp, K. Maziarz, S. Lewis, G. Liu and M. Segler, *NeurIPS 2022 AI for Science: Progress and Promises*, 2022.
- 52 H. Dai, C. Li, C. Coley, B. Dai and L. Song, *Advances in Neural Information Processing Systems*, 2019.
- 53 P. Schwaller, B. Hoover, J.-L. Reymond, H. Strobelt and T. Laino, *Sci. Adv.*, 2021, **7**, eabe4166.
- 54 S. Genheden, O. Engkvist and E. Bjerrum, *J. Chem. Inf. Model.*, 2021, **61**, 3899–3907.
- 55 S. Genheden, O. Engkvist and E. Bjerrum, *Mach. Learn.: Sci. Technol.*, 2022, **3**, 015018.

