

# Digital Discovery

Accepted Manuscript

This article can be cited before page numbers have been issued, to do this please use: F. Duarte, S. Tanovic and E. Wieczorek, *Digital Discovery*, 2025, DOI: 10.1039/D5DD00358J.



This is an Accepted Manuscript, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this Accepted Manuscript with the edited and formatted Advance Article as soon as it is available.

You can find more information about Accepted Manuscripts in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this Accepted Manuscript or any consequences arising from the use of any information it contains.

# An exploration of dataset bias in single-step retrosynthesis prediction

Sara Tanovic,<sup>†</sup> Ewa Wieczorek,<sup>†,‡</sup> and Fernanda Duarte<sup>\*,†</sup>

<sup>†</sup>*Chemistry Research Laboratory, 12 Mansfield Road, Oxford, OX1 3TA*

<sup>‡</sup>*Alzheimer's Research UK Oxford Drug Discovery Institute, Centre for Artificial Intelligence in Precision Medicine, Centre for Medicines Discovery, Nuffield Department of Medicine, University of Oxford, Oxford OX3 7FZ*

E-mail: fernanda.duarte@chem.ox.ac.uk

## Abstract

Single-step retrosynthesis models are integral to the development of computer-aided synthesis planning (CASP) tools, leveraging past reaction data to generate new synthetic pathways. However, it remains unclear how the diversity of reactions within a training set impacts model performance. Here, we assess how dataset size and diversity, as defined using automatically extracted reaction templates, affect accuracy and reaction feasibility of three state-of-the-art architectures – template-based LocalRetro and template-free MEGAN and RootAligned. We show that increasing the diversity of the training set (from 1k to 10k templates) significantly increases top-5 round-trip accuracy while reducing top-10 accuracy, impacting prediction feasibility and recall, respectively. In contrast, increasing dataset size without increasing template diversity yields minimal performance gains for LocalRetro and MEGAN, showing that these architectures are robust even with smaller datasets. Moreover, reaction templates that are less common in the training dataset have significantly lower top-*k* accuracy than more common ones, regardless of the model architecture. Finally, we use an external data source to validate



the drastic difference between top- $k$  accuracies on seen and unseen templates, showing that there is limited capability for generalisation to novel disconnections. Our findings suggest that reaction templates can be used to describe the underlying diversity of reaction datasets and the scope of trained models, and that the task of single-step retrosynthesis suffers from a class imbalance problem.

## Introduction

Retrosynthesis is a key pillar of organic chemistry, requiring expert chemical knowledge to develop a sequence of reactions that lead to the synthesis of a target product. As research has progressed, so too has the space of possible transformations,<sup>1,2</sup> yet organic synthesis remains a bottleneck in drug discovery.<sup>3</sup> The pioneering work of Corey and Wipke<sup>4</sup> has since spawned a plethora of computer-aided synthesis planning programs,<sup>5–7</sup> in which a multi-step algorithm recursively calls on a single-step model to generate potential precursors. These single-step algorithms can be broadly categorised as template-based, where models learn to identify reaction centres and apply rules from an explicitly pre-defined library,<sup>8–11</sup> or template-free, where models learn reaction patterns implicitly from reaction SMILES<sup>12–14</sup> or molecular graphs.<sup>15–17</sup> The latter class of models is unconstrained by reaction templates and is thus expected to be able to propose novel transformations.<sup>12–14,18,19</sup>

These methods, as is the case with machine learning algorithms generally,<sup>20,21</sup> have previously been found to be sensitive to imbalanced data, often reinforcing biases rather than identifying important trends.<sup>22–24</sup> This is most clearly evidenced by template-based models, where retrosynthesis is formulated as a multi-class classification task<sup>25</sup> and thus model performance is heavily affected by the underlying distribution of the reaction templates in the training data. Within retrosynthesis, this bias manifests as preferential prediction of specific reaction classes, regioselectivities, or stereoselectivities which are better represented in the training set.<sup>22–24</sup> The widely used open-source USPTO reaction dataset,<sup>26</sup> derived from US patent data, and its subsets have been extensively used for training and model



comparison,<sup>27–29</sup> however, its underlying biases have been often overlooked during model evaluation.<sup>23</sup> Torren-Peraire *et al.* train and test multiple models on a variety of datasets, but the lack of a common test set means that results and biases cannot be directly compared.<sup>30</sup> Thakkar *et al.* investigate the impact of template library size on the performance of template-based models, but do not use template-free models and do not discuss the impacts of bias.<sup>31</sup> Thus, it is unclear how training data impacts model predictions, and what future reaction databases should look like in terms of size and diversity.<sup>24,32</sup>

Despite many works evaluating and comparing retrosynthesis models, there is little consensus on the best way to realistically evaluate extrapolation to real world scenarios.<sup>30,33</sup> Often models are trained and evaluated on a particular random split of USPTO50k,<sup>27</sup> which is itself a cleaned random subset of the USPTO database,<sup>26</sup> however this relatively small dataset cannot demonstrate how model performance would scale when trained and tested on much larger and more diverse in-house reaction libraries.<sup>30</sup> Recently, Bradshaw *et al.* have shown random splits of patent databases yield overly optimistic results, due to the similarity of reactions within the same patent or published by the same author.<sup>34</sup> Instead, they use patent- and author-based splits to simulate out-of-distribution (OOD) data and measure generalisation to reactions from unseen patents and authors, respectively. Other studies instead define generalisation as the ability to predict novel transformations defined by reaction templates.<sup>35–39</sup> However, these studies focus on how well different model architectures can generalise to new templates, but not how the underlying training data impacts generalisation.

Here, we investigate the effect that dataset size and diversity have on single-step model performance by training and testing on different subsets of a reaction database. We generate USPTO-retro, a retrosynthesis-specific dataset derived from USPTO,<sup>26</sup> analyse its diversity through local reaction templates,<sup>11</sup> and use it to train and test three established single-step architectures: LocalRetro<sup>11</sup> (template-based), MEGAN<sup>17</sup> (graph-based template-free), and RootAligned<sup>14</sup> (SMILES-based template-free). We show that top-*k* accuracy is correlated with the popularity of reaction templates in the training set for all models, regardless of



architecture, suggesting that this metric can serve as a measure of reaction diversity. Finally, we evaluate performance on external test sets extracted from the Pistachio database<sup>40</sup> to demonstrate a protocol for measuring generalisation to seen and unseen reaction templates (Figure 1A).

## Methods

**Data** Two databases are used in this work: the USPTO reaction database<sup>26</sup> for training and testing, and the commercial Pistachio reaction database<sup>40</sup> as an external test set. We apply a retrosynthesis preprocessing pipeline to both datasets based on the cleaning procedure of Gil *et al.*<sup>41</sup> while also removing reagents and uncommon local reaction templates<sup>11</sup> with less than six reactions. The Pistachio database is further filtered to ensure no overlap with the training data. This pipeline removes reagents and erroneous reactions to ensure data quality and is applicable to any reaction database. A detailed description of the data cleaning steps along with the codebase is provided in SI§S1.

This pipeline was applied to the USPTO reaction database<sup>26</sup> to generate USPTO-retro, which includes 1,103,781 atom-mapped reaction SMILES. Reaction templates were extracted using the LocalTemplate<sup>11</sup> algorithm, a modified version of RDChiral,<sup>42</sup> generating a total of 10,028 local reaction templates. This template extraction method was chosen to allow for direct comparison to the LocalRetro model. Two external test sets were created from Pistachio: Pistachio ID, containing 10k reactions with in-distribution templates seen in USPTO-retro, and Pistachio OOD, containing 10k reactions with unseen out-of-distribution templates.

**Splitting** The USPTO-retro dataset was split into training, validation, and test sets using a random 90:5:5 split, consistent with established practice in retrosynthesis studies.<sup>26–29</sup> This is referred to as the **full** split. To prevent data leakage, all reactions sharing the same product were assigned to the same subset.



To investigate the effects of dataset size and diversity, the training set was further split into 10%, 25%, and 50% subsets using two splitting strategies (Figure 1B):

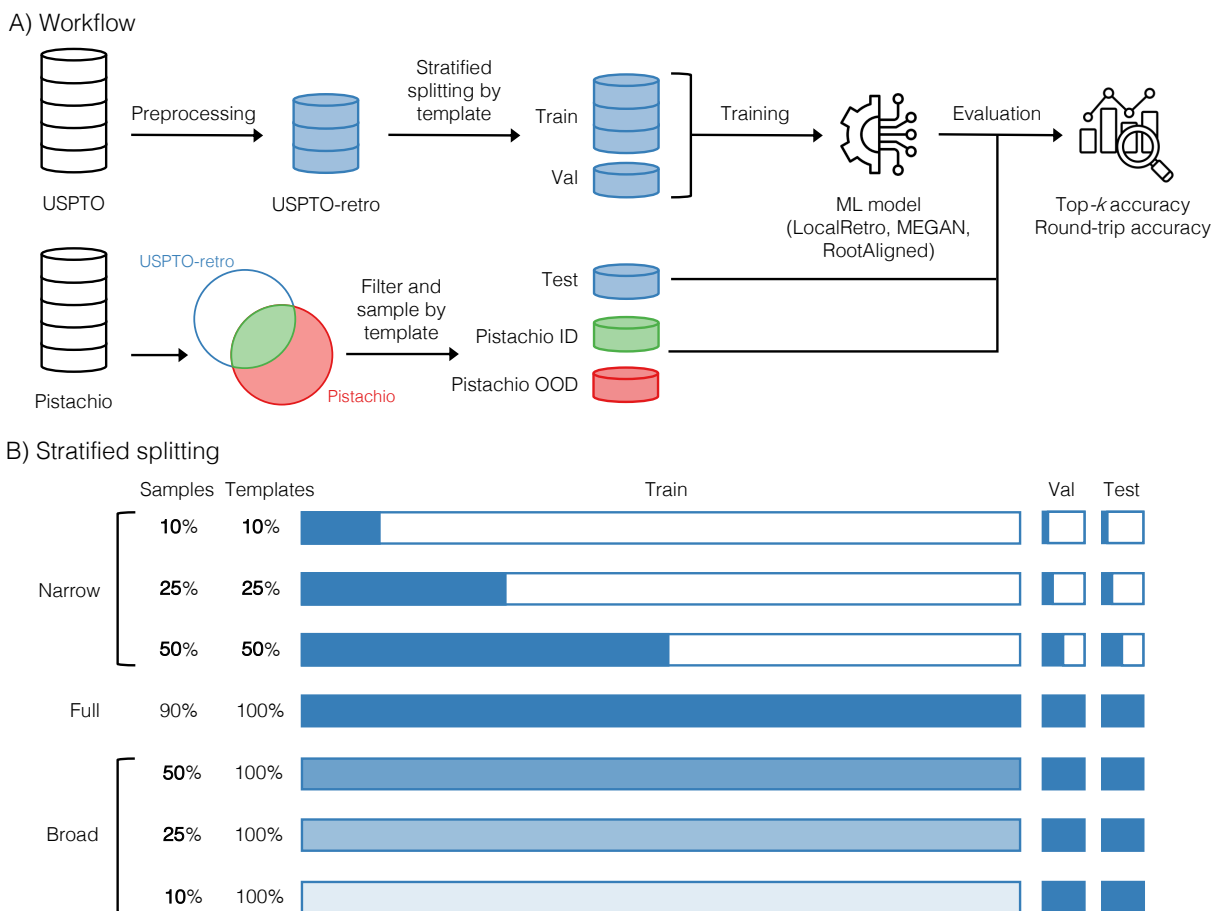
- **Narrow split:** This strategy selects a subset of reaction templates and includes all associated reactions in the training, validation, and test sets, sequentially increasing template diversity with training dataset size. The validation and test sets are similarly filtered to contain only templates seen during training. This split aims to measure how many reaction templates models can learn to predict, and the effect of increasing template diversity on model performance.
- **Broad split:** In contrast, this strategy randomly samples a fraction of reactions from all templates in the full training set while ensuring at least one example of each template is present. The validation and test sets are not altered. This split is designed to measure how much data per template is needed to learn these chemical transformations.

**Models** Three model architectures were evaluated, each representing a distinct class of retrosynthesis algorithms. (i) LocalRetro,<sup>11</sup> a template-based algorithm that learns to choose the most suitable template from an extracted list of templates; (ii) MEGAN,<sup>17</sup> a semi-template algorithm that formulates retrosynthesis as a sequence of graph edits, and (iii) RootAligned,<sup>14</sup> a template-free algorithm that treats retrosynthesis as a sequence-to-sequence translation task, translating product SMILES strings into reactant SMILES. All models were trained using their respective repositories and evaluated using the Syntheseus platform,<sup>33</sup> which automatically removes duplicate and invalid predictions.

**Evaluation** While there are many evaluation metrics available to evaluate the performance of single-step models,<sup>13,33,43</sup> here we employed top- $k$  accuracy and round-trip accuracy, which respectively measure recall and chemical feasibility.<sup>13,33,43</sup>

Top- $k$  accuracy measures the proportion of test reactions for which the ground truth reactants appear among the model's top- $k$  predictions. In this case, the ground truth is





**Figure 1:** A) Workflow of data processing, training, and testing. The USPTO-retro dataset (blue) was randomly split into training, validation, and test sets, and then further split via stratified splitting by template. Two external test sets were created from Pistachio: Pistachio ID (green), containing 10k reactions with templates seen in USPTO-retro, and Pistachio OOD (red), containing 10k reactions with unseen templates. B) Visualisation of the splitting strategies used for training and testing. The sizes of the coloured bars indicate the number of templates sampled, while the opacity represents the proportion of reactions sampled.

the reported reactants from the test set. The top-10 accuracy metric is analysed in all experiments to mimic the desired breadth of a search tree in a multi-step algorithm.<sup>33</sup>

Top-*k* round-trip accuracy evaluates the proportion of top-*k* predicted reactants that satisfy back-translation.<sup>13</sup> This is done by checking whether they regenerate the original product via a forward reaction model (here RootAligned trained on the full USPTO-retro training set) to predict the top-1 product from each set of predicted reactants. If the predicted product matches the original target, the prediction is considered successful. We



report top-1 and top-5 round-trip accuracy metrics to estimate the chemical feasibility of the top predictions.<sup>13</sup> It is important to note that the calculation of round-trip accuracy requires the use of a forward prediction model and is thus not 100% accurate, and should be interpreted as an approximation rather than an absolute measure of chemical validity.

## Results and Discussion

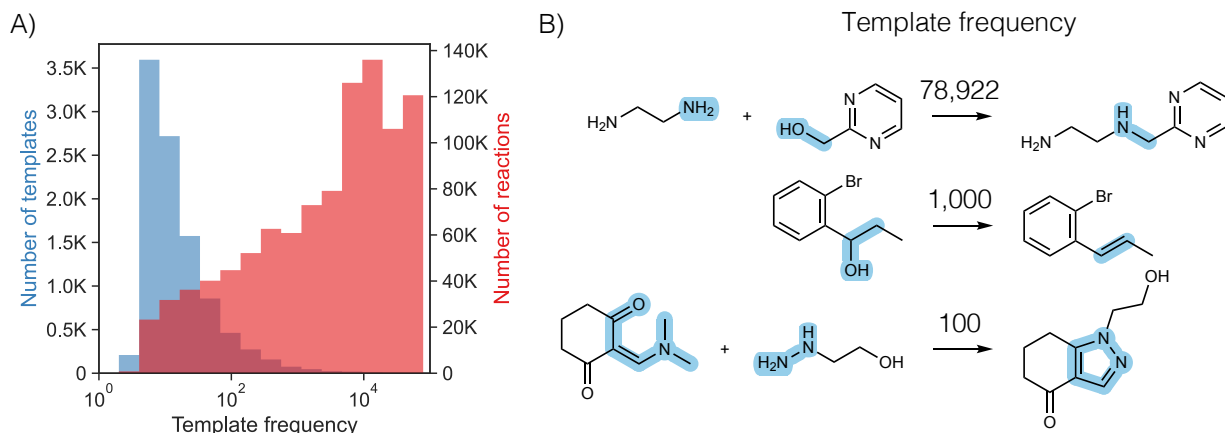
### Data analysis

We started our study by analysing the distribution of reaction templates within the newly generated USPTO-retro dataset, extracted using LocalTemplate.<sup>11</sup> Despite USPTO-retro containing over 1 million atom-mapped reaction SMILES, it shows a significant bias towards a small percentage of templates. Template frequency is used here to quantify the number of reactions a template describes in the training set, and, by extension, reaction classes (Figure 2A). The frequency of a template ranges from 2 to 78,922, with 50% of templates occurring fewer than 12 times. This bias underscores the inherent nature of open-source reaction databases, where certain reactions dominate. For example, the top 10 templates account for just 0.1% of all templates and together describe 30% of the training data.

The most common reaction template, an example of which is shown in Figure 2B, corresponds to a C-N bond-forming S<sub>N</sub>2 reaction, which accounts for >78k (8%) of all reactions in the training set. This template is similar to the next two most popular templates, which differ only in their leaving groups. Conversely, rarer templates include those with uncommon leaving groups or highly specific reaction centres. While these reactions are less common in the dataset, they are not necessarily less effective or harder to apply experimentally. Therefore, understanding the implications of this template imbalance on model performance is key for formulating better training and data curation strategies.







**Figure 2:** A) Histogram of templates (blue) and reactions (red) in the training set grouped by template frequency (on a log scale and with a box width of 0.3). Template frequency refers to the number of reactions in the training set described by a specific template. B) Example reactions from the training set with the template highlighted in blue and the template frequency labelled.

## Impact of template distribution on model performance

To evaluate the impact of template distribution on model performance, we employed two splitting strategies to further partition the training set beyond the initial random split: the narrow and broad split. Both strategies sequentially increase the size of the training data, but differ in the diversity and distributions of their templates. The narrow split increases the number of unique reaction templates in the training set as its size grows, allowing us to isolate the effect of increasing template diversity. In contrast, the broad split maintains template diversity while increasing the number of training examples, allowing us to assess the effect of increasing data volume per template. We analyse the resulting performances from these two strategies in the following subsections.

### Narrow split

The narrow split is designed to evaluate how increasingly template-diverse datasets affect model performance. As expected,<sup>25,31</sup> models trained on less diverse datasets achieve higher top-*k* accuracy, as they have fewer competing reactions to choose from (Figure 3A). Increasing



the number of templates from 1k to 10k results in a decrease in top-10 accuracy of 11.6% for LocalRetro, 14.5% for MEGAN, and 10.4% for RootAligned.

This decrease in top- $k$  accuracy does not imply lower reaction feasibility; rather, it indicates the model's increased vocabulary of reactivity as a broader set of plausible reactions is suggested. Round-trip accuracy is used here to estimate the feasibility of the predicted reactions.<sup>13</sup> The top-1 round-trip accuracy remains roughly consistent across all splits and models, with over 89% of top predictions likely to be feasible reactions. In contrast, the top-5 round-trip accuracy increases by 14-21% across all models as template diversity increases, suggesting that lower-ranked predictions become more feasible when the model is exposed to more reaction types.

This behaviour differs from previous studies wherein top- $k$  accuracy improves with additional randomly split training data.<sup>25,44</sup> In our case, increasing both the volume and diversity of training data leads to a decrease in top- $k$  accuracy. This highlights the importance of explicitly reporting and accounting for reaction template diversity when comparing model performance across datasets with varying levels of diversity.

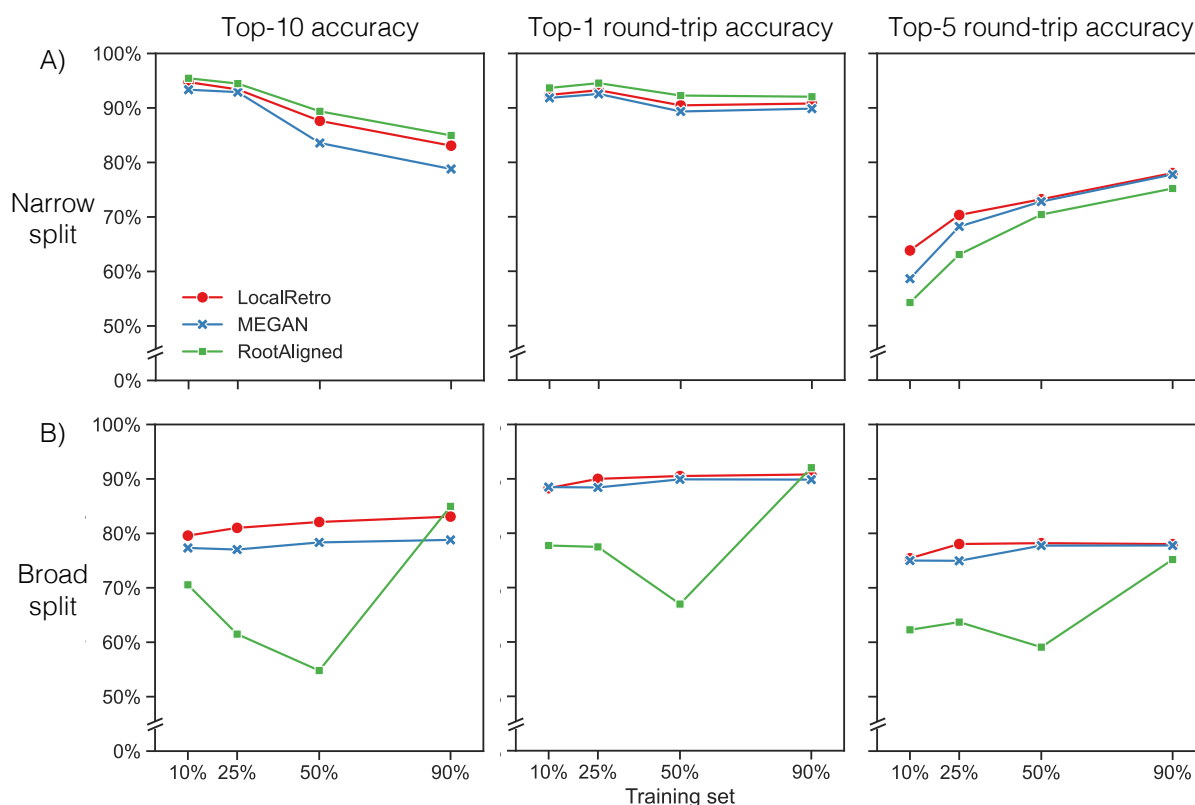
### Broad split

The broad split aims to model the effect of increasing training set size while maintaining reaction diversity by using all available templates. Our results show that performance slightly improves for LocalRetro and MEGAN, with top-10 accuracy increasing by 3.5% for LocalRetro and 1.8% for MEGAN with a ninefold increase in training set size (Figure 3B). These results suggest that, with sufficient reaction diversity, these models are robust against variations in the size of the training set.

In contrast, the RootAligned model exhibits a substantial decrease in performance across the broad split. Its top-10 accuracy degrades by 15.7% between the 10% to 50% training sets, but recovers to 85.0% with the full training set. The consistent performances of LocalRetro and MEGAN indicate that the variations observed for RootAligned arise from the underlying

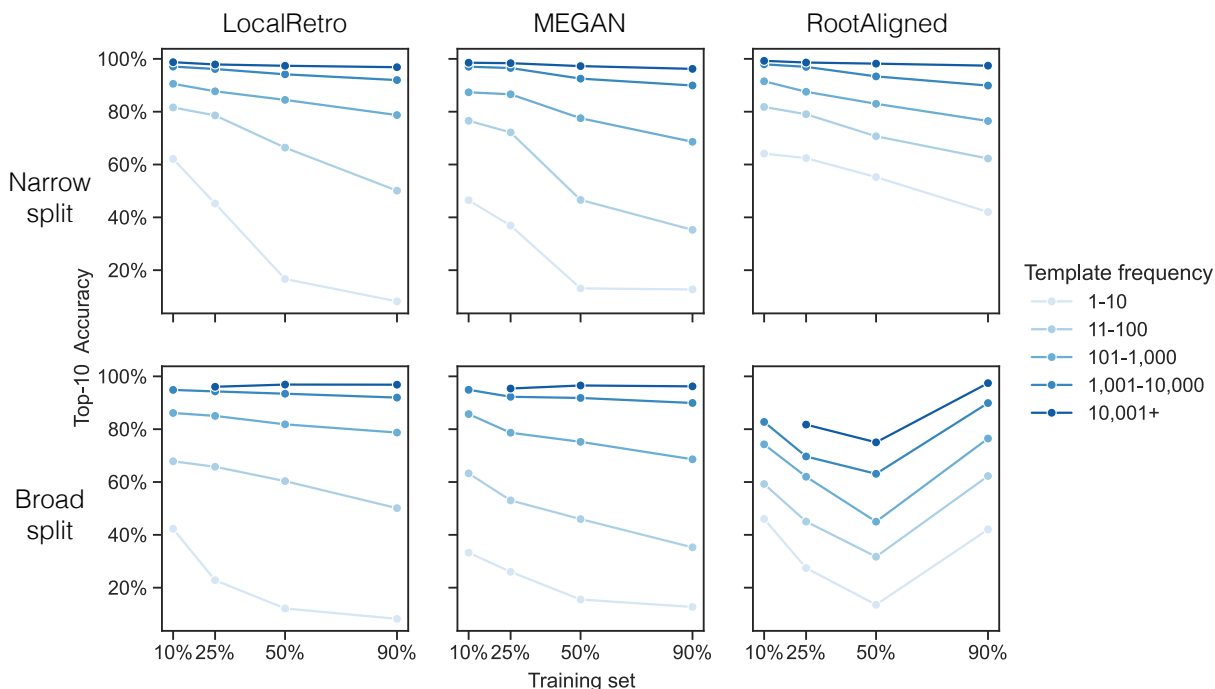


transformer architecture rather than the size or nature of these training sets. This template-free approach attempts to implicitly learn chemistry directly from SMILES strings, whereas the template-based and semi-template methods provide a more structured way of learning reactions through predefined templates and graph edits. Consequently, the learning process of the RootAligned model may require more examples of the same reactions to fully utilise this chemistry. Models may also be more easily overfit on the smaller training datasets, leading to memorisation and pattern matching, which cannot generalise to the test set. Further investigation is needed to determine if this behaviour occurs with other template-free models.



**Figure 3:** Top-10 accuracy (left), top-1 round-trip accuracy (middle) and top-5 round-trip accuracy (right) of models trained on the (A) narrow (increasing template diversity) and (B) broad splits (increasing data volume).





**Figure 4:** Top-10 accuracy of all trained models, as grouped by template frequency in the training set. The template frequency measures the number of times a particular template appears in the training set.

## Accuracy by template

Next, we investigated how template frequency bias in the training data affects model performance, focusing on top-10 accuracy across reaction templates (Figure 4). A clear trend emerges: templates that appear more frequently in the training set are predicted with significantly higher accuracy. The difference in top-10 accuracy between rare templates (frequency of 1-10) and popular templates (frequency of 10,001+) is at most 88.6% for LocalRetro, 83.5% for MEGAN, and 55.4% for RootAligned. A similar, though weaker, correlation is observed when considering Tanimoto similarities between the training and test sets (Figure S4). These trends persist even in models that do not explicitly use reaction templates, such as MEGAN and RootAligned, implying that template frequency reflects the underlying class distribution of reaction data.

In both the narrow and broad splits, increasing the training set size amplifies the spread of top- $k$  accuracies across template frequencies. For the most frequent templates (with



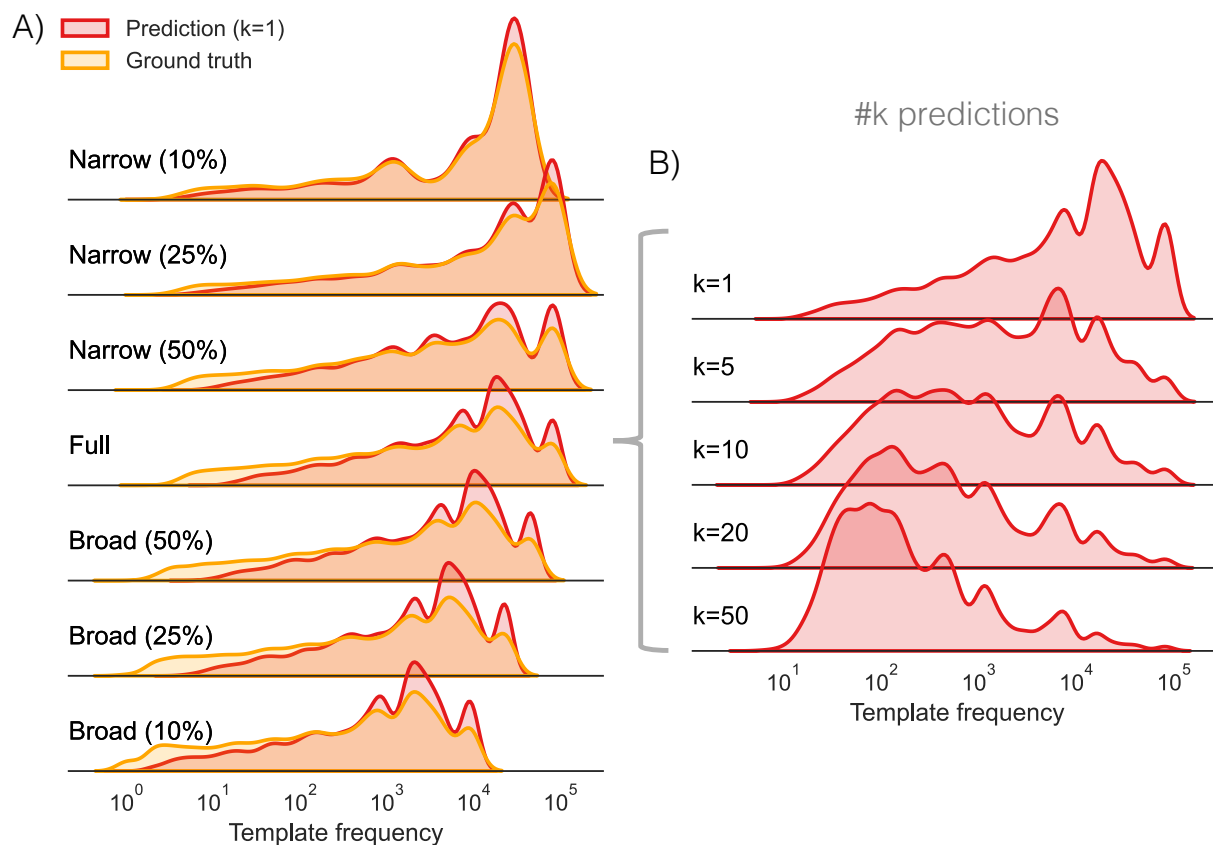
frequency  $> 10,001$ ), LocalRetro and MEGAN consistently achieve top-10 accuracy above 95%, regardless of training set sizes. In contrast, rare templates (with frequency 1-10) show a marked drop in accuracy as training set size increases: top-10 accuracy decreases between the narrow 10% and full 90% training sets by 53.9% for LocalRetro, 33.8% for MEGAN, and 22.1% for RootAligned. This behaviour is most pronounced for LocalRetro, which explicitly considers reaction templates and thus learns to prioritise more frequent classes during training. RootAligned, which implicitly encodes chemistry through SMILES strings, is less sensitive to these class imbalances. These results suggest that increasing both the number and imbalance of reaction templates contributes to performance disparities. To mitigate this, further work is needed to incorporate class balancing strategies during model training.

While the top- $k$  accuracy measures how often a reaction template is correctly predicted, it does not describe how often that type of template is recalled. Thus, it is also important to understand if the models are oversampling from popular reaction classes as a way of mimicking the training set distribution. This behaviour is most easily studied in the LocalRetro model, as its algorithm readily outputs a ranked list of predicted templates. In all splits, the model oversamples the most popular template classes for its highest ranked prediction (Figure 5A). Rarer templates are undersampled compared to the true test distribution, which contributes to their low top-10 accuracy. These rarer templates are instead sampled more often at lower ranks as the model is less confident in their prediction (Figure 5B).

## Generalisation to novel reactions

Generalisability in single-step retrosynthesis refers to a model's predictive capability for novel reactions. This can be assessed in multiple ways, for example, considering the prediction of previously unseen target products using known reaction templates or the prediction of novel disconnections not encountered during training. To systematically evaluate both aspects, we split our external test set from the Pistachio database into Pistachio ID (In-Distribution), which contains novel products with seen templates, and Pistachio OOD (Out-Of-Distribution),



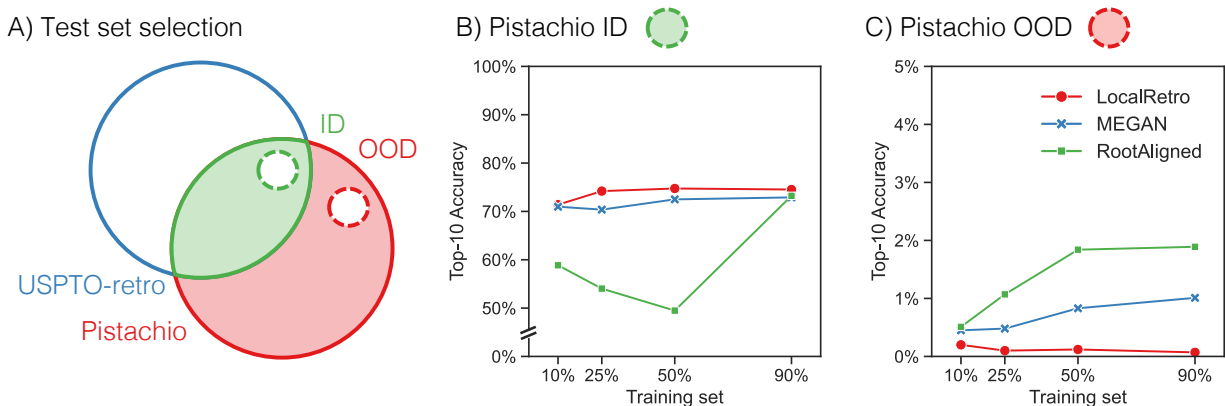


**Figure 5:** (A) Kernel density estimations (KDEs) of the training template frequency of the top prediction from LocalRetro (red) and ground truth (yellow). (B) The KDE distributions of the training template frequency of the #1, #5, #10, #20, and #50 predictions from the LocalRetro model trained on the full training set.

which contains novel products with unseen templates (Figure 6A). We use the broad split to evaluate generalisation to novel products (ID) and the narrow split to evaluate generalisation to novel disconnections (OOD).

On the Pistachio ID test set (Figure 6B), all models exhibit a moderate decline in top-10 accuracies when compared to their performance on the USPTO-retro test set (Section 3.3): 7-9% for LocalRetro, 6% for MEGAN, and 5-12% for RootAligned. This indicates that models successfully generalise to novel products using templates learnt during training, with similar performance trends to previous results. The slightly reduced performance on this test set is likely due to the lower structural similarity between Pistachio ID products and those in the USPTO-retro training sets (Figure S5).





**Figure 6:** A) Diagrammatic representation of the overlap of templates between the USPTO-retro and Pistachio datasets. The Pistachio ID test set is selected from in-distribution templates (the intersection area shown in green), whereas the Pistachio OOD test set is selected from out-of-distribution templates (the exclusive area shown in red). B) Top-10 accuracy of all models trained on the broad split and tested on the Pistachio ID test set. C) Top-10 accuracy of all models trained on the narrow split and tested on the Pistachio OOD test set.

In contrast, performance on the Pistachio OOD test set (Figure 6C) reveals severe limitations in generalisability to novel disconnections, in agreement with previous findings.<sup>35–38</sup> LocalRetro exhibits near-zero top-10 accuracy, which is expected given its reliance on predefined templates. The non-zero accuracy suggests template ambiguity, where different templates from the training and OOD test sets occasionally yield the same sets of reactants. This occurs due to overlapping SMARTS patterns or errors in atom mapping. MEGAN and RootAligned models show modest generalisability, which increases with increased training diversity and peaks at top-10 accuracies of 1% and 2% respectively with the full training sets. Their low but non-zero accuracy implies that models prioritise recognising and applying patterns seen in the training data over utilising underlying chemical principles to generate novel, feasible disconnections.

These results highlight the differences in capabilities between ID and OOD generalisation, emphasising the need for distinct evaluations that distinguish between these two scenarios. Previous studies showing the traditional learning pattern of increasing top-*k* accuracy with increasing training data volume<sup>25,44</sup> may, in fact, be misattributing the effect of additional



template coverage of the test set to additional data. This explanation may also apply to studies showing low generalisability to external datasets<sup>33</sup> or author-/patent-based splits,<sup>34</sup> wherein their test sets possibly contain both seen and unseen templates. Furthermore, the extremely low generalisability of template-free models to novel templates suggests that these models are not yet sufficiently developed to warrant their use for predicting new chemistries.

## Conclusion and future work

In this study, we presented a comprehensive assessment of the accuracy and feasibility of three established single-step retrosynthesis models – template-based LocalRetro and template-free MEGAN and RootAligned – exploring how dataset size and diversity, defined in terms of local reaction templates, affect performance.

Our results have highlighted the critical role of training set diversity in model performance. Increasing the diversity of the training set significantly increases top-5 round-trip accuracy, an indicator of prediction feasibility, while reducing top-10 accuracy, reflecting the ability of the model to recover the ground truth. This trade-off suggests that more diverse datasets enable the prediction of a broader range of plausible reactions, even if they differ from the ground truth. Interestingly, increasing dataset size without increasing template diversity yields minimal performance gains for LocalRetro and MEGAN models, suggesting that template diversity has a greater impact on model performance than volume.

We also examined the impact of template frequency on model performance. All three models, regardless of whether they explicitly use templates, show a strong correlation between a template's frequency in the training set and the model's ability to predict it correctly. This indicates that all models implicitly rely on the distribution of reaction templates learnt during training, with rare templates consistently underperforming compared to more frequent ones.

Finally, to assess real-world applicability, we evaluated model performance on two external test sets derived from the Pistachio database: one containing novel products with known





templates (Pistachio in-distribution (ID)) and another with novel products and unseen templates (Pistachio out-of-distribution (OOD)). While all models generalised reasonably well to new molecules involving known templates, their ability to predict novel disconnections was limited. These results highlight the differences in capabilities between ID and OOD generalisation. LocalRetro failed almost entirely on OOD reactions due to its reliance on predefined templates, while MEGAN and RootAligned achieved only 1–2% top-10 accuracy. These results highlight the need for evaluation protocols that clearly distinguish between in- and out-of-distribution generalisation.

These results also offer a new perspective on recent advances in transfer learning for retrosynthesis prediction, wherein fine-tuning effectively modifies the training template distribution. For instance, our reported mixed fine-tuning approach to bias predictions towards heterocyclic ring disconnections can be viewed as addressing the underlying class imbalance issues present in the initial training set.<sup>45</sup> Our results suggest that similar systematic approaches to class imbalance during training could improve representation across reaction classes. Similar challenges have been addressed in other domains, such as computer vision, through pre-training, data augmentation, and re-weighting strategies,<sup>46</sup> and could be applied to retrosynthesis through the selective augmentation of rare templates or lower weighting of popular templates during the training process.

The performance trends across the narrow and broad splits raise questions about what data should be used to train retrosynthesis models. Ideally, models would learn underlying physical principles to propose feasible reactions; however, evaluation shows that they are more likely to learn to mimic the template distribution of the training set. Further cheminformatic analysis is needed to characterise the biases of common reaction datasets and identify areas for improvement. Furthermore, models do not necessarily exhibit worse accuracy when trained on less data; therefore, data curation efforts should prioritise quality and diversity over quantity. As such, it is clear that as chemists we cannot blindly train models with all available data and not consider the types of chemistry that data represents, and whether



that chemistry suits our synthetic goals and targets.

## Availability of data and materials

**Preprocessing:** The open-source code used to split the training data is provided at <https://github.com/duartegroup/template-splits> (DOI: 10.5281/zenodo.17858529), which also contains a link to download the raw USPTO dataset. The final train, validation, and test splits are provided on FigShare (DOI: 10.6084/m9.figshare.30823988). The proprietary Pistachio dataset (licensed by NextMove Software) is not provided.

**Training:** The open-source packages used to train the machine learning models (using the configuration files provided at <https://github.com/duartegroup/template-splits/tree/main/configs>) can be found at:

- LocalRetro: <https://github.com/kaist-amsg/LocalRetro> (since removed by the authors)
- MEGAN: <https://github.com/molecule-one/megan>
- RootAligned: <https://github.com/otori-bird/retrosynthesis>

**Testing:** The open-source syntheseus package used to analyse the trained models can be found at <https://github.com/microsoft/syntheseus/tree/main>.

## Authors' contributions

ST and EW conceptualised the study. ST carried out all experiments, including the development of a new pipeline. All authors participated in data analyses. ST wrote the first draft of the manuscript, and all authors contributed to further writing. FD supervised the study. All authors approve the final version of the manuscript.



## Conflicts of interest

There are no conflicts to declare.

## Acknowledgements

This work was supported by funding from the Engineering and Physical Sciences Research Council (EPSRC) [grant numbers EP/S024093/1 and EP/Y52878X/1]. We acknowledge useful discussions with Marwin Segler, Krzysztof Maziarz, Jaron Cohen, Veronika Juraskova, and Tom Watts.

## References

- (1) Szymkuć, S.; Badowski, T.; Grzybowski, B. A. Is Organic Chemistry Really Growing Exponentially? *Angew. Chem. Int. Ed.* **2021**, *60*, 26226–26232, Publisher: John Wiley and Sons Inc.
- (2) Brown, D. G.; Boström, J. Analysis of past and present synthetic methodologies on medicinal chemistry: where have all the new reactions gone? *J. Med. Chem.* **2016**, *59*, 4443–4458, Publisher: American Chemical Society.
- (3) Blakemore, D. C.; Castro, L.; Churcher, I.; Rees, D. C.; Thomas, A. W.; Wilson, D. M.; Wood, A. Organic synthesis provides opportunities to transform drug discovery. *Nat. Chem.* **2018**, *10*, 383–394, Publisher: Nature Publishing Group.
- (4) Corey, E. J.; Wipke, T. W. Computer-Assisted Design of Complex Organic Syntheses. *Science* **1969**, *166*, 178–192.
- (5) Zhong, Z.; Song, J.; Feng, Z.; Liu, T.; Jia, L.; Yao, S.; Hou, T.; Song, M. Recent advances in deep learning for retrosynthesis. *WIREs Comput. Mol. Sci.* **2024**, *14*, e1694.



- (6) Coley, C. W.; Green, W. H.; Jensen, K. F. Machine Learning in Computer-Aided Synthesis Planning. *Acc. Chem. Res.* **2018**, *51*, 1281–1289, Publisher: American Chemical Society.
- (7) Jiang, Y.; Yu, Y.; Kong, M.; Mei, Y.; Yuan, L.; Huang, Z.; Kuang, K.; Wang, Z.; Yao, H.; Zou, J.; Coley, C. W.; Wei, Y. Artificial intelligence for retrosynthesis prediction. *Engineering-london.* **2023**, *25*, 32–50, Publisher: Elsevier Ltd.
- (8) Strieth-Kalthoff, F.; Szymkuć, S.; Molga, K.; Aspuru-Guzik, A.; Glorius, F.; Grzybowski, B. A. Artificial intelligence for retrosynthetic planning needs both data and expert knowledge. *J. Am. Chem. Soc.* **2024**, *146*, 11005–11017.
- (9) Shields, J. D.; Howells, R.; Lamont, G.; Leilei, Y.; Madin, A.; Reimann, C. E.; Rezaei, H.; Reuillon, T.; Smith, B.; Thomson, C.; Zheng, Y.; Ziegler, R. E. AiZynth impact on medicinal chemistry practice at AstraZeneca. *RSC Med, Chem*, **2024**, *15*, 1085–1095.
- (10) Seidl, P.; Renz, P.; Dyubankova, N.; Neves, P.; Verhoeven, J.; Wegner, J. K.; Segler, M.; Hochreiter, S.; Klambauer, G. Improving few- and zero-shot reaction template prediction using modern hopfield networks. *J. Chem. Inf. Model.* **2022**, *62*, 2111–2120.
- (11) Chen, S.; Jung, Y. Deep Retrosynthetic Reaction Prediction using Local Reactivity and Global Attention. *JACS Au* **2021**, *1*, 1612–1620.
- (12) Irwin, R.; Dimitriadis, S.; He, J.; Bjerrum, E. J. Chemformer: a pre-trained transformer for computational chemistry. *Mach. Learn.: Sci. Technol.* **2022**, *3*, 15022, Publisher: Institute of Physics.
- (13) Schwaller, P.; Petraglia, R.; Zullo, V.; Nair, V. H.; Haeuselmann, R. A.; Pisoni, R.; Bekas, C.; Iuliano, A.; Laino, T. Predicting retrosynthetic pathways using transformer-based models and a hyper-graph exploration strategy. *Chem. Sci.* **2020**, *11*, 3316–3325, Publisher: Royal Society of Chemistry.



- (14) Zhong, Z.; Song, J.; Feng, Z.; Liu, T.; Jia, L.; Yao, S.; Wu, M.; Hou, T.; Song, M. Root-aligned SMILES: a tight representation for chemical reaction prediction. *Chem. Sci.* **2022**, *13*, 9023–9034, arXiv: 2203.11444 Publisher: Royal Society of Chemistry.
- (15) Tu, Z.; Coley, C. W. Permutation Invariant Graph-to-Sequence Model for Template-Free Retrosynthesis and Reaction Prediction. *J. Chem. Inf. Model.* **2022**, *62*, 3503–3513, arXiv: 2110.09681 Publisher: American Chemical Society.
- (16) Somnath, V. R.; Bunne, C.; Coley, C. W.; Krause, A.; Barzilay, R. Learning Graph Models for Retrosynthesis Prediction. 35th Conference on Neural Information Processing Systems. 2021.
- (17) Sacha, M.; Błaż, M.; Byrski, P.; Dąbrowski-Tumański, P.; Chromiński, M.; Loska, R.; Włodarczyk-Pruszyński, P.; Jastrzębski, S. Molecule Edit Graph Attention Network: Modeling Chemical Reactions as Sequences of Graph Edits. *J. Chem. Inf. Model.* **2021**, *61*, 3273–3284, arXiv: 2006.15426 Publisher: American Chemical Society.
- (18) Schwaller, P.; Laino, T.; Gaudin, T.; Bolgar, P.; Hunter, C. A.; Bekas, C.; Lee, A. A. Molecular Transformer: A Model for Uncertainty-Calibrated Chemical Reaction Prediction. *ACS Cent. Sci.* **2019**, *5*, 1572–1583, arXiv: 1811.02633 Publisher: American Chemical Society.
- (19) Schwaller, P.; Gaudin, T.; Lányi, D.; Bekas, C.; Laino, T. "Found in Translation": predicting outcomes of complex organic chemistry reactions using neural sequence-to-sequence models. *Chem. Sci.* **2018**, *9*, 6091–6098, arXiv: 1711.04810 Publisher: Royal Society of Chemistry.
- (20) van Giffen, B.; Herhausen, D.; Fahse, T. Overcoming the pitfalls and perils of algorithms: a classification of machine learning biases and mitigation methods. *J. Bus. Res.* **2022**, *144*, 93–106.



- (21) Wang, R.; Chaudhari, P.; Davatzikos, C. Bias in machine learning models can be significantly mitigated by careful training: evidence from neuroimaging studies. *Proc. Natl. Acad. Sci.* **2023**, *120*, e2211613120, Publisher: Proceedings of the National Academy of Sciences.
- (22) Thakkar, A.; Vaucher, A. C.; Byekwaso, A.; Schwaller, P.; Toniato, A.; Laino, T. Unbiasing Retrosynthesis Language Models with Disconnection Prompts. *ACS Cent. Sci.* **2023**, *9*, 1488–1498, Publisher: American Chemical Society.
- (23) Kovács, D. P.; McCorkindale, W.; Lee, A. A. Quantitative interpretation explains machine learning models for chemical reaction prediction and uncovers bias. *Nat. Commun.* **2021**, *12*, Publisher: Nature Research.
- (24) Durant, G.; Boyles, F.; Birchall, K.; Deane, C. M. The future of machine learning for small-molecule drug discovery will be driven by data. *Nat. Comput. Sci.* **2024**, 1–9, Publisher: Nature Publishing Group.
- (25) Segler, M. H.; Preuss, M.; Waller, M. P. Planning chemical syntheses with deep neural networks and symbolic AI. *Nature* **2018**, *555*, 604–610, Publisher: Nature Publishing Group.
- (26) Lowe, D. M. Extraction of chemical structures and reactions from the literature. Ph.D. thesis, University of Cambridge, Cambridge, 2012.
- (27) Schneider, N.; Stiefl, N.; Landrum, G. A. What's What: The (Nearly) Definitive Guide to Reaction Role Assignment. *J. Chem. Inf. Model.* **2016**, *56*, 2336–2346, Publisher: American Chemical Society.
- (28) Jin, W.; Coley, C. W.; Barzilay, R.; Jaakkola, T. Predicting Organic Reaction Outcomes with Weisfeiler-Lehman Network. 31st Conference on Neural Information Processing Systems. 2017; arXiv: 1709.04555.



- (29) Dai, H.; Li, C.; Coley, C. W.; Dai, B.; Song, L. Retrosynthesis Prediction with Conditional Graph Logic Network. 33rd Conference on Neural Information Processing Systems. 2019; arXiv: 2001.01408v1.
- (30) Torren-Peraire, P.; Hassen, A. K.; Genheden, S.; Verhoeven, J.; Clevert, D. A.; Preuss, M.; Tetko, I. V. Models Matter: the impact of single-step retrosynthesis on synthesis planning. *Digit. Discov.* **2024**, *3*, 558–572, arXiv: 2308.05522 Publisher: Royal Society of Chemistry.
- (31) Thakkar, A.; Kogej, T.; Reymond, J.-L.; Engkvist, O.; Bjerrum, E. J. Datasets and their influence on the development of computer assisted synthesis planning tools in the pharmaceutical domain. *Chem. Sci.* **2019**, *11*, 154–168, Publisher: The Royal Society of Chemistry.
- (32) Kearnes, S. M.; Maser, M. R.; Wleklinski, M.; Kast, A.; Doyle, A. G.; Dreher, S. D.; Hawkins, J. M.; Jensen, K. F.; Coley, C. W. The open reaction database. *J. Am. Chem. Soc.* **2021**, *143*, 18820–18826, Publisher: American Chemical Society.
- (33) Maziarz, K.; Tripp, A.; Liu, G.; Stanley, M.; Xie, S.; Gainski, P.; Seidl, P.; Segler, M. Re-evaluating retrosynthesis algorithms with syntheseus. *Faraday Discuss*, **2024**, –, arXiv: 2310.19796v1.
- (34) Bradshaw, J.; Zhang, A.; Mahjour, B.; Graff, D. E.; Segler, M. H. S.; Coley, C. W. Challenging reaction prediction models to generalize to novel chemistry. 2025; <http://arxiv.org/abs/2501.06669>, arXiv:2501.06669 [cs].
- (35) Segler, M. H.; Waller, M. P. Modelling Chemical Reasoning to Predict and Invent Reactions. *Chem. Eur. J.* **2017**, *23*, 6118–6128, Publisher: Wiley-VCH Verlag.
- (36) Yu, Y.; Yuan, L.; Wei, Y.; Gao, H.; Ye, X.; Wang, Z.; Wu, F. RetroOOD: understanding out-of-distribution generalization in retrosynthesis prediction. 2023; <http://arxiv.org/abs/2312.10900>, arXiv:2312.10900 [cs] version: 1.



- (37) Tu, H.; Shorewala, S.; Ma, P. T.; Thost, V. *Retrosynthesis Prediction Revisited*.
- (38) Chen, S.; Jung, Y. Assessing the Extrapolation Capability of Template-Free Retrosynthesis Models. 2024; <http://arxiv.org/abs/2403.03960>, arXiv: 2403.03960.
- (39) Westerlund, A. M.; Manohar Koki, S.; Kancharla, S.; Tibo, A.; Saigiridharan, L.; Kabeshov, M.; Mercado, R.; Genheden, S. Do Chemformers Dream of Organic Matter? Evaluating a Transformer Model for Multistep Retrosynthesis. *J. Chem. Inf. Model.* **2024**, *64*, 3021–3033, Publisher: American Chemical Society.
- (40) Mayfield, J.; Lagerstedt, I.; Sayle, R. *Pistachio "Fantastic reactions and how to use them"*; 2021.
- (41) Gil, V. S.; Bran, A. M.; Franke, M.; Schlama, R.; Luterbacher, J. S.; Schwaller, P. Holistic chemical evaluation reveals pitfalls in reaction prediction models. NeurIPS 2023 AI for Science Workshop. 2023; arXiv: 2312.09004v1.
- (42) Coley, C. W.; Green, W. H.; Jensen, K. F. RDChiral: an RDKit wrapper for handling stereochemistry in retrosynthetic template extraction and application. *J. Chem. Inf. Model.* **2019**, *59*, 2529–2537, Publisher: American Chemical Society.
- (43) Hastedt, F.; Bailey, R. M.; Hellgardt, K.; Yaliraki, S. N.; del Rio Chanona, E. A.; Zhang, D. Investigating the reliability and interpretability of machine learning frameworks for chemical retrosynthesis. *Digit. Discov.* **2024**,
- (44) Pang, J.; Vulić, I. Specialising and analysing instruction-tuned and byte-level language models for organic reaction prediction. *Faraday Discuss.* **2025**, *256*, 413–433, Publisher: The Royal Society of Chemistry.
- (45) Wieczorek, E.; Sin, J. W.; Holland, M. T. O.; Wilbraham, L.; Perez, V. S.; Bradley, A.; Miketa, D.; Brennan, P. E.; Duarte, F. Transfer learning for heterocycle synthe-





sis prediction. 2024; <https://chemrxiv.org/engage/chemrxiv/article-details/6617d56321291e5d1d9ef449>.

- (46) Johnson, J. M.; Khoshgoftaar, T. M. Survey on deep learning with class imbalance. *J. Big Data* **2019**, *6*, Publisher: SpringerOpen.



**Data Availability Statement:**

View Article Online  
DOI: 10.1039/D5DD00358J

**Preprocessing:** The open-source code used to split the training data is provided at <https://github.com/duartegroup/template-splits> (DOI: 10.5281/zenodo.17858529), which also contains a link to download the raw USPTO dataset. The final train, validation, and test splits are provided on FigShare (DOI: 10.6084/m9.figshare.30823988). The proprietary Pistachio dataset (licensed by NextMove Software) is not provided.

**Training:** The open-source packages used to train the machine learning models (using the configuration files provided at <https://github.com/duartegroup/template-splits/tree/main/configs>) can be found at:

- LocalRetro: <https://github.com/kaist-amsg/LocalRetro> (since removed by the authors)
- MEGAN: <https://github.com/molecule-one/megan>
- RootAligned: <https://github.com/otori-bird/retrosynthesis>

**Testing:** The open-source syntheseus package used to analyse the trained models can be found at <https://github.com/microsoft/syntheseus/tree/main>.

