

Cite this: *Chem. Sci.*, 2023, 14, 9959 All publication charges for this article have been paid for by the Royal Society of Chemistry

# Multistep retrosynthesis combining a disconnection aware triple transformer loop with a route penalty score guided tree search†

David Kreutter  and Jean-Louis Reymond  \*

Computer-aided synthesis planning (CASP) aims to automatically learn organic reactivity from literature and perform retrosynthesis of unseen molecules. CASP systems must learn reactions sufficiently precisely to propose realistic disconnections, while avoiding overfitting to leave room for diverse options, and explore possible routes such as to allow short synthetic sequences to emerge. Herein we report an open-source CASP tool proposing original solutions to both challenges. First, we use a triple transformer loop (TTL) predicting starting materials (T1), reagents (T2), and products (T3) to explore various disconnection sites defined by combining systematic, template-based, and transformer-based tagging procedures. Second, we integrate TTL into a multistep tree search algorithm (TTLA) prioritizing sequences using a route penalty score (RPScore) considering the number of steps, their confidence score, and the simplicity of all intermediates along the route. Our approach favours short synthetic routes to commercial starting materials, as exemplified by retrosynthetic analyses of recently approved drugs.

Received 27th March 2023  
Accepted 30th August 2023DOI: 10.1039/d3sc01604h  
rsc.li/chemical-science

## Introduction

Retrosynthetic analysis consists in drafting a synthetic sequence to produce a desired product from available starting materials. This analysis is one of the most useful but also difficult tasks in organic chemistry because it requires to integrate the large and complex set of rules that have emerged from millions of reactions reported in almost 200 years of organic synthesis. Computer-aided synthesis planning (CASP), initially conceived by E. J. Corey in the 1960s,<sup>1</sup> aims to harness the power of computers to automate retrosynthesis by exploiting data from experimental reactions collected in databases such as Reaxys<sup>2</sup> or the open-access reaction dataset extracted from US patent office data.<sup>3,4</sup> These databases list reactions of sets of starting materials (SM) and sets of reagents (R) to form one or several products (P).

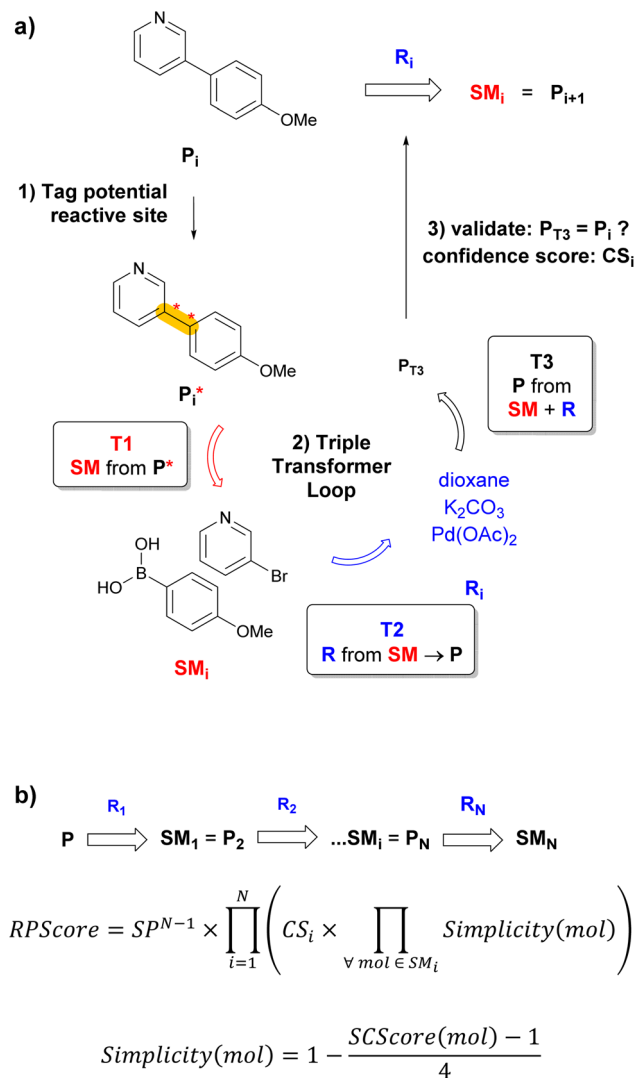
While expert systems based on hand-written rules such as Chematica/Synthia<sup>TM</sup> perform quite well for synthesis planning,<sup>5</sup> CASP ultimately aims to exploit artificial intelligence to automatically learn organic synthesis from reaction examples and propose synthetic routes for new molecules without human intervention.<sup>6–10</sup> Template-based approaches extract reaction rules in the form of substructure transformations and use

machine learning to learn their applicability domain from the structure of P in the training data.<sup>11–14</sup> On the other hand, transformer-based models use the linear SMILES<sup>15,16</sup> notation of chemical reactions and learn to translate the character string of P into the character string of SM + R, or *vice versa*.<sup>17–27</sup> The single-step predictions are then iterated to propose multistep retrosyntheses of target molecules from a selected set of building blocks (BB), which requires prioritizing possible routes using search algorithms such as Monte Carlo tree search,<sup>11,23,28</sup> and-or trees,<sup>25,29</sup> or a multistep graph exploration.<sup>24</sup>

Any CASP system must overcome two critical challenges to propose realistic retrosyntheses. First, the system must learn the context of reactions sufficiently well to propose reactions that make sense, but without overfitting such as to propose diverse retrosynthetic operations on previously unseen molecules. Second, the route-prioritizing algorithm must be designed to allow short sequences to emerge from the multitude of predicted possibilities.<sup>6</sup> Herein, we report a transformer-based retrosynthesis tool which proposes original solutions to both challenges. For single-step retrosynthesis, we use three different transformer models assembled as a triple transformer loop (TTL, Fig. 1a). To broaden the scope of predicted disconnections on a given target molecule, the TTL explores multiple disconnections by using products with tagged reaction centers (P\*) obtained by combining systematic, template-based and transformer-based tagging procedures. Compared to a transformer model trained on predicting SM + R directly from P\*, the TTL achieves better round-trip accuracy for single-step retrosynthesis. For multistep retrosynthesis predictions, we integrate

Department of Chemistry, Biochemistry and Pharmaceutical Sciences, University of Bern, Freiestrasse 3, 3012 Bern, Switzerland. E-mail: david.kreutter@unibe.ch; jean-louis.reymond@unibe.ch

† Electronic supplementary information (ESI) available. See DOI: <https://doi.org/10.1039/d3sc01604h>



and  $Simplicity(mol) = 1$ , if  $Mol \in Commercial\_Database$

**Fig. 1** Multistep retrosynthesis using TTLA. (a) Single-step retrosynthesis. At step  $i$ , each potential reactive site in  $P_i$  is identified systematically, using templates or a tagging transformer, and labelled to produce  $P_i^*$ . Transformer T1 is applied to  $P_i^*$  to predict  $SM_i$  (one or more starting materials), transformer T2 is applied to the top-scoring  $SM_i \rightarrow P_i$  to predict  $R_i$  (one or more reagents), and finally transformer T3 is applied to the top-scoring  $SM_i + R_i$  to produce  $P_{T3}$ . The prediction is finally validated if  $P_{T3} = P_i$  with confidence score  $CS_i$  of T3. Each molecule in the  $SM_i$  set is then used as product  $P_{i+1}$  for the next iteration. The route branches out if  $SM_i$  contains multiple molecules. (b) TTLA sequence and route penalty scoring. All molecules in the  $SM_i$  set of each step are used in the RPScore calculation of a linear sequence. See text for details.

the TTL into a multistep tree search algorithm, here named TTLA, which selects reaction sequences using a new route penalty score (RPScore), which for a route of  $N$  steps, is the product of a step-penalty score  $SP^N$ , the confidence scores of each single-step retrosynthesis ( $CS$ ), and the simplicity scores<sup>24</sup> of all SM along the route (Fig. 1b). This selection scheme favours

short sequences and is exemplified with the prediction of synthetic routes for recently approved drugs.

## Methods

### Dataset

The United States Patent and Trademark Office (USPTO) chemical reaction dataset version shared by Thakkar *et al.*<sup>30</sup> was used for the single-step evaluation as well as for training all transformer models in this study. The dataset is derived from the version of Lowe<sup>3,4</sup> and has been filtered by these authors to include reactions with a single product ( $P$ ) and between 2 and 10 starting materials ( $SM$ ) and reagents ( $R$ ) only. In the present work, we removed the tagging information, and reactions were remapped and retagged using our new SMILES tagging strategy and syntax. The same dataset split for training, validation, and test (90 : 5 : 5), as shared by Thakkar *et al.*<sup>30</sup> was used across all models resulting in 1 139 608, 63 672 and 63 454 reactions respectively.

### Tagging reaction centers

Training the disconnection-aware retrosynthesis model requires a training dataset where all product SMILES have tagged atoms. To tag reacting atoms, we use the atom-mapping tool shared by Schwaller *et al.*<sup>31</sup> to identify the atoms having an environmental change during the reaction, defined as reacting atoms. These reacting atoms are then re-labelled with the atom mapping label “1” while all other atom mapping labels are removed, as described by Byekwaso *et al.*<sup>32</sup> We then replace the atom tagging syntax by its unmapped SMILES notation, *e.g.* replacing “[C:1]” with “C”, and append the atom with another separated tagging token (“!”) using RDkit.<sup>33</sup> This modification allows to maintain an invariant SMILES token usage independent of the neighbouring hydrogen count or stereochemistry.

### Single-step disconnection aware retrosynthesis (T1)

Being able to identify the reaction center of a given reaction, we apply our reaction tagging algorithm on USPTO to obtain a retrosynthesis-tagged training dataset. We remove reagents, catalysts, and solvents, which are identified as the unmapped species in atom-mapped reactions and train the retrosynthesis model to predict the starting materials given as input the tagged product. We use the transformer architecture<sup>18</sup> and train it using the OpenNMT<sup>34,35</sup> library with standard previously-reported hyperparameters for this type of task.<sup>22</sup>

### Automatic tagging of potentially reactive atoms

We use three complementary methods to maximize the tagging possibilities while maintaining a reasonable number of predictions. First, we systematically tag all possible single atoms, pairs of directly connected atoms, and triplets of adjacent atoms (chain or three-membered ring). Secondly, we use templates for tagging the reactive sets of atoms corresponding to the conditional substructure with a variable radius (typically from 1 to 3). Templates occurring more than once and having between 1 and 10 reactive atoms were identified by analyzing



the original USPTO dataset. A given template can contain multiple disconnected sets of reactive atoms. Finally, the transformer model AutoTag reported by Thakkar *et al.*<sup>30</sup> was trained from untagged SMILES to the corresponding tagged molecule to provide additional tagging examples.

### Reagent prediction (T2)

Transformer T2 is trained from the untagged USPTO training set to identify reagents (R) from the combination of SM and P using the same hyperparameters as for T1. Note that R often includes actual reagents and solvents.

### Forward validation (T3)

The third model of the triple-transformer loop is a forward reaction prediction model trained with untagged reactions (molecular transformer).<sup>22</sup> We give this forward validation model the predicted SM<sub>i</sub> (from T1) and the predicted R<sub>i</sub> (from T2) as input separated by the ">" token. If T3 predicts the correct P<sub>i</sub> as its top-1 prediction, those SM<sub>i</sub> and R<sub>i</sub> are stored for the tree search. The confidence score CS<sub>i</sub> for the T3 prediction is used as confidence score for the reaction. T3 serves to filter down a large number of predictions to retain feasible reactions only.

### Single-step TTL tagging strategies study

The performance of individual tagging methods was studied on 500 molecules randomly selected from the USPTO test set for single-step TTL retrosynthesis to which we varied the three strategies over various parameters, changing the template radius from 1 to 3 and the transformer tagging (AutoTag) beam size from 1 to 1000.

### Route penalty score (RPScore)

The RPScore is computed for each predicted linear retrosynthetic sequence of N steps leading from the final product P to starting materials SM<sub>N</sub> (Fig. 1b). To reduce the score of long sequences, we introduce a step penalty SP, with  $0 < SP \leq 1$ , extended to SP<sup>N</sup> for a sequence of N steps. The RPScore is the product of SP<sup>N</sup> with the product of all confidence scores CS<sub>i</sub> (from the T3 prediction) for each individual step and the Simplicity (mol) for all intermediates along the sequence of N steps. By default, the penalty value SP is set to 0.8, but this could be adapted for every search in the configuration file of the multistep exploration. Simplicity(mol)<sup>24</sup> ranges from 0 for complex to 1 for simple molecules and is derived from the molecular synthetic complexity score (SCScore, ranging from 1 to 5) which describes molecular complexity taking synthetic accessibility into account.<sup>36</sup> Here, we assign a value of 1 if the molecule occurs in the BB set of commercial starting materials. In contrast to Schwaller *et al.*,<sup>24</sup> we exclude reagents R<sub>i</sub> from the Simplicity calculation to avoid penalizing steps that use reagents with low calculated Simplicity, which is rarely a measure of their availability or ease of use.

### Multistep exploration strategy

We use a Heuristic Best-First Tree Search algorithm with beam search and iterative expansion to explore retrosynthetic routes as similarly reported for transformer-based retrosynthesis.<sup>24</sup> Once predictions of an iteration are complete, the tree search updates and lists all possible routes, and computes the RPScore. Unsolved routes are sorted by decreasing RPScore. The top 20 unsolved routes, which lead to starting materials absent from the selected set of commercially available building blocks, are selected for expansion by defining them as products P<sub>i</sub> and new SM<sub>i</sub> are predicted by applying a single-step retrosynthesis using TTL. The resulting set of predicted single-step retrosynthesis is updated back to the tree wherever those SMs were present. The tree is updated for the next iteration. The process stops when a chosen minimum number of solved routes or a maximum number of iterations has been reached.

### Building block (BB) set

We combined MolPort (<https://www.molport.com>) and Enamine (<https://www.enamine.net>) databases to build our database of 534 058 commercially available compounds as the building block (BB) set.

## Results and discussion

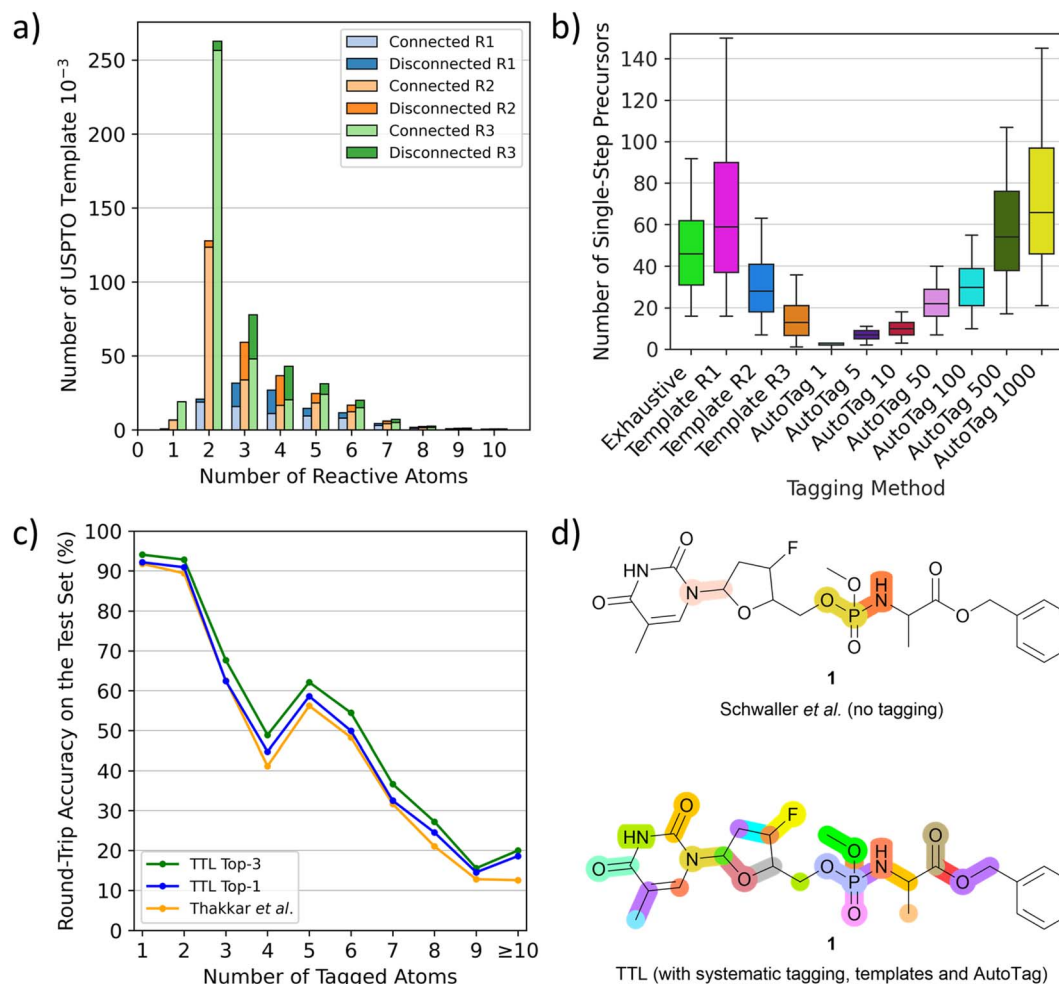
### Training transformer T1 for single-step retrosynthesis

Initially, we use the atom-mapping transformer<sup>31</sup> information to annotate reacting atoms in all products P in the training data, which results in a training dataset containing labelled P\*. Our code is inspired by the recent report by Byekwaso *et al.*,<sup>32</sup> however with a slightly simplified syntax for tagged atoms. Using the tagged P\* data, we then train a transformer model T1 to predict SM from P\*, a task which is simpler than predicting both SM and R from P\* as done by Byekwaso *et al.*<sup>32</sup>

### Tagging potential reactive sites

To use T1 to predict possible SM<sub>i</sub> from a given product P<sub>i</sub> at step i, one must first tag potentially reacting atoms in P<sub>i</sub>. We do this using complementary methods. First, we tag all single atoms as well as pairs and triplets of adjacent atoms systematically in P<sub>i</sub>. Second, we systematically apply templates extracted from tagged P\* in the USPTO dataset. These templates with various conditional radii (from 1 to 3) are substructures containing up to ten tagged atoms, not necessarily connected. Although the most frequent templates are those with two or three connected atoms, which are tags also obtained in the systematic procedure, the templates also include many tags with disconnected atom pairs and triplets as well as tags with four or more atoms, which are missing from the systematic tagging procedure (Fig. 2a). As a third tagging option, we use the tagging approach recently reported by Thakkar *et al.*<sup>30</sup> where reacting atoms are identified using a tagging transformer, here named AutoTag, trained to learn the detailed context from the tagged dataset. The number of predicted tags (sorted by confidence score, called beam size) of AutoTag can be varied to generate a given number of possible tags to extend the retrosynthesis options.





**Fig. 2** TTL and automatic atom tagging. (a) Distribution of the number of tagging templates extracted from USPTO depending on the number of atoms it tags, named “reactive atoms”. Triple bar plot to show the differences between conditional radiuses beyond tagged atoms from 1 to 3 (R1 to R3). Bars are split into a light-coloured part representing the fraction of templates that tags bond-connected atoms and dark-coloured for disconnected atoms. (b) Number of validated single-step starting materials (“precursors”) on the TTL generated depending on the automatic tagging strategy, tested over 500 molecules randomly selected from the TTL test set. (c) Round-trip accuracies of the TTL using the top-1 SM by T1 and the top-1 or top-3 R predicted by T2, compared to the disconnection-aware dual transformer of Thakkar et al.<sup>30</sup> (d) Highlighted disconnection sites of the antiviral molecule 1 using the untagged retrosynthesis and forward validation models of Schwaller et al.,<sup>24</sup> leading to four unique sets of starting materials among three reactive sites (top) and the TTL augmented by systematic tagging, template-based tagging (radius 2) and AutoTag (beam size 50) after forward validation leading to 231 unique sets of starting materials among 26 reactive sites (bottom).

Analyzing the performance of the different tagging methods shows that less restrictive template radius or high AutoTag beam size both lead to an increased number of tagged atoms per molecule (Fig. S1†) as well as a much higher number of generated tagged SMILES (Fig. S2†) and significantly more single-step starting materials (Fig. 2b and S3†), but also to a lower number of high confidence predictions (Fig. S4†), indicating that most of the additionally obtained tags are less chemically meaningful (Fig. S5†). Moreover, the tagging efficiency, evaluated by dividing the number of successful retrosynthetic steps obtained by the number of TTL rounds (number of tags), drops for high AutoTag beam sizes and low radius templates (Fig. S6†). To obtain a good number of validated retrosynthetic steps at reasonable computing cost, we combine three strategies: the systematic tagging (1, 2 and 3 atoms),

templates with a radius of 2, and the AutoTag transformer with a beam size of 50. A Venn diagram analysis of the number of unique retrosynthetic steps obtained shows that 17% of the steps (37.8% of high confidence steps) are predicted by all three methods, while 52.6% of the steps (25.5% of high confidence steps) are coming from only one of the three tagging methods, highlighting their complementarity (Fig. S7 and S8†).

### Triple transformer loop (TTL)

To initiate a validated single-step retrosynthesis prediction for product  $P_i$ , we run T1 on all  $P_i^*$  obtained by the combined selected tagging procedures described above. The transformer outputs a series of possible  $SM_i$ , which are sorted in order of the T1 confidence score. For the top- $B$   $SM_i$  (beam size  $B = 1$  or more), we then apply a second transformer (T2) trained to



predict R from  $SM \rightarrow P$ . For each  $SM_i$ , T2 outputs a series of possible  $R_i$ , from which we retain the top- $B'$  (beam size  $B' = 1$  or more). The TTL is completed with a forward validation<sup>37</sup> transformer (T3) trained to predict P from  $SM + R$  using the same training dataset used for T1 and T2. For all combinations of top  $SM_i$  predicted by T1 and top  $R_i$  predicted by T2, we finally use T3 to predict the most likely product  $P_{T3}$ . The TTL prediction is validated if the top-1 predicted  $P_{T3}$  is identical to the input product  $P_i$  (Fig. 1a). The T3 confidence scores  $CS_i$  of the validated predictions  $SM_i + R_i$  are used to select the best  $R_i$  if  $B' > 1$ , and to calculate the route penalty score (RPScore, see below).

### Performance evaluation

The performance of TTL can be compared with previous single-step retrosynthesis models at three different levels. First, transformer T1, which predicts SM from the tagged product  $P^*$ , can be compared with other single-step retrosynthesis models predicting SM from P, both transformer-based and template-based.<sup>17,19,21,23–27</sup> While these models perform between 40% and 55% top-1 accuracy, our tagged T1 achieves 66% top-1 accuracy, which shows that tagging provides a significant advantage for this task.

Second, the performance of the TTL loop can be compared with the disconnection-aware retrosynthesis model of Thakkar *et al.*<sup>30</sup> in terms of single-step round-trip prediction accuracy from the tagged product  $P^*$ , which is the accuracy of predicting P from the  $SM + R$  initially predicted from  $P^*$ . TTL using only the top-1 predictions for T1 and T2 performs comparably to Thakkar's disconnection-aware retrosynthesis model (80.44% *vs.* 79.09% accuracy). The TTL performance increases to 83.04% when considering the top-1 prediction of T1 and the top-3 predictions of T2. Similar to the observation by Thakkar *et al.*,<sup>30</sup> we furthermore find that the prediction accuracy strongly decreases as a function of the number of tagged atoms (Fig. 2c). Subsequently to our preprinted report, a separate study has investigated the performance of the reagent prediction transformer.<sup>38</sup>

Thirdly, one can compare the single-step round trip accuracy of TTL with that of the non-tagged retrosynthesis model of Schwaller *et al.*,<sup>24,37</sup> who evaluated if a forward prediction model predicted the correct product P from the  $SM + R$  predicted by their model from the non-tagged P. As discussed by Thakkar *et al.*,<sup>30</sup> the untagged transformer may sometimes choose a different and easier to predict disconnection than that recorded in the test set, and therefore performs slightly better (82.4% top-1 accuracy) than the tagged transformer, which is forced by tagging to apply the retrosynthesis of the test set. Here, we find that the top-1 round-trip prediction accuracy ( $P \rightarrow P$ ), obtained by applying our multiple tagging procedure followed by the TTL, reaches 99.9%, which means that our approach is almost always able to propose at least one forward-validated possible retrosynthetic step from any product molecule.

Furthermore, a critical feature of any single-step retrosynthesis model in view of multi-step retrosynthesis concerns the diversity of possible disconnections proposed. We find that this

diversity is greatly enhanced by the multiple tagging approach. For instance, when tested on unseen molecules, the TTL combined multiple tagging provides validated disconnections at several possible reactive sites. By contrast, the baseline transformer, trained as reported by Schwaller *et al.*<sup>24</sup> to produce SM directly from P using the unannotated data for training, chooses fewer disconnection points, as exemplified here for the pro-nucleotide **1** (Fig. 2d).<sup>39</sup>

### Multistep retrosynthesis

By integrating the single-step retrosynthesis TTL into a multi-step tree search, we obtain a multistep retrosynthesis algorithm, here named TTLA. In each retrosynthesis iteration, TTLA runs the TTL exhaustively on all SM of the preceding iteration, newly defined as P, and ranks the routes to the newly predicted SM using a composite route penalty score RPScore (Fig. 1b, see Methods for details).

When prioritizing multiple retrosynthesis options during the tree search, TTLA uses the RPScore to rank the different routes leading to the SM produced in the latest iteration of TTL, and only extends retrosynthesis on a small number (typically 20) of SM taken from the top RPScore routes. Because each additional step imposes a penalty (usually  $P = 0.8$ ), lengthy routes and unproductive loops involving protection/deprotection cycles of the same functional group are rapidly falling down the RPScore priority list, which leads the algorithm to explore alternative routes, so that short synthetic sequences are eventually prioritized even if their first retrosynthetic steps were initially not top scoring.

As commonly observed with CASP tools as well as with transformer models in general, the top-scoring outputs of TTLA must be inspected to identify relevant predictions. While the RPScore is used in the tree search, we find relevant routes by inspecting both the top-RPScore route and the top-CScore routes ( $CScore(\text{route}) = \text{the product of } CS_i \text{ for all steps}$ ) in the TTLA output, as discussed below with examples.

TTLA is exemplified here for predicting the synthesis of two drug molecules approved in 2020, namely fostemsavir (**2**, Fig. 3), a prodrug which upon phosphatase cleavage releases the anti-retroviral agent temsavir as HIV entry inhibitor,<sup>40</sup> and ozanimod (**10**, Fig. 4), a sphingosine-1-phosphate receptor antagonist used as an immunomodulatory agent to treat multiple sclerosis.<sup>41</sup> The commercial process for both drugs was recently reviewed.<sup>42</sup> None of the synthetic steps involved in these two processes occur in the USPTO dataset used for training TTL, making them a good test case for TTLA. For these examples, we challenged TTLA to predict synthetic routes starting from a list of 534 058 commercially available BB.

The reported commercial process for the antiviral drug fostemsavir (**2**, Fig. 3a, details in Fig. S9†) is a linear sequence involving the sequential *C*-acylation of pyrrolopyridine **3** with oxalyl monochloride **4a** (step a) and benzoylpiperazine (**5a**, step b), followed by coupling of with triazole **6** (step c), *N*-alkylation of the pyrrole with the protected chloromethylphosphate **7a** (step d), and finally deprotection of the *tert*-butyl ester protecting groups (step e).



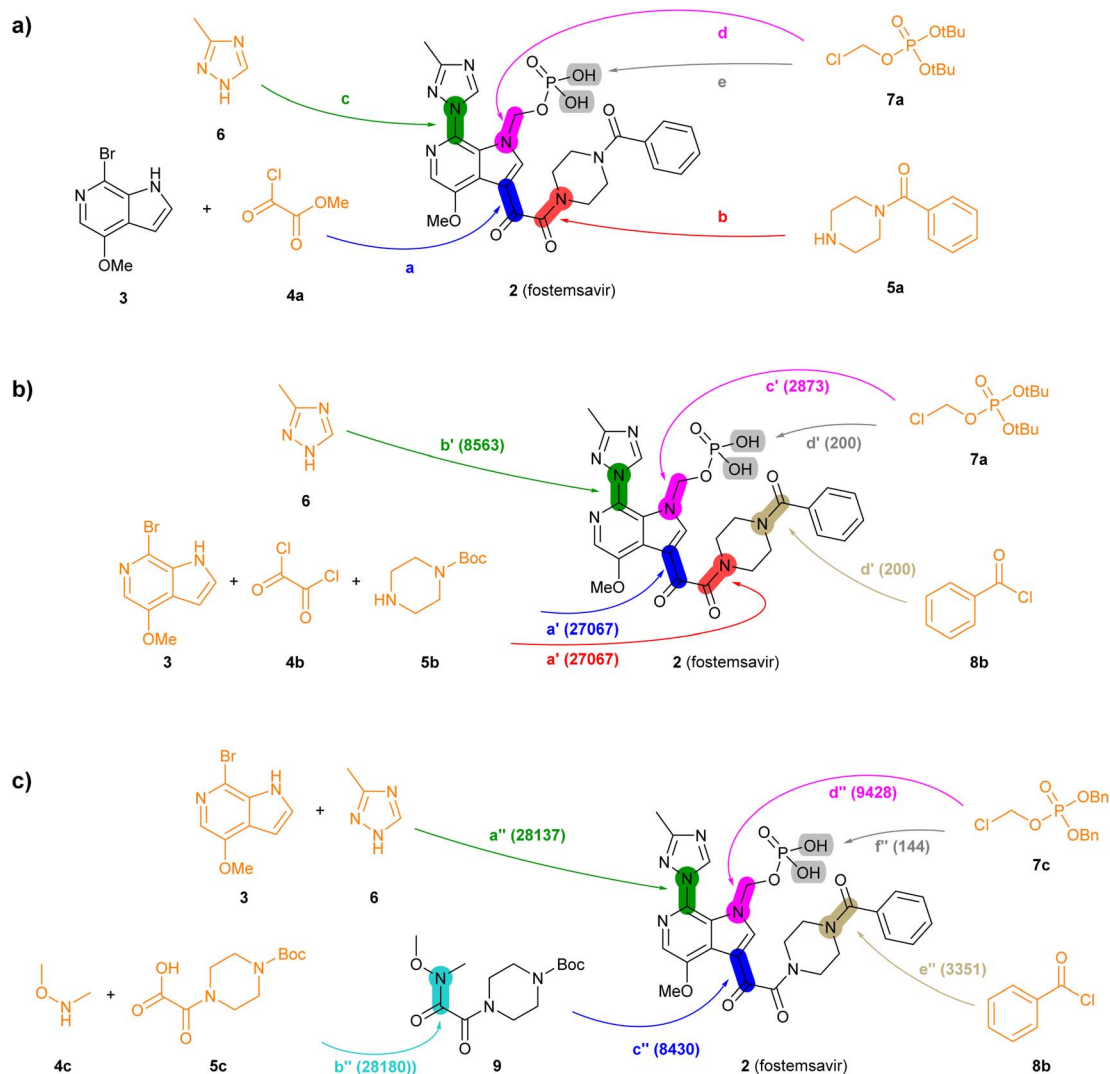


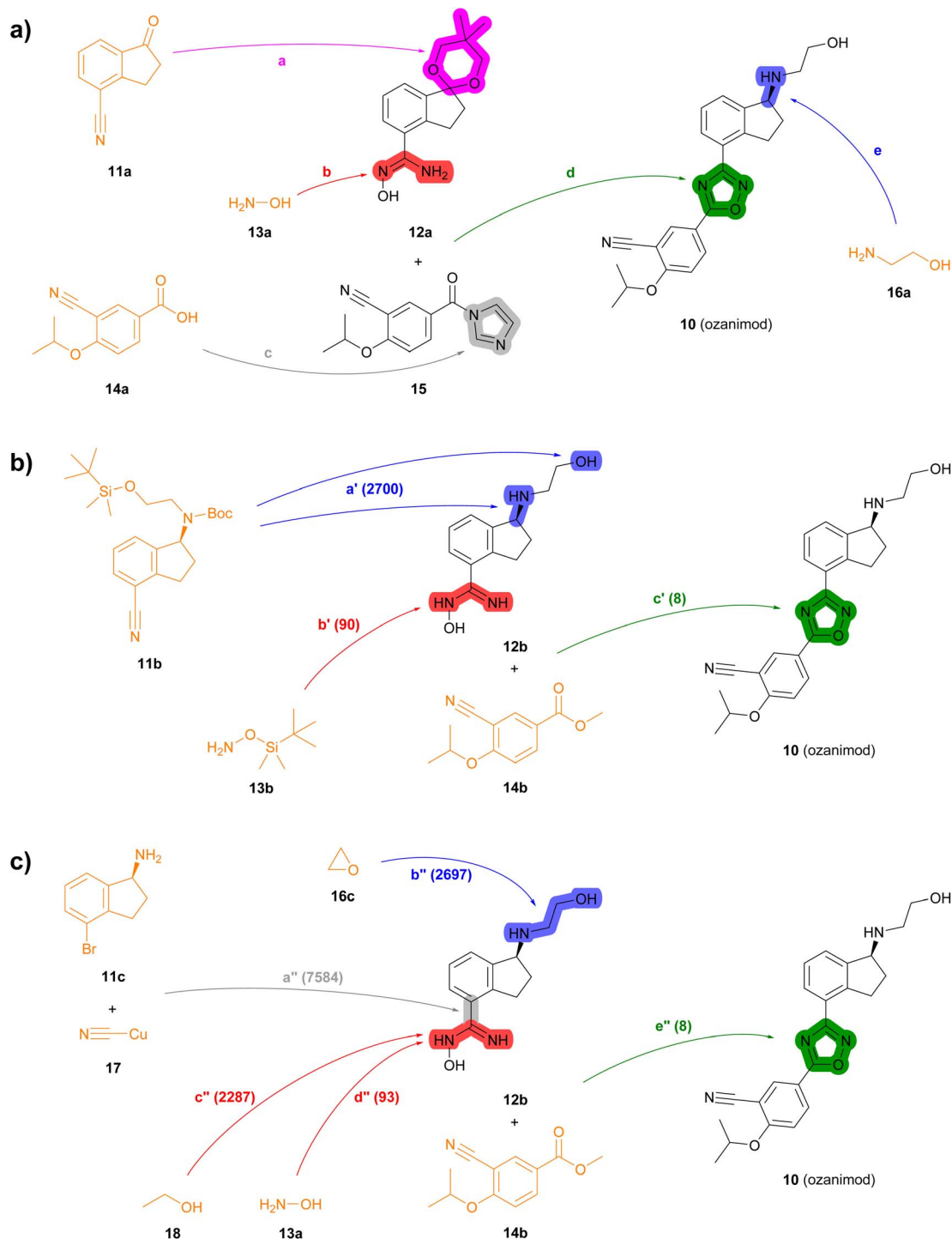
Fig. 3 Summary of reported and TTLA predicted routes for fostemsavir 2. Bonds formed in each step are highlighted in colour. Numbers in parenthesis correspond to the order in which the multistep tree search prioritized predictions. The full retrosynthesis routes are drawn out in the ESI† Fig. S9–S11.† (a) Commercial process. Reported reagents: (a)  $\text{AlCl}_3$ ,  $\text{Bu}_4\text{NHSO}_4$ ,  $\text{CH}_2\text{Cl}_2$ , then  $\text{KOH}$ , then  $\text{H}_3\text{PO}_4$ ; (b)  $\text{Ph}_2\text{POCl}$ ,  $\text{NMM}$ ,  $\text{NMP}$ ; (c)  $\text{KOH}$ ,  $\text{CuI}$ , then  $\text{KOH}$ ,  $\text{EtOH}$ ,  $\text{LiI}$ ; (d)  $\text{Et}_4\text{NI}$ ,  $\text{K}_2\text{CO}_3$ ,  $\text{CH}_3\text{CN}/\text{H}_2\text{O}$ ; (e)  $\text{AcOH}$ ,  $\text{H}_2\text{O}$ . (b) Highest TTLA RPScore route. Predicted reagents: (a')  $\text{Et}_3\text{N}$ ,  $\text{CH}_2\text{Cl}_2$ ; (b')  $\text{K}_2\text{CO}_3$ ,  $\text{CuI}$ , toluene; (c')  $\text{K}_2\text{CO}_3$ ,  $\text{DMF}$ ; (d')  $\text{HCl}$ ,  $N,N$ -diisopropylethylamine,  $\text{H}_2\text{O}$ , dioxane. (c) Highest TTLA CScore route. Predicted reagents: (a'') (2S)-pyrrolidine-2-carboxylic acid,  $\text{K}_2\text{CO}_3$ ,  $\text{CuI}$ ,  $\text{EtOAc}$ ,  $\text{DMSO}$ ; (b'') no reagent predicted; (c'')  $n\text{-BuLi}$ ,  $\text{THF}$ ; (d'')  $\text{K}_2\text{CO}_3$ ,  $\text{DMF}$ ; (e'')  $\text{TFA}$ ,  $\text{DMAP}$ ,  $\text{CH}_2\text{Cl}_2$ , (f'')  $\text{Pd}$ ,  $\text{EtOH}$ .

When challenged with 2, TTLA proposes many possible routes from similar starting materials as the commercial process, but in a different order. The highest RPScore route is a linear sequence starting from the double *C*- and *N*-alkylation of oxalyl chloride (4b) with pyrrolopyridine (3) and 1-boc-piperazine (5b) in one pot (step a', Fig. 3b, details in Fig. S10†). The aryl bromide of the resulting intermediate is then substituted with triazole 6 (step b'), and its pyrrole NH group is alkylated with *tert*-butyl chloromethyl phosphate 7a, similarly to the commercial route (step c'). In the final step, the phosphate and the piperazine groups are deprotected with acid, followed by benzoylation of the free piperazine with benzoylchloride (8b) to form fostemsavir 2 (step d').

On the other hand, the highest CScore route is a convergent sequence starting with alkylation of triazole 6 with pyrrolopyridine 3 on the one hand (step a'', Fig. 3c, details in Fig. S11†), and the preparation of the Weinreb amide 9 from boc-oxalylpiperazine 5c and *N,O*-dimethylhydroxylamine 4c on the other hand (step b''). The resulting intermediates are then coupled (step c''), and the product is *N*-alkylated on the pyrrole nitrogen with benzyl-protected chloromethyl phosphate 7c (step d''). Deprotection of the piperazine group allows the acylation with benzoylchloride (8b, step e''). Reductive deprotection of the benzyl phosphate esters finally gives the product 2 (step f'').

In the second example, the drug ozanimod 10 is synthesized commercially in a convergent sequence of 7 steps from ketone 11a and benzoic acid 14a (Fig. 4a, details in Fig. S12†). After





**Fig. 4** Summary of reported and TTLA predicted routes for ozanimod **10**. Bonds formed in each step are highlighted in colour. Numbers in parenthesis correspond to the order in which the multistep tree search prioritized predictions. The full retrosynthesis routes are drawn out in the ESI† Fig. S12–S14.† (a) Commercial process. Reported reagents: (a)  $\text{HC}(\text{OMe})_3$ ,  $p$ -TsOH,  $\text{PhCH}_3$ ; (b)  $\text{NH}_2\text{OH} \cdot \text{HCl}$ ,  $\text{Et}_3\text{N}$ ; (c) carbonyl diimidazole; (d)  $\text{NaOH}$ ; (e) (i)  $p$ -TsOH, acetone, (ii)  $\text{NH}_2\text{CH}_2\text{CH}_2\text{OH}$ ,  $p$ -TsOH,  $\text{PhCH}_3$ , (iii) chiral Ru-complex,  $\text{Et}_3\text{N}/\text{HCO}_2\text{H}$ . (b) Highest TTLA RPScore route. Predicted reagents: (a')  $\text{HCl}$ , dioxane; (b')  $\text{ZnCl}_2$ ,  $\text{AcOEt}$ , toluene; (c')  $\text{HCl}$ ,  $t$ -BuOK, THF. (c) Highest TTLA CScore route. Predicted reagents: (a'') 1-Methylpyrrolidin-2-one; (b'') no reagent predicted; (c'')  $\text{HCl}$ ,  $\text{Et}_2\text{O}$ ; (d'')  $\text{HCl}$ ,  $\text{NaHCO}_3$ ,  $\text{EtOH}$ ; (e'')  $\text{HCl}$ ,  $t$ -BuOK, THF.

initial protection of ketone **11a** as an acetal (step a), its nitrile group is reacted with hydroxylamine **13a** to form the *N*-hydroxyamidine intermediate **12a** (step b). In parallel, benzoic acid **14a** is activated to the corresponding benzoyl imidazole **15** (step c).

Intermediates **12a** and **15** are then condensed to form the oxazole ring (step d). The acetal group of the resulting intermediate is then deprotected and condensed with ethanolamine (**16a**) to the



corresponding imine, which is reduced enantioselectively using a chiral ruthenium catalyst to form **10** (step e).

Many of the high-scoring routes identified with TTLA are extremely short sequences starting with commercially available close analogs of the drug and were removed from the list of top-scoring routes. Interestingly, TTLA also proposes routes that resemble the commercial process but start from chiral starting materials such as aminoinindanes **11b** and **11c**, which avoids the enantioselective reaction used for the commercial process. For example, the best RPScore route is a linear synthesis from **11b** starting with the removal of the Boc and TBS protecting groups of the ethanolamine side chain and conversion of the cyano group to the corresponding *N*-hydroxyamidine by reaction with TBS-hydroxylamine (**13b**) to form intermediate **12b** (steps a' and b', Fig. 4b, details in Fig. S13†). The third and final step of this short sequence is the condensation of *N*-hydroxyamidine **12b** with cyanobenzoate **14b** yielding ozanimod **10** (step c').

The best CScoring route is a somewhat longer linear sequence employing the same condensation of **12b** and **14b** as the final step (step e'', Fig. 4c, details in Fig. S14†). In this proposed sequence however, intermediate **12b** requires four steps from the chiral aminobromointhane **11c** as follows. First, the cyano group is installed by reaction of the aryl bromide with copper cyanide (step a''). Second, the primary amine reacts with ethylene oxide **16c** to form the *N*-hydroxyethyl side chain (step b'). Third, the cyano group introduced in step a'' reacts with ethanol (**18**) to form an ethyl imidate intermediate (step c''), which further reacts with ethanolamine (**13a**) in a fourth step to form the *N*-hydroxyamidine group in **12b** (step d'').

Analyzing the details of the TTLA collective output shows that, although TTLA did not formulate routes identical to the commercial processes, the set of commercial starting materials used by TTLA are very similar to those used in the reported commercial processes for both drugs (Fig. S15 and S16†). In fact, all starting materials used in the commercial process for fostemsavir are present in the set for this drug.

In terms of individual reaction steps, we find that TTLA explores a large number of single reactions to arrive at the top-scoring short routes proposed in the above retrosynthesis. In the case of fostemsavir, the key retrosynthetic *C*- and *N*-acylation of the oxalyl starting material is discovered after 27 067 single predicted steps (Fig. 3b, step a'), probably because this step is rather complex and unusual. In the case of ozanimod, TTLA performed 7594 individual single-step predictions to arrive at the proposed retrosyntheses, with the best scoring route being discovered after 2700 iterations. Interestingly, the formation of the oxadiazole ring is discovered already at iteration 8 (Fig. 4b, step c'). It should be noted that the order of iterations and therefore the number of attempts necessary to identify high-scoring routes depends on the scoring function used to prioritize node expansion, here the RPScore, which takes the simplicity and number of steps into account.

The output of TTLA can be visualized by representing the collective predicted single steps in a TMAP<sup>43</sup> computed using the differential reaction fingerprint (DRFP)<sup>44</sup> as a similarity measure. As illustrated for ozanimod, colour-coding by step iteration number indicates that TTLA explores a broad diversity of steps directly from the beginning of the retrosynthesis exploration, which we attribute to our diverse reaction center

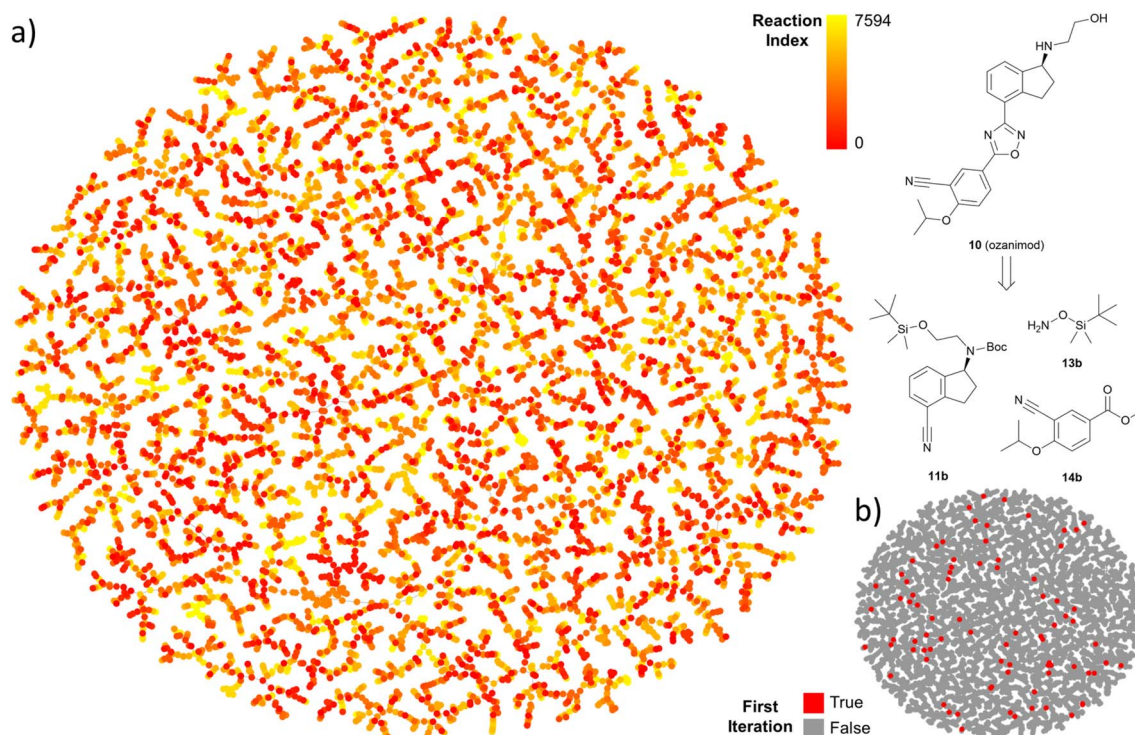


Fig. 5 TMAP representation of iterated predictions for the multistep search of ozanimod. (a) Predicted reactions from the target molecule (low indexes) to end nodes. (b) Highlighted first iteration of the TTLA search. Interactive map available at <https://tm.gdb.tools/TTLA/ozanimod>.





tagging approach used (Fig. 5a). This diversity is also visible when colour-coding all steps involving the final product, corresponding to the initial retrosynthesis, which are broadly distributed on the map (Fig. 5b). A similar pattern is visible in the TMAP of the predicted single steps for fostemsavir (Fig. S17†).

### Comparing TTLA with other retrosynthesis tools

Previous retrosynthesis tools, template-based or transformer-based, predict starting material from products by applying the most probable retrosynthetic operation according to a training set. Here we combined exhaustive and template-based methods to label many potential reactive sites, which lead us to test many possible disconnections (Fig. 2d). These potential reactive sites were then challenged with the TTL, which produced detailed predictions including starting materials and reagents. In the examples discussed above TTLA identified short routes comparable to the reported processes, which were all examples of optimized production routes.

By comparison, a currently available version of AiZynth-Finder (v3.7.0),<sup>14</sup> a templated-based retrosynthesis tool, fails to propose a synthesis for fostemsavir due to its inability to find a synthesis for a bis-*tert*-butyl phosphate starting material (Fig. S18†). AiZynthFinder furthermore proposes a short route similar to TTLA for ozanimod, although including somewhat less realistic steps, for example, an alkylation of a primary amine with 2-bromoethyl acetate which would probably rather lead to acetyl transfer, and no indication of reagents (Fig. S19†). On the other hand, the online portal of IBM RXN for chemistry,<sup>45</sup> which uses a transformer model, predicts essentially the same route as TTLA for fostemsavir (Fig. S20†). For ozanimod however, this tool settles on an eight-step route which, although containing realistic steps, is simply much longer than the commercial process or the route proposed by TTLA (Fig. S21†). For both of these retrosynthesis tools, whether the routes are part of their training sets is not known.

To statistically evaluate our TTLA, we selected target molecules from the retrosynthesis benchmark dataset shared by Genheden *et al.* which were absent from our training dataset.<sup>46</sup> Due to the high computing time of our method, a random subset of 240 target molecules was selected. Solved routes involving reaction steps present in our training dataset were removed from the evaluation. TTLA proposed retrosyntheses to commercially available starting materials for 97.5% of the target molecules, which is comparable to the performance of other retrosynthetic tools reported in the original paper.<sup>46</sup> Selected examples are shown in Fig. S22–S31.†

## Conclusion

In summary, our data shows that a triple transformer loop (TTL) operating on products with tagged reactive atoms achieves efficient single-step retrosynthesis predictions. TTL was integrated into a tree-exploration strategy using a route penalty scoring scheme to form the multistep retrosynthesis tool TTLA, which can predict short synthetic routes for drug molecules.

Since our approach uses transformer models, it should be possible to specialize TTLA for specific reaction classes by transfer learning similar to transformer models for forward prediction.<sup>47</sup> Furthermore, predicting SM from P and R from SM + P separately might be potentially adapted to reactions with more complex reagents such as enzymes<sup>48–50</sup> and help expand the scope of CASP systems. It should however be noted that the use of multiple transformer models and the detailed analysis of many possible disconnections renders our approach relatively slow, requiring up to several hours of computing time for a full retrosynthetic analysis. Efficiency increases might be possible in the future by fine-tuning the selection of potential disconnections and improving the tree search.

## Data availability

Code and instructions to compute multistep retrosynthesis as well as the code to tag reactive sites can be found on our GitHub repository:

<https://github.com/reymond-group/MultiStepRetrosynthesisTTL>. The original USPTO dataset can be found at <https://doi.org/10.6084/m9.figshare.5104873.v1>. The derived version of USPTO of Thakkar *et al.* could be found in their Zenodo repository.<sup>30,51</sup>

## Author contributions

DK designed and carried out the study and wrote the paper, JLR designed and supervised the study and wrote the paper.

## Conflicts of interest

The authors declare that they have no competing interests.

## Acknowledgements

This work was supported financially by Novartis. We would like to thank Dr Thierry Schlama, Dr Daniel Kaufmann, Dr Radka Snajdrova, Dr John Lopez, Dr Frederic Stanger, Dr Fabrice Gallou and Dr Thomas Ruch, for helpful discussions. Calculations were performed on UBELIX (<http://www.id.unibe.ch/hpc>), the HPC cluster at the University of Bern.

## References

- 1 E. J. Corey and W. T. Wipke, Computer-Assisted Design of Complex Organic Syntheses, *Science*, 1969, **166**(3902), 178–192, DOI: [10.1126/science.166.3902.178](https://doi.org/10.1126/science.166.3902.178).
- 2 A. J. Lawson, J. Swienty-Busch, T. Géoui and D. Evans, The Making of Reaxys—Towards Unobstructed Access to Relevant Chemistry Information, in *The Future of the History of Chemical Information*, ACS Symposium Series, American Chemical Society, 2014, vol. 1164, pp. 127–148, DOI: [10.1021/bk-2014-1164.ch008](https://doi.org/10.1021/bk-2014-1164.ch008).
- 3 D. M. Lowe, Extraction of Chemical Structures and Reactions from the Literature, PhD thesis, University of Cambridge, 2012, DOI: [10.17863/CAM.16293](https://doi.org/10.17863/CAM.16293).



- 4 D. M. Lowe, Chemical Reactions from US Patents (1976-Sep2016), *Figshare Dataset*, 2017, DOI: [10.6084/m9.figshare.5104873.v1](#).
- 5 S. Szymkuć, E. P. Gajewska, T. Klucznik, K. Molga, P. Dittwald, M. Startek, M. Bajczyk and B. A. Grzybowski, Computer-Assisted Synthetic Planning: The End of the Beginning, *Angew. Chem., Int. Ed.*, 2016, 55(20), 5904–5937, DOI: [10.1002/anie.201506101](#).
- 6 C. W. Coley, W. H. Green and K. F. Jensen, Machine Learning in Computer-Aided Synthesis Planning, *Acc. Chem. Res.*, 2018, 51(5), 1281–1289, DOI: [10.1021/acs.accounts.8b00087](#).
- 7 F. Strieth-Kalthoff, F. Sandfort, M. H. S. Segler and F. Glorius, Machine Learning the Ropes: Principles, Applications and Directions in Synthetic Chemistry, *Chem. Soc. Rev.*, 2020, 49(17), 6154–6168, DOI: [10.1039/C9CS00786E](#).
- 8 A. Thakkar, S. Johansson, K. Jorner, D. Buttar, J.-L. Reymond and O. Engkvist, Artificial Intelligence and Automation in Computer Aided Synthesis Planning, *React. Chem. Eng.*, 2021, 6(1), 27–51, DOI: [10.1039/D0RE00340A](#).
- 9 K. Molga, S. Szymkuć and B. A. Grzybowski, Chemist Ex Machina: Advanced Synthesis Planning by Computers, *Acc. Chem. Res.*, 2021, 54(5), 1094–1106, DOI: [10.1021/acs.accounts.0c00714](#).
- 10 P. Schwaller, A. C. Vaucher, R. Laplaza, C. Bunne, A. Krause, C. Corminboeuf and T. Laino, Machine Intelligence for Chemical Reaction Space, *Wiley Interdiscip. Rev. Comput. Mol. Sci.*, 2022, 12(5), e1604, DOI: [10.1002/wcms.1604](#).
- 11 M. H. S. Segler, M. Preuss and M. P. Waller, Planning Chemical Syntheses with Deep Neural Networks and Symbolic AI, *Nature*, 2018, 555(7698), 604–610, DOI: [10.1038/nature25978](#).
- 12 C. W. Coley, R. Barzilay, T. S. Jaakkola, W. H. Green and K. F. Jensen, Prediction of Organic Reaction Outcomes Using Machine Learning, *ACS Cent. Sci.*, 2017, 3(5), 434–443, DOI: [10.1021/acscentsci.7b00064](#).
- 13 A. Thakkar, T. Kogej, J.-L. Reymond, O. Engkvist and E. J. Bjerrum, Datasets and Their Influence on the Development of Computer Assisted Synthesis Planning Tools in the Pharmaceutical Domain, *Chem. Sci.*, 2019, 11(1), 154–168, DOI: [10.1039/C9SC04944D](#).
- 14 S. Genheden, A. Thakkar, V. Chadimová, J.-L. Reymond, O. Engkvist and E. Bjerrum, AiZynthFinder: A Fast, Robust and Flexible Open-Source Software for Retrosynthetic Planning, *J. Cheminf.*, 2020, 12(1), 70, DOI: [10.1186/s13321-020-00472-1](#).
- 15 D. Weininger, SMILES, a Chemical Language and Information System. 1. Introduction to Methodology and Encoding Rules, *J. Chem. Inf. Comput. Sci.*, 1988, 28(1), 31–36, DOI: [10.1021/ci00057a005](#).
- 16 D. Weininger, A. Weininger and J. L. Weininger, SMILES. 2. Algorithm for Generation of Unique SMILES Notation, *J. Chem. Inf. Comput. Sci.*, 1989, 29(2), 97–101, DOI: [10.1021/ci00062a008](#).
- 17 J. Nam and J. Kim, Linking the Neural Machine Translation and the Prediction of Organic Chemistry Reactions, *arXiv*, 2016, preprint, arXiv:1612.09529, DOI: [10.48550/arXiv.1612.09529](#).
- 18 A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser and I. Polosukhin, Attention Is All You Need, in *Advances in Neural Information Processing Systems 30*, ed. I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan and R. Garnett, Curran Associates, Inc., 2017, pp. 5998–6008.
- 19 B. Liu, B. Ramsundar, P. Kawthekar, J. Shi, J. Gomes, Q. Luu Nguyen, S. Ho, J. Sloane, P. Wender and V. Pande, Retrosynthetic Reaction Prediction Using Neural Sequence-to-Sequence Models, *ACS Cent. Sci.*, 2017, 3(10), 1103–1113, DOI: [10.1021/acscentsci.7b00303](#).
- 20 P. Schwaller, T. Gaudin, D. Lányi, C. Bekas and T. Laino, “Found in Translation”: Predicting Outcomes of Complex Organic Chemistry Reactions Using Neural Sequence-to-Sequence Models, *Chem. Sci.*, 2018, 9(28), 6091–6098, DOI: [10.1039/C8SC02339E](#).
- 21 P. Karpov, G. Godin and I. V. Tetko, A Transformer Model for Retrosynthesis, in *Artificial Neural Networks and Machine Learning – ICANN 2019: Workshop and Special Sessions*, ed. I. V. Tetko, V. Kůrková, P. Karpov and F. Theis, Lecture Notes in Computer Science; Springer International Publishing, Cham, 2019, pp. 817–830, DOI: [10.1007/978-3-030-30493-5\\_78](#).
- 22 P. Schwaller, T. Laino, T. Gaudin, P. Bolgar, C. A. Hunter, C. Bekas and A. A. Lee, Molecular Transformer: A Model for Uncertainty-Calibrated Chemical Reaction Prediction, *ACS Cent. Sci.*, 2019, 5(9), 1572–1583, DOI: [10.1021/acscentsci.9b00576](#).
- 23 K. Lin, Y. Xu, J. Pei and L. Lai, Automatic Retrosynthetic Route Planning Using Template-Free Models, *Chem. Sci.*, 2020, 11(12), 3355–3364, DOI: [10.1039/C9SC03666K](#).
- 24 P. Schwaller, R. Petraglia, V. Zullo, V. H. Nair, R. A. Haeuselmann, R. Pisoni, C. Bekas, A. Iuliano and T. Laino, Predicting Retrosynthetic Pathways Using Transformer-Based Models and a Hyper-Graph Exploration Strategy, *Chem. Sci.*, 2020, 11(12), 3316–3325, DOI: [10.1039/C9SC05704H](#).
- 25 S. Zheng, T. Zeng, C. Li, B. Chen, C. W. Coley, Y. Yang and R. Wu, Deep Learning Driven Biosynthetic Pathways Navigation for Natural Products with BioNavi-NP, *Nat. Commun.*, 2022, 13(1), 3342, DOI: [10.1038/s41467-022-30970-9](#).
- 26 I. V. Tetko, P. Karpov, R. Van Deursen and G. Godin, State-of-the-Art Augmented NLP Transformer Models for Direct and Single-Step Retrosynthesis, *Nat. Commun.*, 2020, 11(1), 5575, DOI: [10.1038/s41467-020-19266-y](#).
- 27 R. Irwin, S. Dimitriadis, J. He and E. J. Bjerrum, Chemformer: A Pre-Trained Transformer for Computational Chemistry, *Mach. learn.: sci. technol.*, 2022, 3(1), 015022, DOI: [10.1088/2632-2153/ac3ffb](#).
- 28 X. Wang, Y. Qian, H. Gao, W. C. Coley, Y. Mo, R. Barzilay and K. F. Jensen, Towards Efficient Discovery of Green Synthetic Pathways with Monte Carlo Tree Search and Reinforcement Learning, *Chem. Sci.*, 2020, 11(40), 10959–10972, DOI: [10.1039/D0SC04184J](#).
- 29 B. Chen, C. Li, H. Dai and L. Song, Retro\*: Learning Retrosynthetic Planning with Neural Guided A\* Search, in



- Proceedings of the 37th International Conference on Machine Learning*, PMLR, 2020, pp. 1608–1616.
- 30 A. Thakkar, A. C. Vaucher, A. Byekwaso, P. Schwaller, A. Toniato and T. Laino, Unbiasing Retrosynthesis Language Models with Disconnection Prompts, *ACS Cent. Sci.*, 2023, **9**(7), 1488–1498, DOI: [10.1021/acscentsci.3c00372](https://doi.org/10.1021/acscentsci.3c00372).
  - 31 P. Schwaller, B. Hoover, J.-L. Reymond, H. Strobelt and T. Laino, Extraction of Organic Chemistry Grammar from Unsupervised Learning of Chemical Reactions, *Sci. Adv.*, 2021, **7**(15), eabe4166, DOI: [10.1126/sciadv.abe4166](https://doi.org/10.1126/sciadv.abe4166).
  - 32 A. Byekwaso, A. C. Vaucher, P. Schwaller, A. Toniato and T. Laino, A Sequence-to-Sequence Transformer Model for Disconnection Aware Retrosynthesis, 2021, DOI: [10.26434/chemrxiv-2021-7hp1s](https://doi.org/10.26434/chemrxiv-2021-7hp1s).
  - 33 G. Landrum, *RDKit: Open-Source Cheminformatics*, 2006.
  - 34 G. Klein, Y. Kim, Y. Deng, J. Senellart and A. Rush, OpenNMT: Open-Source Toolkit for Neural Machine Translation, in *Proceedings of ACL 2017, System Demonstrations*, Association for Computational Linguistics, Vancouver, Canada, 2017, pp. 67–72.
  - 35 OpenNMT/OpenNMT-py. GitHub, <https://github.com/OpenNMT/OpenNMT-py> accessed 2020-07-28.
  - 36 C. W. Coley, L. Rogers, W. H. Green and K. F. S. C. S. Jensen, Synthetic Complexity Learned from a Reaction Corpus, *J. Chem. Inf. Model.*, 2018, **58**(2), 252–261, DOI: [10.1021/acs.jcim.7b00622](https://doi.org/10.1021/acs.jcim.7b00622).
  - 37 P. Schwaller, R. Petraglia, V. H. Nair and T. Laino *Evaluation Metrics for Single-Step Retrosynthetic Models*, Second Workshop on Machine Learning and the Physical Sciences (NeurIPS 2019), 2019.
  - 38 M. Andronov, V. Voinarovska, N. Andronova, M. Wand, D.-A. Clevert and J. Schmidhuber, Reagent Prediction with a Molecular Transformer Improves Reaction Data Quality, *Chem. Sci.*, 2023, **14**(12), 3235–3246, DOI: [10.1039/D2SC06798F](https://doi.org/10.1039/D2SC06798F).
  - 39 W. Velanguparackel, N. Hamon, J. Balzarini, C. McGuigan and A. D. Westwell, Synthesis, Anti-HIV and Cytostatic Evaluation of 3'-Deoxy-3'-Fluorothymidine (FLT) pro-Nucleotides, *Bioorg. Med. Chem. Lett.*, 2014, **24**(10), 2240–2243, DOI: [10.1016/j.bmcl.2014.03.092](https://doi.org/10.1016/j.bmcl.2014.03.092).
  - 40 T. Wang, Y. Ueda, Z. Zhang, Z. Yin, J. Matiskella, B. C. Pearce, Z. Yang, M. Zheng, D. D. Parker, G. A. Yamanaka, Y.-F. Gong, H.-T. Ho, R. J. Colonno, D. R. Langley, P.-F. Lin, N. A. Meanwell and J. F. Kadow, Discovery of the Human Immunodeficiency Virus Type 1 (HIV-1) Attachment Inhibitor Tamsavir and Its Phosphonoxyethyl Prodrug Postemsavir, *J. Med. Chem.*, 2018, **61**(14), 6308–6327, DOI: [10.1021/acs.jmedchem.8b00759](https://doi.org/10.1021/acs.jmedchem.8b00759).
  - 41 F. L. Scott, B. Clemons, J. Brooks, E. Brahmachary, R. Powell, H. Dedman, H. G. Desale, G. A. Timony, E. Martinborough, H. Rosen, E. Roberts, M. F. Boehm and R. J. Peach, Ozanimod (RPC1063) Is a Potent Sphingosine-1-Phosphate Receptor-1 (S1P1) and Receptor-5 (S1P5) Agonist with Autoimmune Disease-Modifying Activity, *Br. J. Pharmacol.*, 2016, **173**(11), 1778–1792, DOI: [10.1111/bph.13476](https://doi.org/10.1111/bph.13476).
  - 42 A. C. Flick, C. A. Leverett, H. X. Ding, E. L. McInturff, S. J. Fink, S. Mahapatra, D. W. Carney, E. A. Lindsey, J. C. DeForest, S. P. France, S. Bertritt, S. V. Bigi-Butterill, T. S. Gibson, R. B. Watson, Y. Liu and C. J. O'Donnell, Synthetic Approaches to the New Drugs Approved During 2020, *J. Med. Chem.*, 2022, **65**(14), 9607–9661, DOI: [10.1021/acs.jmedchem.2c00710](https://doi.org/10.1021/acs.jmedchem.2c00710).
  - 43 D. Probst and J.-L. Reymond, Visualization of Very Large High-Dimensional Data Sets as Minimum Spanning Trees, *J. Cheminf.*, 2020, **12**(1), 12, DOI: [10.1186/s13321-020-0416-x](https://doi.org/10.1186/s13321-020-0416-x).
  - 44 D. Probst, P. Schwaller and J.-L. Reymond, Reaction Classification and Yield Prediction Using the Differential Reaction Fingerprint DRFP, *Digital Discovery*, 2022, **1**(2), 91–97, DOI: [10.1039/D1DD00006C](https://doi.org/10.1039/D1DD00006C).
  - 45 IBM RXN for Chemistry, <https://rxn.res.ibm.com> accessed 2022-09-05.
  - 46 S. Genheden and E. PaR. Bjerrum, Towards a Framework for Benchmarking Retrosynthesis Route Predictions, *Digital Discovery*, 2022, **1**(4), 527–539, DOI: [10.1039/D2DD00015F](https://doi.org/10.1039/D2DD00015F).
  - 47 G. Pesciullesi, P. Schwaller, T. Laino and J.-L. Reymond, Transfer Learning Enables the Molecular Transformer to Predict Regio- and Stereoselective Reactions on Carbohydrates, *Nat. Commun.*, 2020, **11**(1), 4874, DOI: [10.1038/s41467-020-18671-7](https://doi.org/10.1038/s41467-020-18671-7).
  - 48 W. Finnigan, L. J. Hepworth, S. L. Flitsch and N. J. Turner, RetroBioCat as a Computer-Aided Synthesis Planning Tool for Biocatalytic Reactions and Cascades, *Nat. Catal.*, 2021, **4**(2), 98–104, DOI: [10.1038/s41929-020-00556-z](https://doi.org/10.1038/s41929-020-00556-z).
  - 49 D. Kreutter, P. Schwaller and J.-L. Reymond, Predicting Enzymatic Reactions with a Molecular Transformer, *Chem. Sci.*, 2021, **12**(25), 8648–8659, DOI: [10.1039/D1SC02362D](https://doi.org/10.1039/D1SC02362D).
  - 50 D. Probst, M. Manica, Y. G. Nana Teukam, A. Castrogiovanni, F. Paratore and T. Laino, Biocatalysed Synthesis Planning Using Data-Driven Learning, *Nat. Commun.*, 2022, **13**(1), 964, DOI: [10.1038/s41467-022-28536-w](https://doi.org/10.1038/s41467-022-28536-w).
  - 51 A. Thakkar, A. Vaucher, A. Byekwaso, P. Schwaller, A. Toniato and T. Laino, Disconnection Labelled Reaction Data, *Zenodo*, 2022, DOI: [10.5281/zenodo.7101695](https://doi.org/10.5281/zenodo.7101695).

