

ASKCOS: an open source software suite for synthesis planning

Zhengkai Tu¹, Sourabh J. Choure², Mun Hong Fong²,
Jihye Roh², Itai Levin³, Kevin Yu⁴, Joonyoung F. Joung²,
Nathan Morgan², Shih-Cheng Li², Xiaoqi Sun², Huiqian Lin²,
Mark Murnin², Jordan P. Liles², Thomas J. Struble⁵,
Michael E. Fortunato⁶, Mengjie Liu^{2,†}, William H. Green²,
Klavs F. Jensen², Connor W. Coley^{1,2*}

¹Department of Electrical Engineering and Computer Science,
Massachusetts Institute of Technology, 77 Massachusetts Ave,
Cambridge, MA, 02139, USA.

²Department of Chemical Engineering, Massachusetts Institute of
Technology, 77 Massachusetts Ave, Cambridge, MA, 02139, USA.

³Department of Biological Engineering, Massachusetts Institute of
Technology, 77 Massachusetts Ave, Cambridge, MA, 02139, USA.

⁴Center for Computational Science and Engineering, Massachusetts
Institute of Technology, 77 Massachusetts Ave, Cambridge, MA, 02139,
USA.

⁵Bristol Myers Squibb, 250 Water Street, Cambridge, MA, 02141, USA.

⁶Novartis Institutes for BioMedical Research, Inc., 250 Massachusetts
Avenue, Cambridge, MA, 02139, USA.

*Corresponding author(s). E-mail(s): ccoley@mit.edu;
Contributing authors: ztu@mit.edu; sjchoure@mit.edu;
fong410@mit.edu; jroh99@mit.edu; itail@mit.edu; kyu3@mit.edu;
jjoung@mit.edu; knathan@mit.edu; scli@mit.edu; xiaoqis@mit.edu;
linhq@mit.edu; murninm@mit.edu; jliles24@mit.edu;
Thomas.Struble@bms.com; mike.fortunato@novartis.com;
mjliu@mit.edu; whgreen@mit.edu; kfjensen@mit.edu;

[†]Current affiliation: AstraZeneca. Work done while at MIT.

Abstract

The advancement of machine learning and the availability of large-scale reaction datasets have accelerated the development of data-driven models for computer-aided synthesis planning (CASP) in the past decade. Here, we detail the newest version of ASKCOS, an open source software suite for synthesis planning that makes available several research advances in a freely available, practical tool. Four one-step retrosynthesis models form the basis of both interactive planning and automatic planning modes. Retrosynthetic planning is complemented by other modules for feasibility assessment and pathway evaluation, including reaction condition recommendation, reaction outcome prediction, and auxiliary capabilities such as solubility prediction and quantum mechanical descriptor prediction. ASKCOS has assisted hundreds of medicinal, synthetic, and process chemists in their day-to-day tasks, complementing expert decision making. It is our belief that CASP tools like ASKCOS are an important part of modern chemistry research, and that they offer ever-increasing utility and accessibility.

Keywords: synthesis planning, open source software, data-driven predictive chemistry

1 Introduction

Synthesis planning describes a broad category of approaches for selecting experimental pathways and procedures during target-oriented synthesis. While planning a synthesis campaign may require significant chemistry expertise and benefits from the years of training that many experienced chemists undergo, the well-defined yet combinatorially complex nature of synthesis planning renders this task particularly amenable to algorithmic reasoning. Formally, computer-aided synthesis planning (CASP) integrates a variety of computational methodologies that assist chemists with different tasks in this process, including identifying viable synthetic routes through retrosynthetic analysis, recommending reaction conditions, and predicting reaction outcomes.

Since the 1960s, chemists have sought to encode the rules of organic synthesis into automated computational systems [1, 2]. Early CASP tools generally relied on expert-encoded reaction rules and heuristics for making suggestions. In particular, for one-step retrosynthetic analysis [3], expert systems such as LHASA [4] and SECS [5] made use of reaction templates that encode chemical reaction rules based on molecular pattern matching; AIPHOS [6] and WODCA [7] were among the first to combine retrosynthesis, reaction condition suggestion, and product prediction into an integrated system. More recent advancements include tools like Chematica (now Synthia) [8], which leverages modern computational capabilities alongside expert-curated rules and heuristics to propose transformations, generating synthetic pathways for complex molecules that have been successfully implemented in the laboratory [9, 10].

With the development of machine learning and the availability of reaction datasets containing millions of entries, there has been renewed interest in CASP with data-driven approaches [11, 12]. Many data-driven models for one-step retrosynthesis have been developed with formulations based on reaction template prediction [13, 14] or retrieval [15, 16], machine translation [17–19], graph edit prediction [20, 21], as well

as other generative models [22, 23]. These one-step predictors have been integrated with various tree-search algorithms to navigate through the network of hypothetical reactions to identify synthetic pathways in which all starting materials are accessible [24–28]. Machine learning models have also been applied to other elements of synthesis planning such as reaction condition recommendation [29–33] and reaction outcome prediction [17, 20, 34–37]. Like in typical supervised learning settings, these approaches are trained on historical reaction data and try to generalize to unseen targets.

Parallel to the development of *algorithms* for CASP has been the emergence of many software tools. Among many proprietary examples are LHASA [4], Synthia [8], Chemical.ai [38], IBM’s RXN for Chemistry [39], Spaya by Iktos [40], Manifold by PostEra [41], Molecule.one [42], SciFinder (integrating technology from ChemPlanner [43] and Molecule.one), and Reaxys (integrating technology from Iktos and Pending.AI [44]; these are complemented by ASKCOS [45], AiZynthFinder [46], and Syntheseus [47] as open source offerings. ASKCOS is distinct in its breadth and attempt to cover a wide range of different tasks in synthesis planning, rather than focusing mostly on retrosynthetic analysis.

The ASKCOS software has been in development since 2016. It was originally designed to use algorithmically-extracted templates and simple tree search algorithms such as depth-first search or best-first search using heuristics. Over the past eight years, it has been expanded significantly, deployed, used, and evaluated by the community and, in particular, at dozens of pharmaceutical and chemical companies within and beyond the Machine Learning for Pharmaceutical Discovery and Synthesis (MLPDS) consortium [48]. Chemists have integrated various modules of ASKCOS into their workflows for molecule and route ideation, at times adapting individual components into proprietary design tools [49]. While ASKCOS has proven useful for many chemists and researchers [45, 49–60], it has only been formally (and briefly) described in the context of robotic flow chemistry in Coley et al. [45].

Here, we report the latest version of the open source ASKCOS software suite, reflecting the variety of new functionalities and improvements that have contributed to its growth. At a high level, ASKCOS has two modes of operation for its core retrosynthesis functionality: an interactive mode with the Interactive Path Planner (IPP) and an automatic mode with the Tree Builder. Both modes can use multiple one-step strategies individually or simultaneously to combine the strengths of each. In addition to the retrosynthesis functionality, ASKCOS provides partial solutions to many adjacent tasks within synthesis planning. The open source and modular nature of ASKCOS also allows for easy access to various prediction modules, either via the web-based interface or via application programming interfaces (APIs), which are available under permissive MIT licenses. In the rest of this article, we present these functionalities in terms of both their scientific aspects and their usage in ASKCOS (Figure 1).

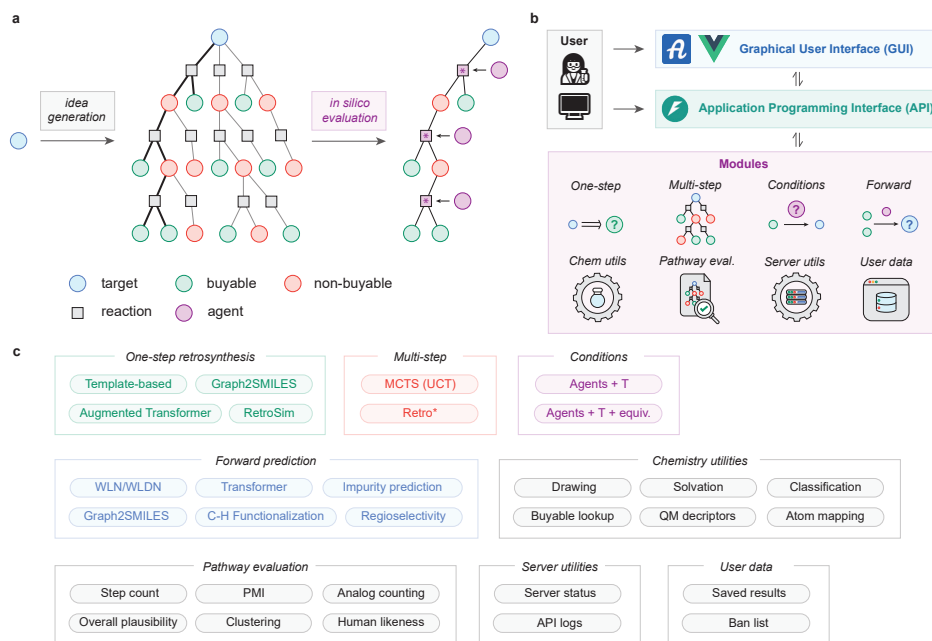


Fig. 1 ASKCOS overview. **a)** The typical target-oriented synthesis workflow for which ASKCOS is designed. A target molecule (blue circle) is recursively expanded retrosynthetically until buyable starting materials (green circles) are reached. Agents/conditions can be predicted for each proposed reaction, which can then undergo further evaluation, e.g., to anticipate reaction products. **b)** High-level flowchart for ASKCOS usage. A user can either interact via the graphical user interface or programmatic endpoints, which call various prediction modules. **c)** Summary of modules available in ASKCOS, which fall under themes including one-step retrosynthesis (green), multi-step search (red), condition recommendation (purple), reaction outcome prediction (blue), as well as utilities and supplementary capabilities (black). The modularity of the software design enables the straightforward extension of functionality (e.g., to new data-driven models) as they are developed in a research setting and made production-ready.

2 Results

2.1 One-step retrosynthetic expansion

The one-step retrosynthetic expansion engine lies at the core of retrosynthetic analysis in ASKCOS. Given a target molecule, a list of candidate precursors is first predicted with user-specified one-step model(s). A fast plausibility filter based on the in-scope filter described by Segler et al. [25] removes unlikely precursors, and the remaining list is reranked based on buyability and complexity (e.g., by heavy atom count, ring count, or by SCScore [61]). Thereafter, the candidates undergo a series of optional post-processing steps. Because several predictions from the list of candidate precursors may correspond to highly similar strategies (e.g., differing only by leaving groups), precursors can be clustered based on their structures or reaction classes to return a diverse set of suggestions to the user; atom mapping and template extraction can

be performed on reactions proposed by template-free model(s), which help check the chemical validity of the suggestions; precursors can be filtered or grouped based on which atoms are involved in the reaction centers; selectivity checks can automatically flag reactions with potential regioselectivity issues.

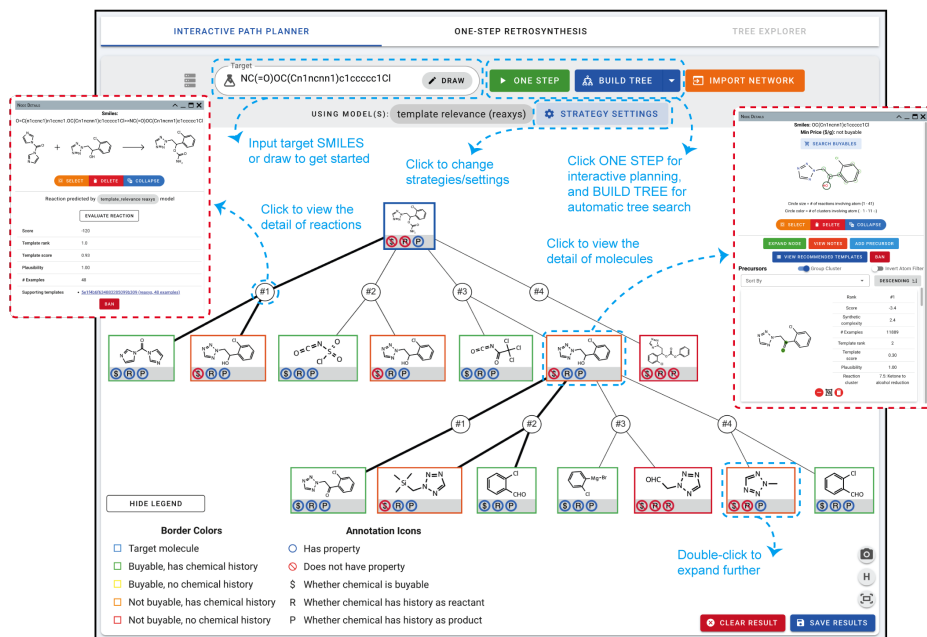


Fig. 2 Annotated screenshot of the interactive planning canvas in ASKCOS. The input target molecule shown is cenobamate, a drug for partial-onset seizures, defined by entering the SMILES string NC(=O)OC(Cn1ncnn1)c1ccccc1Cl in the search bar or using an external name resolver if permitted by security policies. An interactive retrosynthetic search can be initiated with the **ONE STEP** button. An automated multi-step retrosynthetic search can be initiated with the **BUILD TREE** button. Search settings for both modes can be set with the **STRATEGY SETTINGS** button. Here, results are shown for an interactive search performed using the template-based single step model trained on the Reaxys dataset. The top 4 suggestions by the model are displayed on the canvas. A legend for the frame colors is displayed at the bottom left of the page. The displayed two-tier tree was generated by following the initial single step expansion of the target molecule with the expansion of one of the suggested precursors. Additional retrosynthetic expansions can be performed by double clicking on a chemical node. The chemical node info window provides users with additional information and actions (e.g., adding a comment for that chemical or banning it from future suggestions). If a one-step model has already been used to suggest precursors for that chemical, the node info window displays (right inset, red dashed frame) suggested precursors and allows users to sort or filter these suggestions using a variety of metrics (e.g., score, synthetic complexity, number of rings). A reaction node info window (left inset, red dashed frame) provides users with information about the suggested reaction, including predicted scores and database precedents if they exist. The window also allows users to perform actions on the node such as removing or hiding it. Results can be saved, exported, or reorganized into different graphical views using additional tools in the bottom-right of the canvas.

2.2 Interactive planning with the template relevance model

The use of reaction templates to suggest retrosynthetic disconnections has remained a popular choice since its origination in the early CASP tools of the 1960s [1, 4, 5]. Template-based models within ASKCOS follow the neural-symbolic approach of Segler and Waller [62] wherein a policy network is trained to rank which templates appear most *strategic* and *chemically plausible* given the target. ASKCOS contains a variety of such models trained to use template sets derived from reactions in Pistachio [63], CAS Content [64], USPTO [65], and Reaxys [66] using RDChiral [67] (see Section [Methods](#)). Specialized models trained on enzymatic reaction data in BKMS [68] and a specialized "ring-breaker" Pistachio model are also available and described in Levin et al. [50] and Thakkar et al. [69], respectively. Because template-based models propose candidate precursors using templates extracted from published reactions, model suggestions can be traced back to reaction precedents and are therefore somewhat explainable.

Figure 2 shows a sample planning step with the original template-based model trained on Reaxys in 2016 in the *Interactive Path Planner (IPP)*. The target is specified by a simplified molecular input line entry system (SMILES) string [70] but can also be defined by drawing its structure with Ketcher or by common name using the PubChem API [71]. One-step expansion can then be triggered by pressing the green button. The top few (5 by default) suggestions will be added to the canvas as the child nodes of the target, with circles representing the reaction nodes and rectangles representing the molecule nodes. Clicking on the nodes display more context-specific details (shown as inset figures in red dotted lines in Figure 2) and provide access to additional features, such as banning the reaction/molecule (thereby preventing them from appearing in future searches) or deleting/collapsing all child nodes of that node. Reaction nodes report the template scores (i.e., the template probabilities returned by the model), plausibility (evaluated by the fast binary filter as mentioned in Section [One-step retrosynthetic expansion](#)), and links to template details (which include links to the reaction precedents that the templates were extracted from).

Each reaction can be analyzed further with the **EVALUATE REACTION** button, such as recommending conditions for the reaction or predicting reaction outcomes (see Sections [Reaction condition recommendation](#) and [Reaction outcome prediction](#) for details). Molecule nodes report their price and are additionally color-coded to reflect whether they are buyable or appear in known reactions stored in a reference database, as explained in the bottom-left IPP legend. Molecule node-specific features are also available, including expanding the nodes, adding and saving notes, and viewing recommended templates for the molecule. Once molecule nodes are expanded, their node details additionally list *all* predicted precursors. These precursors can then be sorted according to different criteria (e.g., heuristic scores, synthetic complexity, number of precedents), grouped into clusters (if already clustered as mentioned in Section [One-step retrosynthetic expansion](#)), or filtered by reaction center by selecting the green-circled atom in the rendering of the molecule. Specific precursors can be added or removed from the canvas with the green + button or the red - button, respectively.

After this first expansion of the target molecule, it is up to the user to decide which molecule node(s) to expand further, hence the *interactive* nature of this view. For instance, users can choose a non-buyable molecule to expand next and continue

this process until an appropriate pathway is identified. Precursors can be added to the molecule manually with the `ADD PRECURSOR` button if a user wishes to supplement model predictions with their own ideas. The explored network and notes can be saved to a user’s profile for future access or exported as a JSON file for offline processing.

2.3 Interactive planning with multiple models, including template-free models

There is an abundance of one-step retrosynthetic models reported in the past several years that forgo the use of templates and instead learn to predict reactant structures from product structures in a more flexible end-to-end manner. ASKCOS currently contains four categories of one-step strategies, including Transformer [18, 34, 72], Graph2SMILES [17], Retrosim [15], and the aforementioned template relevance strategy [62]. Transformer and Graph2SMILES use template-free approaches and model retrosynthesis as SMILES-to-SMILES and graph-to-SMILES translation tasks, respectively. Retrosim offers a retrieval-based, learning-free approach in which reactions are suggested based on analogy to most-similar precedents. Each model may exhibit distinct strengths and failure modes not reflected in their quantitative performance on standard benchmark tasks (Table 1). For this reason, ASKCOS supports the consolidation of recommendations from multiple strategies; when multiple models recommend the same precursors, this can be interpreted as a sign of confidence in that recommendation.

Mixing and matching different one-step strategies is performed within the `STRATEGY SETTINGS` menu (Figure 2). Different strategies are queried in sequence and combined results are de-duplicated before appearing on the page. Each strategy has its own settings, such as the maximum number of templates to be applied for the template relevance model and the training set to use. When viewing node details, suggested reactions from template-free and Retrosim models contain different metadata (e.g., Retrosim results display the reaction precedents if that information is included on deployment). When multiple models predict the same precursor, the metadata from all models are merged to show all template and/or reaction precedent information where applicable. Most of these strategies can be retrained using new reaction databases, e.g., proprietary collections from an internal electronic lab notebook system, which may provide better coverage of different reaction and substrate types.

2.4 Automatic multi-step planning with the Tree Builder

In addition to interactive planning where users guide the selection of which molecule nodes to expand, *automatic* planning can be more convenient, particularly when there are many target molecules of interest. Retrosynthetic searches can be run for thousands of targets through asynchronous requests, albeit at the expense of losing explicit control on the direction of expansion. Formally, automatic multi-step planning has been formulated as tree search or graph search problems. In each iteration of the search, a molecule is selected for one-step retrosynthetic expansion. New hypothetical reactions and their corresponding reactants are added to the search tree. This process is repeated until some termination criterion is reached, for example, until a synthetic pathway is

Table 1 Previously reported top-k accuracies (%) of one-step models on two commonly-used benchmark datasets: USPTO-50k [73] and USPTO-full [13].

Model name	Type	USPTO-50k			USPTO-full		
		Top 1	Top 10	Ref.	Top 1	Top 10	Ref.
Template relevance	Template-based	45.2	83.5	[74]	35.8	60.8	[13]
Retrosim	Template-based	37.3	74.1	[15]	32.8	56.1	[13]
Transformer w/o aug.	Template-free	43.1	78.7	[72]	42.9	66.8	[75]
Transformer w/ aug.	Template-free	53.2	85.2	[18]	44.4	73.3	[18]
Graph2SMILES	Template-free	52.9	79.5	[17]	45.7	63.4	[17]

Note: for the template relevance model, accuracies on the original NeuralSym model are reported. Best values from multiple references are recorded. Transformer w/ aug. refers to the variant where equivalent reaction SMILES are used to *augment* the training set, which is a known empirical technique for boosting accuracy [76]. Due to the vast amount of literature related to Transformer and SMILES augmentation, we did not attempt to do an exhaustive literature survey for Transformer-based models.

found in which all starting materials are buyable. The well-known baseline algorithms include Monte Carlo tree search (MCTS) [25, 72], Proof Number Search (PNS) [24, 26] and A* Search [28] which tend to be inspired by other artificial intelligence applications such as AlphaGo [77]. More recent work in the field on such search algorithms has mainly focused on improvement of selection policies, for example, with reinforcement learning [78–81], supervised learning [82–84], or chemical heuristics [27, 85]. A number of studies have also sought to improve the expansion policy by focusing on the one-step model in the context of multi-step planning [86, 87]. Others have approached synthesis planning outside of the typical tree search formulation, such as through sequence generation [88] or as a bidirectional search [89].

ASKCOS supports automatic planning through the *Tree Builder* module which currently provides two search algorithms, MCTS [25] and Retro* [28]. Either can be selected from STRATEGY SETTINGS. After specifying the target in the IPP canvas, the user can click the BUILD TREE button in Figure 2 to initiate a Tree Builder job, which will run asynchronously in the background. A pop-up window will inform the user when the results are ready and available for viewing. All tree results can be found under the My Results page as shown on the top left corner of Figure 3, where the status of any result entry will change from **started** to **completed** once the Tree Builder job is finished. From here, the user can visualize the result either as a full retrosynthesis tree in the IPP canvas, or as individual routes in the Tree Explorer (see Section Pathway scoring and ranking for further analysis). The settings used for the job can be inspected, and the result can be shared via the green sharing button in Figure 3.

2.5 Reaction condition recommendation

The prediction of reaction conditions, including the identity of agents (catalysts, reagents, solvents) and operating conditions such as temperature and equivalence ratios, is an often overlooked aspect of synthesis planning. Condition recommendation is essential for any suggested reactions to be experimentally validated, and reaction

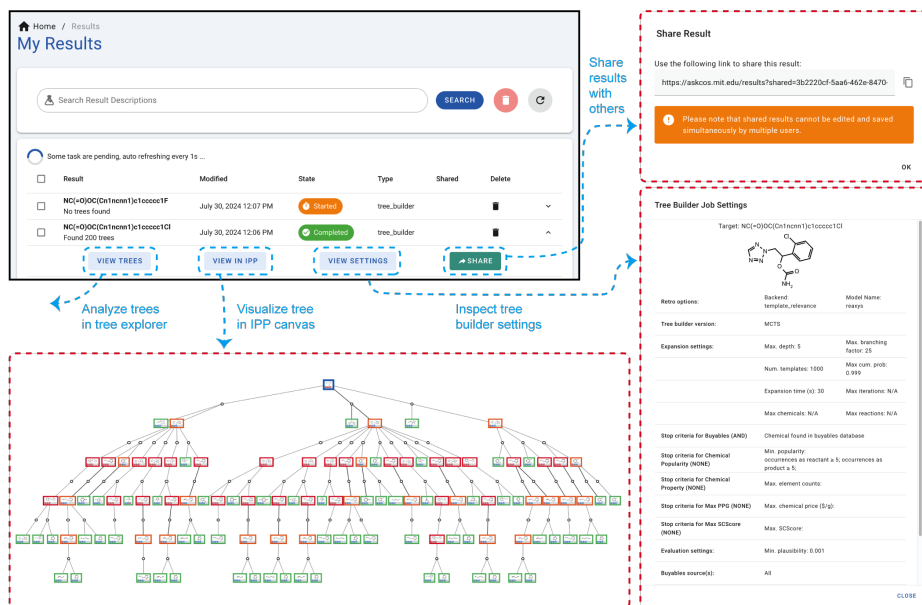


Fig. 3 Annotated screenshot of tree search results in ASKCOS. The Tree Builder job can be initiated with the BUILD TREE button on the interactive path planning page. In My Results page, the status for the Tree Builder job is shown as Started. Once the job is complete, the status changes to Completed, and the trees found can be visualized/analyzed in the Tree Explorer (VIEW TREES) or the IPP canvas (VIEW IN IPP). The settings used for the Tree Builder job can be viewed with the VIEW SETTINGS button, and the results can be shared with others with the SHARE button. Here, the Tree Builder job for the target molecule NC(=O)OC(Cn1ncnn1)c1cccc1Cl is performed using the MCTS algorithm with the template-based single step model trained on the Reaxys dataset, using a maximum depth of 5, a maximum branching factor of 25, and an expansion time of 30s (right bottom inset, red dashed frame). The results can be visualized in the IPP canvas (left bottom inset, red dashed frame).

outcomes are strongly influenced by reaction conditions. Reaction type specific data-driven models have been built, for example, to predict classes of solvent and catalyst for Michael additions [90], ligands for Pd-catalyzed C-N coupling [91], and reaction contexts including temperature and pressure for four families of substrate-specific cross-coupling reactions [30, 31]. Global models which are not specific to reaction families have also been developed using a variety of approaches, including multi-class classifier chains [29] and its variant using a Transformer encoder with SMILES inputs [92], retrieval-augmented prediction using text descriptions of similar reactions [32], and a two-stage model with candidate generation and ranking as distinct steps [33].

ASKCOS includes a data-driven condition recommendation model based on Gao et al. [29] as well as a second version being developed in ongoing work which additionally predicts equivalence ratios. This model formulates condition recommendation in terms of four sub-problems (1) predicting agent identities as multi-label classification; (2) predicting temperature as binned classification; (3) predicting reactant equivalence

CONDITION RECOMMENDATION PRODUCT PREDICTION IMPURITY PREDICTION REGIOSELECTIVITY PREDICTION AROMATIC C-H FUNCTIONALIZATION

Reactants O=C=NC(=O)C(Cl)(Cl)Cl.ClC1OC(Cn1ncnn1)c1ccccc1Cl Product NC(=O)OC(Cn1ncnn1)c1ccccc1Cl

Click to change models and settings

Reactants (Amount) Reagents (Amount) Temperature Predict with conditions

1 O=C=NC(=O)C(Cl)(Cl)Cl (1.24) ClC1OC(Cn1ncnn1)c1ccccc1Cl (1.00) C1CCOC1 (100.51) -3 °C

Shortcut to run product prediction

CONDITION RECOMMENDATION PRODUCT PREDICTION IMPURITY PREDICTION REGIOSELECTIVITY PREDICTION AROMATIC C-H FUNCTIONALIZATION

Reactants O=C=NC(=O)C(Cl)(Cl)Cl.ClC1OC(Cn1ncnn1)c1ccccc1Cl

Agents C1CCOC1 Solvent

Click to change models and settings

Rank	Product	Probability	Max. Score	Molecular Weight	Predict Impurities	Predict regio-selectivities
1		0.9726	-17.504	267.1	<input type="button" value="→"/>	<input type="button" value="→"/>
2		0.0273	-21.077	410.9	<input type="button" value="→"/>	<input type="button" value="→"/>
3		0.0001	-27.052	268.1	<input type="button" value="→"/>	<input type="button" value="→"/>

Fig. 4 Annotated screenshot of the condition recommender (top) and the forward predictor (bottom) in ASKCOS. In the condition recommender, the SMILES strings of reactants and products are input in the **Reactants** and **Products** panels. Condition recommendation can be initiated with the **GET RESULTS** button. Prediction models and their training set can be set with the **MODEL** and **SETTINGS** buttons, respectively. Each proposed condition has shortcuts to run product prediction. In the forward predictor, the SMILES strings of reactants and agents are pre-populated using the shortcut from the condition recommendation page. Product prediction can be initiated with the **GET RESULTS** button. Prediction models and their training sets are chosen via **SETTINGS** button. Here, the results are shown for three predicted products with their probabilities, which should be interpreted only qualitatively. Each prediction has shortcuts to run impurity and regioselectivity prediction as additional evaluations.

ratios as multi-target regression; and (4) predicting agent equivalence ratios as multi-target regression. A screenshot of the condition recommendation page in ASKCOS is shown at the top of Figure 4. As with other pages, the reactants and product can be specified either by SMILES strings or by drawing. After generating model predictions by clicking on `GET RESULTS`, several different condition settings are displayed in a tabular format. The general layout of the reaction condition recommendation page is used across several other modules for consistency: the top panel is for the inputs and the bottom panel is for the prediction results; the model and settings can be modified by clicking on the blue buttons for `MODEL` and `SETTINGS`.

2.6 Reaction outcome prediction

Another component of synthesis planning beyond retrosynthesis is the prediction of reaction outcome(s). In our workflows, these predictions mostly serve to identify chemically infeasible or unfavorable reactions, which the user can choose to prune from the synthesis tree in an interactive setting. Many data-driven approaches simplify this task as predicting the identity of the major product, making it equivalent to a molecule-to-molecule transformation. Similar to one-step retrosynthesis, some studies have formulated it as forward template prediction [93, 94], whereas later developments have been dominated by graph-edit based [20, 35], electron-flow based [36, 37], and translation-based [17–19, 34] template-free methods. Outcome prediction can answer more fine-grained questions about reaction outcomes, such as the site selectivity of aromatic C-H functionalization reactions [95] or other situations where multiple regioisomers appear possible based on reaction templates [96]. Another application of outcome prediction is the analysis of potential impurities. When impurities are defined as the minor products of the main reactions or the products of side reactions, they can be predicted by considering several predicted reaction outcomes with lower probabilities or by predicting the outcomes of new reactant sets containing a product of the original reactions (i.e., anticipating potential over-reaction).

The major product prediction page in ASKCOS is shown in the bottom of Figure 4. In this specific example, the top outcome is carbamate formation arising from nucleophilic attack of the alcohol into the isocyanate (followed by hydrolysis [97]) with a probability of 0.9726. Non-hydrolyzed product and nucleophilic attack by a tetrazole nitrogen are predicted to be the second and third most likely products, but with lower probabilities of 0.0273 and 0.0001, respectively. Past studies have shown these probabilities to correlate with accuracy *on average* but may not be a robust measure of confidence for a particular result [34, 35, 98]. Modules for impurity prediction, regioselectivity, and C-H site-selectivity are accessible via other tabs on the same page. The impurity prediction module relies on the major product predictor and considers minor products, over-reaction, dimerization, solvent adducts, and subsets of reactants; the details of regio- and site-selectivity predictions are reported in Guan et al. [96] and Struble et al. [95], respectively.

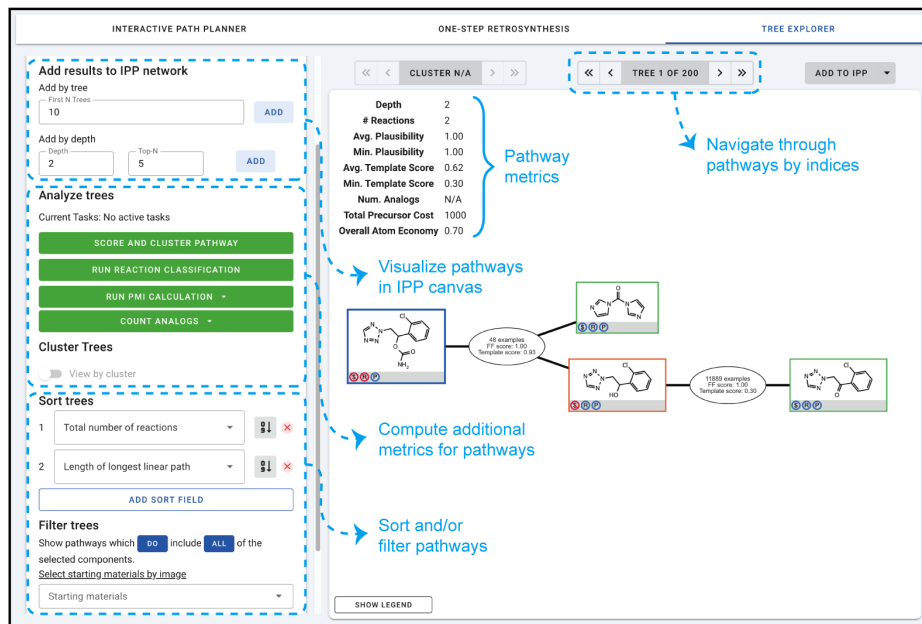


Fig. 5 Annotated screenshot of the Tree Explorer in ASKCOS. One of 200 routes returned for the target molecule cenobamate is displayed in the main panel (right) of the window. Pathway metrics for the route are shown in the top left corner of the main panel. The left panel provides users with additional actions to perform on the returned synthetic routes. Users can transfer information from the Tree Builder to the Interactive Path Planner to continue exploring retrosynthetic suggestions beyond what was returned (top); users can initiate longer-running pathway-level scoring calculations (middle); finally, users can sort or filter the discovered routes based on various calculated metrics or the presence or absence of specific starting materials and intermediates (bottom).

2.7 Pathway scoring and ranking

When many putative synthetic pathways are found for a target molecule, a new challenge arises: to identify the routes that best satisfy a *chemist's* goals, not merely to identify *any* route. It is impractical to triage routes manually when the number of suggestions becomes too large. Strategically similar pathways can first be clustered or grouped based on trained pathway embeddings [99] or the reactions types they involve (e.g., using NameRxn [100] categories or approximations thereof). Pathway-level evaluations then help prioritize promising synthetic routes, though there are many criteria by which a synthetic pathway could be judged [101]. Simple, readily-calculable metrics include step count, longest linear sequence, atom economy, cost of starting materials, or diversifiability based on an estimate of the size of an analog space achieved through building block enumeration [102]. More complex metrics that may be derived from other predictive models include estimates of human likeness [99], overall perceived likelihood of feasibility, pathway greenness based on solvent usage [81],

and process mass intensity (PMI). Assessing purification/isolation requirements in a robust manner remains elusive [103].

Pathway evaluation in ASKCOS is intended to be performed on pathways returned by the Tree Builder and is accessible from the results page by clicking **VIEW TREES** as shown in Figure 3. Since none of the aforementioned evaluation criteria is perfect or sufficient on its own, ASKCOS makes all of them available as part of the *Tree Explorer* (Figure 5), with operations and options on the left panel and the canvas on the right for displaying the pathway with its metrics. In this specific example, 200 pathways have been returned (the default maximum); the first pathway is shown along with automatically-calculated metrics such as the depth, the average plausibility, and the analog count of the pathway. The left panel is organized into three sections that allow the user to visualize multiple best-ranked pathways in the IPP canvas; calculate additional evaluation metrics on-request (rather than automatically due to their computational cost); and sort or filter pathways, e.g., based on certain starting materials of interest.

2.8 Utilities and supplementary predictive models

Beyond these “core” synthesis planning capabilities, ASKCOS contains additional complementary tools. These include basic drawing functionality (**Drawing**) and buyable building block search by SMILES or SMARTS (**Buyable Look-up**). The latter makes use of a predefined commercial catalog that is easily customizable when ASKCOS is deployed. The utilities page also offers two additional machine learning models for solvation prediction (**Solubility Prediction** and **Solvent Screening**) [104] and the prediction of atom- and bond-level descriptors calculated by DFT (**QM Descriptor**) [96, 105]. Solubility prediction is part of a long-term goal of improving the relevance of CASP for process chemistry [106] as it helps guide the selection of solvents for reactions, liquid-liquid extractions, or crystallizations. ML-estimated QM descriptors can be used as features by other predictive models [107] or standalone, e.g., for human assessment of selectivity.

The pages for solubility prediction, solvent screening, and QM feature prediction are organized under the **Utilities** tab, each with a standard layout as in Figure 4. Solubility prediction requires a solute, solvent, and temperature as inputs; solvent screening expects a single solute, a list of solvents, and a list of temperatures; QM prediction needs the SMILES of the target molecule. The screenshots of these pages are shown in Supplementary Figure S1, S2 and S3, with more detailed explanation on the theoretical aspect and usage in Supplementary Sections [Details of solubility prediction and solvent screening](#) and [Details of QM descriptor prediction](#).

3 Discussion

As a broad, extensible software suite for synthesis planning, ASKCOS has been adopted by various organizations within and beyond the context of the MLPDS consortium. While not all usage of ASKCOS is publicly described, several use cases where ASKCOS has aided chemists’ workflows have been discussed in the 2020 review by Struble et al. [49]. Particularly well-received features include the interactive planning

mode when automatic tree search fails, as well as the possibility for retrosynthetic suggestions to be linked back to literature precedents. ASKCOS has helped serve as a foundation for components of AiZynthFinder [46, 108] (e.g., in its template extraction strategy). Discovery chemists from Janssen have made use of various modules in ASKCOS at the lead optimization stage. In particular, the one-step retrosynthesis API helped narrow a library of 222k alcohols to 15.7k on the basis of synthesizability [58]. Pfizer has reported the use of ASKCOS to augment human ideas in their internal graph databases [52], and Syngenta has incorporated ASKCOS as one of several idea generation tools when comparing synthetic routes [60]. Beyond the industrial setting, various researchers have used ASKCOS to propose synthetic routes for candidate protease inhibitors [51] and other potential anti-COVID-19 drugs [59].

As is true of other computational tools, most if not all of the functionalities in ASKCOS aim to *assist* and not replace expert chemists. Ultimately, these models are influenced by the data on which they are trained and may recapitulate popular patterns and trends in data (albeit in useful ways) without understanding the underlying physical sciences. The interpretability of our models fall on a wide spectrum. For the one-step retrosynthesis models, for example, template-based approaches ensure traceability to literature precedents. In contrast, translation-based approaches make predictions in a more black-box manner, which can be creative but to the extent of being alchemical, e.g., by inappropriately adding atoms in the generated SMILES. We refer the reader to individual manuscripts for each prediction module for more in-depth discussions on their strengths and limitations.

As an illustration of how ASKCOS can be used as an assistive tool, we conduct a synthesis planning exercise, which is described further in Supplementary Section [Illustration of applying ASKCOS to FDA-approved small molecule drugs from 2019 to 2023](#). We start with automatic retrosynthetic planning using the Tree Builder for all targets. Using typical search settings with template relevance models trained on Reaxys and Pistachio, the default buyable database (consisting of a few hundred thousand molecules from eMolecules [109], Sigma Aldrich [110], LabNetwork [111], Mcule [112], and ChemBridge [113]), and a limit on the number of chemical nodes in the search tree of 5000, hypothetical retrosynthetic pathways are found for many of the targets. Sample shortest routes are presented in Supplementary Figures [S4](#), [S5](#), [S6](#), and [S7](#) *exactly as returned* along with the top-1 proposed conditions from the V1 condition recommender. We then demonstrate the use of additional modules in ASKCOS for further analysis when proposed steps are counter-intuitive or appear implausible, e.g., by cross-referencing literature precedents or examining lower-ranked suggested conditions. Thereafter, we show how the flexibility of ASKCOS helps us handle the cases where automatic planning with typical settings fails, by providing a variety of tree search options, and more importantly, a user-friendly interface to directly modify proposed routes. Additional example pathways proposed by ASKCOS from three re-runs with manual edits where appropriate are shown in Supplementary Figures [S8](#), [S9](#), and [S10](#). Examples of targets for which ASKCOS could not automatically find pathways even after these re-runs are shown in Supplementary Figure [S11](#), which may require interactive planning starting from the targets as described in

Sections [Interactive planning with the template relevance model](#) and [Interactive planning with multiple models, including template-free models](#). Synthesis planning tools can produce a range of suggestions, some of which can be high-risk while others are high-confidence. The ideal balance between creativity and conservatism is a matter of personal judgment, and we find that cross-referencing proposals with literature within or outside of ASKCOS is often fruitful.

The utility of a CASP tool depends on not only the modules and tools it contains, but also how it is deployed and customized within an organization. The open source nature of ASKCOS allows for local deployment and, in particular, deployment behind a firewall when working with proprietary data. Deployment is easily customized so that specific modules can be enabled or disabled as needed to save computational resources. Other possible customization includes retrosynthetic model retraining and integration, as well as replacing the default building block database, which are elaborated in Sections [Model retraining and integration](#), and [User customization](#), respectively. Retraining of translation-based forward predictors with in-house data is similarly possible for the Transformer [18] and Graph2SMILES [17] models, which has been shown to boost prediction performance for company-specific reactions [114].

The development of ASKCOS has been guided by a combination of suggestions from collaborators, colleagues, and the community. With this refreshed open source release, we envision a gradual shift to more community-driven development. We laid the groundwork for easy future extension with a major backend refactor in late-2023, in which a microservice-based architecture was formalized and functionalities were re-modularized, as discussed in Supplementary Section [Technical details of software engineering](#). The refactor has made new model addition straightforward, not only for the ASKCOS team but also for advanced users who would like to replace or extend ASKCOS modules with their own.

We believe that CASP—and computer-assisted chemistry more broadly—is an important part of modern chemistry research that deserves to be precompetitive and accessible to all. At the outset of our work to develop ASKCOS, the landscape of open source solutions was stark; even today, the vast majority of solutions remain commercial even when based on methods described in the open literature. With its ease of use and deployment, customizability, and extensibility, we hope that ASKCOS will find increased adoption and continue to provide a sustainable framework for open source yet production-ready CASP tools.

4 Methods

4.1 Overview

The majority of the usage of ASKCOS from end-users’ perspectives has been covered in Section [Results](#). In the rest of Section [Methods](#), we will elaborate on the details of modeling and computation, which were only briefly discussed in Section [Results](#). Considerations for software development and for the 2023 refactor are elaborated in Supplementary Section [Technical details of software engineering](#). Other non-central but very useful features for more advanced users including model retraining and customization are described in Supplementary Section [Advanced features](#).

The cheminformatics-related computations throughout all ASKCOS modules rely heavily upon RDKit [115]. Machine learning capabilities are implemented with commonly used packages including but not limited to PyTorch, Tensorflow, Numpy, and Pandas. NetworkX is used for modeling most trees and/or graphs. Clustering is done with scikit-learn, hdbscan, or reaction types from a baseline reaction classification model trained on NameRxn labels [100]. Atom mapping is mostly done with RXNMapper [116] or Indigo [117].

4.2 Technical details of the template relevance model and MCTS tree search

Our implementations of the template relevance model and of MCTS tree search deviate from what is described in Segler et al. [25] with several modifications. Specifically, most of the trained template relevance models we provide use simple feedforward neural networks, which we found to have comparable performance to the original but more complex *highway networks* [118] on larger datasets (e.g., with hundreds of thousands of training reactions). We use RDKit for computing Morgan fingerprints of the input targets, and RDChiral [67] for reaction template extraction. The template classification model with feedforward networks is then implemented, trained, and evaluated using PyTorch [119].

Our MCTS tree search differs significantly from Segler et al. [25] and we use a simplified formulation for Upper Confidence bound applied to Trees (UCT) [120]. In particular, we do not use a *rollout* phase. The UCT score for a given reaction node in the search tree is calculated as

$$a_r = Q_r + c * U_r \tag{1}$$

$$= \frac{s_r * v_r}{n_r} + c * \sqrt{\frac{\ln N_r}{n_r}} \tag{2}$$

where the score (a_r) takes into consideration an exploitation term (Q_r) and an exploration term ($c * U_r$) with c being the weight for exploration. Q_r can be interpreted as a heuristic score with s_r being the *reaction score* from model output (e.g., the template probability for template relevance model), v_r being the average buyability score of all children (1.0 for buyables and 0.0 for non-buyables), and n_r having its typical definition of node visit counts. The N_r in the exploration term is the visit counts of the parent chemical node. The tree search operates in a select-expand-update loop starting from the root node. The search network will keep expanding until the termination criteria is reached, for example, if reaching the time limit. A path enumeration phase identifies all synthesis pathways in the search tree. Optionally, the search can be configured to terminate once the first viable pathway is found.

4.3 Technical details of other modules

The implementations of other modules are summarized below.

- **Augmented Transformer [18] for retrosynthesis and forward prediction:** we re-implement using PyTorch, the OpenNMT [121] package, and a regex tokenizer based on previous work by Schwaller [34, 122] to tokenize SMILES strings into input and output tokens.
- **Graph2SMILES [17] for outcome prediction and retrosynthesis:** no deviation from the published version.
- **Retrosim [15] for one-step retrosynthesis:** instead of extracting the templates on-the-fly *after* retrieving similar targets in the original implementation, we pre-extract all templates and store them in the database for later use, which speeds up inference at the expense of storage.
- **Retro* [28] for multi-step search:** while remaining faithful to the original algorithm, the code structure of our implementation is heavily tailored towards that of the MCTS for consistency.
- **WLDN5 [35] for reaction outcome prediction:** no deviation from the published version.
- **The reaction condition recommender [29]:** no deviation from the published version for the V1 model. The V2 model is experimental and undergoing active development with publication underway.
- **The analog counting module from Levin et al. [102]:** no deviation from the published version.
- **The regio-selectivity predictor from Guan et al. [96]:** no deviation from the published version.
- **The site-selectivity predictor from Struble et al. [95]:** no deviation from the published version.
- **The synthesis pathway scorer from Mo et al. [99]:** no deviation from the published version.
- **The SCScore from Coley et al. [61]:** no deviation from the published version.
- **The solubility prediction module from Vermeire et al. [104]:** no deviation from the published version.
- **The QM descriptor predictor from Li et al. [105]:** no deviation from the published version.

Supplementary information. Supplementary information is available for this manuscript, which provides more details on solubility prediction, solvent screening, and QM descriptor prediction, as well as other advanced features including model retraining and customization. It also includes sections to elaborate on software engineering considerations and to describe in detail the case study on FDA-approved small molecule drugs mentioned in Section [Discussion](#).

Acknowledgements. The authors thank the Machine Learning for Pharmaceutical Discovery and Synthesis consortium for numerous discussions over the years. We thank CAS and NextMove Software for providing large-scale reaction data on which various prediction models have been trained. We thank Itlize Global, LLC and TOC Research for providing software development services in the 2023 refactor. We thank all the past and current contributors to ASKCOS who are too numerous to name. The full

contributor list is included in our public instance at <https://askcos.mit.edu> and will be continuously updated.

Declarations

4.4 Funding

This work was supported by the DARPA Make-It program under Contract ARO W911NF-16-2-0023, the Machine Learning for Pharmaceutical Discovery and Synthesis (MLPDS) consortium, and the National Institutes of Health under grant 1U18TR004149. The continued development of ASKCOS is supported by the MLPDS consortium. Z.T. received additional funding from the MolSSI Fellowship Program and the NSERC PGS-D fellowship under Application No. 577866-2023.

4.5 Competing interests

The authors declare no competing interests

4.6 Ethics approval and consent to participate

Not applicable.

4.7 Consent for publication

Not applicable.

4.8 Data availability

All data and trained model weights are shared under MIT licenses, with the only exceptions being models trained on data from the CAS Content Collection [64] which are only accessible to MLPDS members, as well as the template relevance model trained on the Reaxys dataset ca. 2016 and its associated template set, which are released under the non-commercial CC BY-NC 4.0 license.

We refer the reader to the respective manuscripts for various models for the availability of original training data. In particular, data derived from US patents is generally openly available, including but not limited to USPTO_50k and USPTO_full (from the GLN repository [123]), USPTO_480k (from the WLN repository [124]), as well as USPTO_STEREO (from the Molecular Transformer repository [125]). Proprietary data from the CAS Content Collection, Pistachio, or Reaxys are not able to be shared.

4.9 Materials availability

Not applicable.

4.10 Code availability

All of the code associated with ASKCOS is fully open sourced under MIT licenses and is available at https://gitlab.com/mlpds_mit/askcosv2, with the main entry

point being the [askcos2_core](https://doi.org/10.5281/zenodo.13929900) repository. A snapshot of all repositories including all data and model checkpoints at the time of writing has been archived and is available at <https://doi.org/10.5281/zenodo.13929900>. The ASKCOS wiki is accessible at https://gitlab.com/mlpds_mit/askcosv2/askcos-docs/-/wikis/home.

4.11 Author contributions

Z.T., C.W.C., M.L., M.E.F., T.J.S., and M.M. conceived the idea of the 2023 ASKCOS refactor, on top of historical development efforts led by C.W.C., M.E.F., and M.L. Z.T. led the design and implementation of microservice-based ASKCOS. Z.T., S.J.C., M.H.F., and H.L. developed various codes for the refactor. Z.T. and C.W.C. led the manuscript writing. Every author contributed to and approved the manuscript. In particular, J.R. and K.Y. made the figures and polished the introduction. The main contributors for the result sections include J.R. and I.L. (for retrosynthesis and pathway evaluation), X.S. (for reaction condition recommendation), and J.F.J. (for reaction outcome prediction). N.M. and S.-C.L. led the discussion of the solubility modules and of the QM descriptor module, respectively. S.J.C., M.H.F., H.L., and M.M. contributed to various sections for software engineering details in the SI. J.R., Z.T., and J.P.L. conducted the synthesis planning exercise and discussed the results in the SI. W.H.G., K.F.J., and C.W.C. have provided oversight, organization, and fundraising to support the development of ASKCOS over the years.

References

- [1] Corey, E.J., Wipke, W.T.: Computer-Assisted Design of Complex Organic Syntheses. *Science* **166**(3902), 178–192 (1969) <https://doi.org/10.1126/science.166.3902.178> . Publisher: American Association for the Advancement of Science
- [2] Warr, W.A.: A Short Review of Chemical Reaction Database Systems, Computer-Aided Synthesis Design, Reaction Prediction and Synthetic Feasibility. *Molecular Informatics* **33**(6-7), 469–476 (2014) <https://doi.org/10.1002/minf.201400052>
- [3] Corey, E.J.: Robert Robinson Lecture. Retrosynthetic thinking—essentials and examples. *Chemical Society Reviews* **17**(0), 111–133 (1988) <https://doi.org/10.1039/CS9881700111> . Publisher: The Royal Society of Chemistry
- [4] Corey, E.J., Cramer, R.D.I., Howe, W.J.: Computer-assisted synthetic analysis for complex molecules. Methods and procedures for machine generation of synthetic intermediates. *Journal of the American Chemical Society* **94**(2), 440–459 (1972) <https://doi.org/10.1021/ja00757a022> . Publisher: American Chemical Society
- [5] Wipke, W.T., Ouchi, G.I., Krishnan, S.: Simulation and evaluation of chemical synthesis—SECS: An application of artificial intelligence techniques. *Artificial Intelligence* **11**(1), 173–193 (1978) [https://doi.org/10.1016/0004-3702\(78\)90016-4](https://doi.org/10.1016/0004-3702(78)90016-4)

- [6] Funatsu, K., Sasaki, S.-I.: Computer-assisted organic synthesis design and reaction prediction system, "AIPHOS". *Tetrahedron Computer Methodology* **1**(1), 27–37 (1988) [https://doi.org/10.1016/0898-5529\(88\)90006-1](https://doi.org/10.1016/0898-5529(88)90006-1)
- [7] Gasteiger, J., Ihlenfeldt, W.D.: The WODCA System. In: Gasteiger, J. (ed.) *Software Development in Chemistry* 4, pp. 57–65. Springer, Berlin, Heidelberg (1990). https://doi.org/10.1007/978-3-642-75430-2_7
- [8] Grzybowski, B.A., Szymkuć, S., Gajewska, E.P., Molga, K., Dittwald, P., Wołos, A., Klucznik, T.: Chematica: A Story of Computer Code That Started to Think like a Chemist. *Chem* **4**(3), 390–398 (2018) <https://doi.org/10.1016/j.chempr.2018.02.024> . Publisher: Elsevier
- [9] Klucznik, T., Mikulak-Klucznik, B., McCormack, M.P., Lima, H., Szymkuć, S., Bhowmick, M., Molga, K., Zhou, Y., Rickershauser, L., Gajewska, E.P., Toutchkine, A., Dittwald, P., Startek, M.P., Kirkovits, G.J., Roszak, R., Adamski, A., Sieredzińska, B., Mrksich, M., Trice, S.L.J., Grzybowski, B.A.: Efficient Syntheses of Diverse, Medicinally Relevant Targets Planned by Computer and Executed in the Laboratory. *Chem* **4**(3), 522–532 (2018) <https://doi.org/10.1016/j.chempr.2018.02.002> . Publisher: Elsevier
- [10] Mikulak-Klucznik, B., Gołebowska, P., Bayly, A.A., Popik, O., Klucznik, T., Szymkuć, S., Gajewska, E.P., Dittwald, P., Staszewska-Krajewska, O., Beker, W., Badowski, T., Scheidt, K.A., Molga, K., Mlynarski, J., Mrksich, M., Grzybowski, B.A.: Computational planning of the synthesis of complex natural products. *Nature* **588**(7836), 83–88 (2020) <https://doi.org/10.1038/s41586-020-2855-y> . Number: 7836 Publisher: Nature Publishing Group
- [11] Tu, Z., Stuyver, T., Coley, C.W.: Predictive chemistry: machine learning for reaction deployment, reaction development, and reaction discovery. *Chemical Science* **14**(2), 226–244 (2023) <https://doi.org/10.1039/D2SC05089G> . Publisher: The Royal Society of Chemistry
- [12] Schwaller, P., Vaucher, A.C., Laplaza, R., Bunne, C., Krause, A., Corminboeuf, C., Laino, T.: Machine intelligence for chemical reaction space. *WIREs Computational Molecular Science* **12**(5), 1604 (2022) <https://doi.org/10.1002/wcms.1604>
- [13] Dai, H., Li, C., Coley, C., Dai, B., Song, L.: Retrosynthesis Prediction with Conditional Graph Logic Network. In: *Advances in Neural Information Processing Systems*, vol. 32. Curran Associates, Inc., Vancouver, British Columbia, Canada (2019)
- [14] Chen, S., Jung, Y.: Deep Retrosynthetic Reaction Prediction using Local Reactivity and Global Attention. *JACS Au* **1**(10), 1612–1620 (2021) <https://doi.org/10.1021/jacsau.1c00246> . Publisher: American Chemical Society

- [15] Coley, C.W., Rogers, L., Green, W.H., Jensen, K.F.: Computer-Assisted Retrosynthesis Based on Molecular Similarity. *ACS Central Science* **3**(12), 1237–1245 (2017) <https://doi.org/10.1021/acscentsci.7b00355> . Publisher: American Chemical Society
- [16] Xie, S., Yan, R., Guo, J., Xia, Y., Wu, L., Qin, T.: Retrosynthesis prediction with local template retrieval. In: Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence and Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence and Thirteenth Symposium on Educational Advances in Artificial Intelligence. AAAI'23/IAAI'23/EAAI'23, vol. 37, pp. 5330–5338. AAAI Press, Washington DC, US (2023). <https://doi.org/10.1609/aaai.v37i4.25664>
- [17] Tu, Z., Coley, C.W.: Permutation Invariant Graph-to-Sequence Model for Template-Free Retrosynthesis and Reaction Prediction. *Journal of Chemical Information and Modeling* **62**(15), 3503–3513 (2022) <https://doi.org/10.1021/acs.jcim.2c00321> . Publisher: American Chemical Society
- [18] Tetko, I.V., Karpov, P., Van Deursen, R., Godin, G.: State-of-the-art augmented NLP transformer models for direct and single-step retrosynthesis. *Nature Communications* **11**(1), 5575 (2020) <https://doi.org/10.1038/s41467-020-19266-y> . Number: 1 Publisher: Nature Publishing Group
- [19] 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. *Chemical Science* **13**(31), 9023–9034 (2022) <https://doi.org/10.1039/D2SC02763A> . Publisher: Royal Society of Chemistry
- [20] Sacha, M., Błaż, M., Byrski, P., Dabrowski-Tumański, P., Chromiński, M., Loska, R., Włodarczyk-Pruszyński, P., Jastrzebski, S.: Molecule Edit Graph Attention Network: Modeling Chemical Reactions as Sequences of Graph Edits. *Journal of Chemical Information and Modeling* **61**(7), 3273–3284 (2021) <https://doi.org/10.1021/acs.jcim.1c00537> . Publisher: American Chemical Society
- [21] Somnath, V.R., Bunne, C., Coley, C., Krause, A., Barzilay, R.: Learning Graph Models for Retrosynthesis Prediction. In: Advances in Neural Information Processing Systems, vol. 34, pp. 9405–9415. Curran Associates, Inc., virtual (2021)
- [22] Igashov, I., Schneuing, A., Segler, M., Bronstein, M.M., Correia, B.: RetroBridge: Modeling Retrosynthesis with Markov Bridges. In: The Twelfth International Conference on Learning Representations (2023)
- [23] Gaiński, P., Koziarski, M., Maziarz, K., Segler, M., Tabor, J., Śmieja, M.: RetroGFN: Diverse and Feasible Retrosynthesis using GFlowNets. In: ICLR 2024 Workshop on Generative and Experimental Perspectives for Biomolecular Design (2024)

- [24] Heifets, A., Jurisica, I.: Construction of new medicines via game proof search. In: Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence. AAAI'12, pp. 1564–1570. AAAI Press, Toronto, Ontario, Canada (2012)
- [25] Segler, M.H.S., Preuss, M., Waller, M.P.: Planning chemical syntheses with deep neural networks and symbolic AI. *Nature* **555**(7698), 604–610 (2018) <https://doi.org/10.1038/nature25978> . Number: 7698 Publisher: Nature Publishing Group
- [26] Kishimoto, A., Buesser, B., Chen, B., Botea, A.: Depth-First Proof-Number Search with Heuristic Edge Cost and Application to Chemical Synthesis Planning. In: Advances in Neural Information Processing Systems, vol. 32. Curran Associates, Inc., Vancouver, British Columbia, Canada (2019)
- [27] 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. *Chemical Science* **11**(12), 3316–3325 (2020) <https://doi.org/10.1039/C9SC05704H> . Publisher: The Royal Society of Chemistry
- [28] Chen, B., Li, C., Dai, H., Song, L.: Retro*: Learning Retrosynthetic Planning with Neural Guided A* Search. In: Proceedings of the 37th International Conference on Machine Learning, pp. 1608–1616. PMLR, virtual (2020). ISSN: 2640-3498
- [29] Gao, H., Struble, T.J., Coley, C.W., Wang, Y., Green, W.H., Jensen, K.F.: Using Machine Learning To Predict Suitable Conditions for Organic Reactions. *ACS Central Science* **4**(11), 1465–1476 (2018) <https://doi.org/10.1021/acscentsci.8b00357> . Publisher: American Chemical Society
- [30] Maser, M.R., Cui, A.Y., Ryou, S., DeLano, T.J., Yue, Y., Reisman, S.E.: Multilabel Classification Models for the Prediction of Cross-Coupling Reaction Conditions. *Journal of Chemical Information and Modeling* **61**(1), 156–166 (2021) <https://doi.org/10.1021/acs.jcim.0c01234> . Publisher: American Chemical Society
- [31] Kwon, Y., Kim, S., Choi, Y.-S., Kang, S.: Generative Modeling to Predict Multiple Suitable Conditions for Chemical Reactions. *Journal of Chemical Information and Modeling* (2022) <https://doi.org/10.1021/acs.jcim.2c01085> . Publisher: American Chemical Society
- [32] Qian, Y., Li, Z., Tu, Z., Coley, C., Barzilay, R.: Predictive Chemistry Augmented with Text Retrieval. In: Bouamor, H., Pino, J., Bali, K. (eds.) Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pp. 12731–12745. Association for Computational Linguistics, Singapore (2023). <https://doi.org/10.18653/v1/2023.emnlp-main.784>
- [33] Chen, L.-Y., Li, Y.-P.: Enhancing chemical synthesis: a two-stage deep neural

network for predicting feasible reaction conditions. *Journal of Cheminformatics* **16**(1), 11 (2024) <https://doi.org/10.1186/s13321-024-00805-4>

- [34] 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 Central Science* **5**(9), 1572–1583 (2019) <https://doi.org/10.1021/acscentsci.9b00576> . Publisher: American Chemical Society
- [35] Coley, C.W., Jin, W., Rogers, L., Jamison, T.F., Jaakkola, T.S., Green, W.H., Barzilay, R., Jensen, K.F.: A graph-convolutional neural network model for the prediction of chemical reactivity. *Chemical Science* **10**(2), 370–377 (2019) <https://doi.org/10.1039/C8SC04228D> . Publisher: The Royal Society of Chemistry
- [36] Bradshaw, J., Kusner, M.J., Paige, B., Segler, M.H.S., Hernández-Lobato, J.M.: A Generative Model For Electron Paths. In: *The Seventh International Conference on Learning Representations* (2019)
- [37] Bi, H., Wang, H., Shi, C., Coley, C., Tang, J., Guo, H.: Non-Autoregressive Electron Redistribution Modeling for Reaction Prediction. In: *Proceedings of the 38th International Conference on Machine Learning*, pp. 904–913. PMLR, virtual (2021). ISSN: 2640-3498
- [38] Chemical.ai. <https://chemical.ai/> Accessed 2024-07-08
- [39] RXN for Chemistry. <https://rxn.app.accelerate.science/rxn> Accessed 2024-07-08
- [40] Spaya. <https://spaya.ai/> Accessed 2024-07-08
- [41] MANIFOLD. <https://app.postera.ai/> Accessed 2024-08-27
- [42] Molecule.one. <https://www.molecule.one/> Accessed 2024-08-27
- [43] Speed synthetic planning. <https://www.cas.org/solutions/cas-scifinder-discovery-platform/cas-scifinder/synthesis-planning> Accessed 2024-07-08
- [44] Reaxys - Predictive Retrosynthesis. <https://www.elsevier.com/products/reaxys/predictive-retrosynthesis> Accessed 2024-11-18
- [45] Coley, C.W., Thomas, D.A., Lummiss, J.A.M., Jaworski, J.N., Breen, C.P., Schultz, V., Hart, T., Fishman, J.S., Rogers, L., Gao, H., Hicklin, R.W., Plehiers, P.P., Byington, J., Piotti, J.S., Green, W.H., Hart, A.J., Jamison, T.F., Jensen, K.F.: A robotic platform for flow synthesis of organic compounds informed by AI planning. *Science* **365**(6453), 1566 (2019) <https://doi.org/10.1126/science.aax1566> . Publisher: American Association for the Advancement of Science

- [46] Genheden, S., Thakkar, A., Chadimová, V., Reymond, J.-L., Engkvist, O., Bjer-rum, E.: AiZynthFinder: a fast, robust and flexible open-source software for retrosynthetic planning. *Journal of Cheminformatics* **12**(1), 70 (2020) <https://doi.org/10.1186/s13321-020-00472-1>
- [47] Maziarz, K., Tripp, A., Liu, G., Stanley, M., Xie, S., Gainski, P., Seidl, P., Segler, M.: Re-evaluating Retrosynthesis Algorithms with Syntheseus. *Faraday Discussions* (2024) <https://doi.org/10.1039/D4FD00093E> . Publisher: The Royal Society of Chemistry
- [48] MLPDS, Machine Learning for Pharmaceutical Discovery and Synthesis Con-sortium. <https://mlpds.mit.edu> Accessed 2024-07-08
- [49] Struble, T.J., Alvarez, J.C., Brown, S.P., Chytil, M., Cisar, J., DesJarlais, R.L., Engkvist, O., Frank, S.A., Greve, D.R., Griffin, D.J., Hou, X., Johannes, J.W., Kreatsoulas, C., Lahue, B., Mathea, M., Mogk, G., Nicolaou, C.A., Palmer, A.D., Price, D.J., Robinson, R.I., Salentin, S., Xing, L., Jaakkola, T., Green, W.H., Barzilay, R., Coley, C.W., Jensen, K.F.: Current and Future Roles of Artificial Intelligence in Medicinal Chemistry Synthesis. *Journal of Medicinal Chem-istry* **63**(16), 8667–8682 (2020) <https://doi.org/10.1021/acs.jmedchem.9b02120> . Publisher: American Chemical Society
- [50] Levin, I., Liu, M., Voigt, C.A., Coley, C.W.: Merging enzymatic and synthetic chemistry with computational synthesis planning. *Nature Communications* **13**(1), 7747 (2022) <https://doi.org/10.1038/s41467-022-35422-y> . Number: 1 Publisher: Nature Publishing Group
- [51] Soukaina, E., Wissal, L., Yassine, K., Achraf, E.A., Guenoun, F., Bouachrine, M.: Design of new dipeptide inhibitors against SARS-CoV 3CLpro: 3D-QSAR, molecular docking, MD simulation, ADMET studies and retrosynthesis strategy. *Arabian Journal of Chemistry* **17**(2), 105584 (2024) <https://doi.org/10.1016/j.arabjc.2023.105584>
- [52] Avila, C., West, A., C. Vicini, A., Waddington, W., Brearley, C., Clarke, J., M. Derrick, A.: Chemistry in a graph: modern insights into commercial organic synthesis planning. *Digital Discovery* (2024) <https://doi.org/10.1039/D4DD00120F> . Publisher: Royal Society of Chemistry
- [53] Fromer, J.C., Coley, C.W.: An algorithmic framework for synthetic cost-aware decision making in molecular design. *Nature Computational Science* **4**(6), 440–450 (2024) <https://doi.org/10.1038/s43588-024-00639-y> . Publisher: Nature Publishing Group
- [54] Sankaranarayanan, K., Jensen, K.F.: Computer-assisted multistep chemoen-zymatic retrosynthesis using a chemical synthesis planner. *Chemical Science* **14**(23), 6467–6475 (2023) <https://doi.org/10.1039/D3SC01355C> . Publisher: The Royal Society of Chemistry

- [55] Koscher, B.A., Canty, R.B., McDonald, M.A., Greenman, K.P., McGill, C.J., Bilodeau, C.L., Jin, W., Wu, H., Vermeire, F.H., Jin, B., Hart, T., Kulesza, T., Li, S.-C., Jaakkola, T.S., Barzilay, R., Gómez-Bombarelli, R., Green, W.H., Jensen, K.F.: Autonomous, multiproperty-driven molecular discovery: From predictions to measurements and back. *Science* **382**(6677), 1407 (2023) <https://doi.org/10.1126/science.adi1407> . Publisher: American Association for the Advancement of Science
- [56] Nambiar, A.M.K., Breen, C.P., Hart, T., Kulesza, T., Jamison, T.F., Jensen, K.F.: Bayesian Optimization of Computer-Proposed Multistep Synthetic Routes on an Automated Robotic Flow Platform. *ACS Central Science* **8**(6), 825–836 (2022) <https://doi.org/10.1021/acscentsci.2c00207> . Publisher: American Chemical Society
- [57] Mahjour, B., Lu, J., Fromer, J., Casetti, N., Coley, C.: Ideation and Evaluation of Novel Multicomponent Reactions via Mechanistic Network Analysis and Automation (2024). <https://doi.org/10.26434/chemrxiv-2024-qfjh9-v2>
- [58] Seierstad, M., Tichenor, M.S., DesJarlais, R.L., Na, J., Bacani, G.M., Chung, D.M., Mercado-Marin, E.V., Steffens, H.C., Mirzadegan, T.: Novel Reagent Space: Identifying Unorderable but Readily Synthesizable Building Blocks. *ACS Medicinal Chemistry Letters* **12**(11), 1853–1860 (2021) <https://doi.org/10.1021/acsmchemlett.1c00340> . Publisher: American Chemical Society
- [59] Qi, W., Zhai, D., Song, D., Liu, C., Yang, J., Sun, L., Li, Y., Li, X., Deng, W.-Q.: Optimized synthesis of the anti-COVID-19 drugs aided by retrosynthesis software. *RSC Medicinal Chemistry* (2023) <https://doi.org/10.1039/D2MD00444E> . Publisher: RSC
- [60] Pasquini, M., Stenta, M.: Linchemin: Syngraph—a data model and a toolkit to analyze and compare synthetic routes. *Journal of Cheminformatics* **15**(1), 41 (2023)
- [61] Coley, C.W., Rogers, L., Green, W.H., Jensen, K.F.: SCScore: Synthetic Complexity Learned from a Reaction Corpus. *Journal of Chemical Information and Modeling* **58**(2), 252–261 (2018) <https://doi.org/10.1021/acs.jcim.7b00622> . Publisher: American Chemical Society
- [62] Segler, M.H.S., Waller, M.P.: Neural-Symbolic Machine Learning for Retrosynthesis and Reaction Prediction. *Chemistry – A European Journal* **23**(25), 5966–5971 (2017) <https://doi.org/10.1002/chem.201605499>
- [63] The Pistachio dataset. <https://www.nextmovesoftware.com/pistachio.html> Accessed 2024-07-08
- [64] The CAS reactions collection. <https://www.cas.org/cas-data/cas-reactions> Accessed 2024-07-08

- [65] Lowe, D.M.: Extraction of chemical structures and reactions from the literature. Thesis, University of Cambridge (October 2012). <https://doi.org/10.17863/CAM.16293> . Accepted: 2013-07-23T08:23:10Z
- [66] Reaxys. <https://www.elsevier.com/promotions/chemistry-database> Accessed 2024-07-08
- [67] Coley, C.W., Green, W.H., Jensen, K.F.: RDChiral: An RDKit Wrapper for Handling Stereochemistry in Retrosynthetic Template Extraction and Application. *Journal of Chemical Information and Modeling* **59**(6), 2529–2537 (2019) <https://doi.org/10.1021/acs.jcim.9b00286> . Publisher: American Chemical Society
- [68] BKMS-react: an integrated and non-redundant biochemical reaction database. <https://bkms.brenda-enzymes.org/index.php> Accessed 2024-08-27
- [69] Thakkar, A., Selmi, N., Reymond, J.-L., Engkvist, O., Bjerrum, E.J.: “Ring Breaker”: Neural Network Driven Synthesis Prediction of the Ring System Chemical Space. *Journal of Medicinal Chemistry* **63**(16), 8791–8808 (2020) <https://doi.org/10.1021/acs.jmedchem.9b01919> . Publisher: American Chemical Society
- [70] Weininger, D.: SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules. *Journal of Chemical Information and Computer Sciences* **28**(1), 31–36 (1988) <https://doi.org/10.1021/ci00057a005> . Publisher: American Chemical Society
- [71] PUG REST documentation. <https://pubchem.ncbi.nlm.nih.gov/docs/pug-rest#section=Compound-Property-Tables> Accessed 2024-07-09
- [72] Lin, K., Xu, Y., Pei, J., Lai, L.: Automatic retrosynthetic route planning using template-free models. *Chemical Science* **11**(12), 3355–3364 (2020) <https://doi.org/10.1039/C9SC03666K> . Publisher: The Royal Society of Chemistry
- [73] Schneider, N., Stiefl, N., Landrum, G.A.: What’s What: The (Nearly) Definitive Guide to Reaction Role Assignment. *Journal of Chemical Information and Modeling* **56**(12), 2336–2346 (2016) <https://doi.org/10.1021/acs.jcim.6b00564> . Publisher: American Chemical Society
- [74] 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. *Journal of Chemical Information and Modeling* **62**(9), 2111–2120 (2022) <https://doi.org/10.1021/acs.jcim.1c01065> . Publisher: American Chemical Society
- [75] Zhu, J., Xia, Y., Wu, L., Xie, S., Zhou, W., Qin, T., Li, H., Liu, T.-Y.: Dual-view Molecular Pre-training. In: Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining. KDD ’23, pp. 3615–3627. Association

for Computing Machinery, New York, NY, USA (2023). <https://doi.org/10.1145/3580305.3599317>

- [76] Bjerrum, E.J.: SMILES Enumeration as Data Augmentation for Neural Network Modeling of Molecules. arXiv. arXiv:1703.07076 [cs] (2017). <https://doi.org/10.48550/arXiv.1703.07076>
- [77] Silver, D., Huang, A., Maddison, C.J., Guez, A., Sifre, L., Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., Dieleman, S., Grewe, D., Nham, J., Kalchbrenner, N., Sutskever, I., Lillicrap, T., Leach, M., Kavukcuoglu, K., Graepel, T., Hassabis, D.: Mastering the game of Go with deep neural networks and tree search. *Nature* **529**(7587), 484–489 (2016) <https://doi.org/10.1038/nature16961> . Number: 7587 Publisher: Nature Publishing Group
- [78] Liu, G., Xue, D., Xie, S., Xia, Y., Tripp, A., Maziarz, K., Segler, M., Qin, T., Zhang, Z., Liu, T.-Y.: Retrosynthetic Planning with Dual Value Networks. In: Proceedings of the 40th International Conference on Machine Learning, pp. 22266–22276. PMLR, Honolulu, Hawaii, US (2023). ISSN: 2640-3498
- [79] Yu, Y., Wei, Y., Kuang, K., Huang, Z., Yao, H., Wu, F.: GRASP: Navigating Retrosynthetic Planning with Goal-driven Policy. *Advances in Neural Information Processing Systems* **35**, 10257–10268 (2022)
- [80] Schreck, J.S., Coley, C.W., Bishop, K.J.M.: Learning Retrosynthetic Planning through Simulated Experience. *ACS Central Science* **5**(6), 970–981 (2019) <https://doi.org/10.1021/acscentsci.9b00055> . Publisher: American Chemical Society
- [81] Wang, X., Qian, Y., Gao, H., W. Coley, C., Mo, Y., Barzilay, R., F. Jensen, K.: Towards efficient discovery of green synthetic pathways with Monte Carlo tree search and reinforcement learning. *Chemical Science* **11**(40), 10959–10972 (2020) <https://doi.org/10.1039/D0SC04184J> . Publisher: Royal Society of Chemistry
- [82] Xie, S., Yan, R., Han, P., Xia, Y., Wu, L., Guo, C., Yang, B., Qin, T.: RetroGraph: Retrosynthetic Planning with Graph Search. In: Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining. KDD '22, pp. 2120–2129. Association for Computing Machinery, New York, NY, USA (2022). <https://doi.org/10.1145/3534678.3539446>
- [83] Zhao, D., Tu, S., Xu, L.: Efficient retrosynthetic planning with MCTS exploration enhanced A* search. *Communications Chemistry* **7**(1), 1–12 (2024) <https://doi.org/10.1038/s42004-024-01133-2> . Publisher: Nature Publishing Group
- [84] Hong, S., Zhuo, H.H., Jin, K., Shao, G., Zhou, Z.: Retrosynthetic planning with experience-guided Monte Carlo tree search. *Communications Chemistry* **6**(1), 1–14 (2023) <https://doi.org/10.1038/s42004-023-00911-8> . Publisher: Nature Publishing Group

- [85] Kreutter, D., Reymond, J.-L.: Multistep retrosynthesis combining a disconnection aware triple transformer loop with a route penalty score guided tree search. *Chemical Science* **14**(36), 9959–9969 (2023) <https://doi.org/10.1039/D3SC01604H> . Publisher: The Royal Society of Chemistry
- [86] Liu, S., Tu, Z., Xu, M., Zhang, Z., Lin, L., Ying, R., Tang, J., Zhao, P., Wu, D.: FusionRetro: Molecule Representation Fusion via In-Context Learning for Retrosynthetic Planning. In: *Proceedings of the 40th International Conference on Machine Learning*, pp. 22028–22041. PMLR, Honolulu, Hawaii, US (2023). ISSN: 2640-3498
- [87] Kim, J., Ahn, S., Lee, H., Shin, J.: Self-Improved Retrosynthetic Planning. In: *Proceedings of the 38th International Conference on Machine Learning*, pp. 5486–5495. PMLR, virtual (2021). ISSN: 2640-3498
- [88] Shee, Y., Li, H., Morgunov, A., Batista, V.: DirectMultiStep: Direct Route Generation for Multi-Step Retrosynthesis. *arXiv* (2024). <https://doi.org/10.48550/arXiv.2405.13983>
- [89] Yu, K., Roh, J., Li, Z., Gao, W., Wang, R., Coley, C.W.: Double-Ended Synthesis Planning with Goal-Constrained Bidirectional Search. In: *Advances in Neural Information Processing Systems*. Curran Associates, Inc., Vancouver, British Columbia, Canada (2024)
- [90] Marcou, G., Sousa, J., Latino, D.A.R.S., Luca, A., Horvath, D., Rietsch, V., Varnek, A.: Expert System for Predicting Reaction Conditions: The Michael Reaction Case. *Journal of Chemical Information and Modeling* **55**(2), 239–250 (2015) <https://doi.org/10.1021/ci500698a> . Publisher: American Chemical Society
- [91] Li, J., Eastgate, M.D.: Making better decisions during synthetic route design: leveraging prediction to achieve greenness-by-design. *Reaction Chemistry & Engineering* **4**(9), 1595–1607 (2019) <https://doi.org/10.1039/C9RE00019D> . Publisher: The Royal Society of Chemistry
- [92] Wang, X., Hsieh, C.-Y., Yin, X., Wang, J., Li, Y., Deng, Y., Jiang, D., Wu, Z., Du, H., Chen, H., Li, Y., Liu, H., Wang, Y., Luo, P., Hou, T., Yao, X.: Generic Interpretable Reaction Condition Predictions with Open Reaction Condition Datasets and Unsupervised Learning of Reaction Center. *Research* **0**(ja) (2023) <https://doi.org/10.34133/research.0231> . Publisher: American Association for the Advancement of Science
- [93] Coley, C.W., Barzilay, R., Jaakkola, T.S., Green, W.H., Jensen, K.F.: Prediction of Organic Reaction Outcomes Using Machine Learning. *ACS Central Science* **3**(5), 434–443 (2017) <https://doi.org/10.1021/acscentsci.7b00064> . Publisher: American Chemical Society

- [94] Chen, S., Jung, Y.: A generalized-template-based graph neural network for accurate organic reactivity prediction. *Nature Machine Intelligence*, 1–9 (2022) <https://doi.org/10.1038/s42256-022-00526-z> . Publisher: Nature Publishing Group
- [95] Struble, T.J., Coley, C.W., Jensen, K.F.: Multitask prediction of site selectivity in aromatic C–H functionalization reactions. *Reaction Chemistry & Engineering* **5**(5), 896–902 (2020) <https://doi.org/10.1039/D0RE00071J> . Publisher: The Royal Society of Chemistry
- [96] Guan, Y., Coley, C.W., Wu, H., Ranasinghe, D., Heid, E., Struble, T.J., Pattanaik, L., Green, W.H., Jensen, K.F.: Regio-selectivity prediction with a machine-learned reaction representation and on-the-fly quantum mechanical descriptors. *Chemical Science* **12**(6), 2198–2208 (2021) <https://doi.org/10.1039/D0SC04823B> . Publisher: The Royal Society of Chemistry
- [97] Hiramama, M., Uei, M.: Carbamate mediated 1,3-asymmetric induction. A stereoselective synthesis of acyclic 1,3-diol systems. *Tetrahedron Letters* **23**(50), 5307–5310 (1982) [https://doi.org/10.1016/S0040-4039\(00\)85825-6](https://doi.org/10.1016/S0040-4039(00)85825-6)
- [98] Neves, P., McClure, K., Verhoeven, J., Dyubankova, N., Nugmanov, R., Gedich, A., Menon, S., Shi, Z., Wegner, J.K.: Global reactivity models are impactful in industrial synthesis applications. *Journal of Cheminformatics* **15**(1), 20 (2023) <https://doi.org/10.1186/s13321-023-00685-0>
- [99] Mo, Y., Guan, Y., Verma, P., Guo, J., Fortunato, M.E., Lu, Z., Coley, C.W., Jensen, K.F.: Evaluating and clustering retrosynthesis pathways with learned strategy. *Chemical Science* **12**(4), 1469–1478 (2021) <https://doi.org/10.1039/D0SC05078D> . Publisher: The Royal Society of Chemistry
- [100] NameRxn: Expert System for Named Reaction Identification and Classification. <https://www.nextmovesoftware.com/namerxn.html> Accessed 2024-07-10
- [101] Hoffmann, R.W.: Ranking of Synthesis Plans. In: Hoffmann, R.W. (ed.) *Elements of Synthesis Planning*, pp. 133–144. Springer, Berlin, Heidelberg (2009). https://doi.org/10.1007/978-3-540-79220-8_8
- [102] Levin, I., Fortunato, M.E., Tan, K.L., Coley, C.W.: Computer-aided evaluation and exploration of chemical spaces constrained by reaction pathways. *AIChE Journal* **69**(12), 18234 (2023) <https://doi.org/10.1002/aic.18234>
- [103] Kuznetsov, A., Sahinidis, N.V.: ExtractionScore: A Quantitative Framework for Evaluating Synthetic Routes on Predicted Liquid–Liquid Extraction Performance. *Journal of Chemical Information and Modeling* **61**(5), 2274–2282 (2021) <https://doi.org/10.1021/acs.jcim.0c01426> . Publisher: American Chemical Society
- [104] Vermeire, F.H., Chung, Y., Green, W.H.: Predicting Solubility Limits of Organic

- Solutes for a Wide Range of Solvents and Temperatures. *Journal of the American Chemical Society* **144**(24), 10785–10797 (2022) <https://doi.org/10.1021/jacs.2c01768> . Publisher: American Chemical Society
- [105] Li, S.-C., Wu, H., Menon, A., Spiekermann, K.A., Li, Y.-P., Green, W.H.: When Do Quantum Mechanical Descriptors Help Graph Neural Networks to Predict Chemical Properties? *Journal of the American Chemical Society* **146**(33), 23103–23120 (2024) <https://doi.org/10.1021/jacs.4c04670> . Publisher: American Chemical Society
- [106] Griffin, D.J., Coley, C.W., Frank, S.A., Hawkins, J.M., Jensen, K.F.: Opportunities for Machine Learning and Artificial Intelligence to Advance Synthetic Drug Substance Process Development. *Organic Process Research & Development* (2023) <https://doi.org/10.1021/acs.oprd.3c00229> . Publisher: American Chemical Society
- [107] Stuyver, T., Coley, C.W.: Quantum chemistry-augmented neural networks for reactivity prediction: Performance, generalizability, and explainability. *The Journal of Chemical Physics* **156**(8) (2022)
- [108] 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 Medicinal Chemistry* **15**(4), 1085–1095 (2024) <https://doi.org/10.1039/D3MD00651D> . Publisher: RSC
- [109] eMolecules. <https://www.emolecules.com/> Accessed 2024-11-05
- [110] MilliporeSigma. <https://www.sigmaaldrich.com/> Accessed 2024-11-05
- [111] WuXi LabNetwork. <https://www.labnetwork.com/> Accessed 2024-11-05
- [112] Mcule: Advanced Tools to Find and Order Molecules Online. <https://mcule.com/> Accessed 2024-11-05
- [113] ChemBridge: The Gold Standard in Small Molecule Libraries and Building Blocks. <https://chembridge.com/> Accessed 2024-11-05
- [114] Lee, A.A., Yang, Q., Sresht, V., Bolgar, P., Hou, X., Klug-McLeod, J.L., Butler, C.R.: Molecular Transformer unifies reaction prediction and retrosynthesis across pharma chemical space. *Chemical Communications* **55**(81), 12152–12155 (2019) <https://doi.org/10.1039/C9CC05122H> . Publisher: The Royal Society of Chemistry
- [115] RDKit: Open-Source Cheminformatics Software. <https://www.rdkit.org/> Accessed 2024-07-25

- [116] Schwaller, P., Hoover, B., Reymond, J.-L., Strobelt, H., Laino, T.: Extraction of organic chemistry grammar from unsupervised learning of chemical reactions. *Science Advances* **7**(15), 4166 (2021) <https://doi.org/10.1126/sciadv.abe4166> . Publisher: American Association for the Advancement of Science
- [117] Indigo Toolkit. <https://lifescience.opensource.epam.com/indigo/api/index.html> Accessed 2024-07-25
- [118] Srivastava, R.K., Greff, K., Schmidhuber, J.: Training Very Deep Networks. In: *Advances in Neural Information Processing Systems*, vol. 28. Curran Associates, Inc., Montréal, Québec, Canada (2015)
- [119] PyTorch. <https://pytorch.org/> Accessed 2024-07-25
- [120] Kocsis, L., Szepesvári, C.: Bandit based monte-carlo planning. In: *Proceedings of the 17th European Conference on Machine Learning. ECML'06*, pp. 282–293. Springer, Berlin, Heidelberg (2006). https://doi.org/10.1007/11871842_29
- [121] Klein, G., Kim, Y., Deng, Y., Senellart, J., Rush, A.: OpenNMT: Open-Source Toolkit for Neural Machine Translation. In: Bansal, M., Ji, H. (eds.) *Proceedings of ACL 2017, System Demonstrations*, pp. 67–72. Association for Computational Linguistics, Vancouver, Canada (2017)
- [122] 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. *Chemical Science* **9**(28), 6091–6098 (2018) <https://doi.org/10.1039/C8SC02339E> . Publisher: The Royal Society of Chemistry
- [123] The GLN repository. <https://github.com/Hanjun-Dai/GLN/> Accessed 2024-10-14
- [124] The WLN repository. <https://github.com/wengong-jin/nips17-rexgen/> Accessed 2024-10-14
- [125] The Molecular Transformer repository. <https://github.com/pschwlr/MolecularTransformer/> Accessed 2024-10-14

Supplementary information for ASKCOS: an open source software suite for synthesis planning

Zhengkai Tu¹, Sourabh J. Choure², Mun Hong Fong²,
Jihye Roh², Itai Levin³, Kevin Yu⁴, Joonyoung F. Joung²,
Nathan Morgan², Shih-Cheng Li², Xiaoqi Sun², Huiqian Lin²,
Mark Murnin², Jordan P. Liles², Thomas J. Struble⁵,
Michael E. Fortunato⁶, Mengjie Liu^{2,†}, William H. Green²,
Klavs F. Jensen², Connor W. Coley^{1,2*}

¹Department of Electrical Engineering and Computer Science,
Massachusetts Institute of Technology, 77 Massachusetts Ave,
Cambridge, MA, 02139, USA.

²Department of Chemical Engineering, Massachusetts Institute of
Technology, 77 Massachusetts Ave, Cambridge, MA, 02139, USA.

³Department of Biological Engineering, Massachusetts Institute of
Technology, 77 Massachusetts Ave, Cambridge, MA, 02139, USA.

⁴Center for Computational Science and Engineering, Massachusetts
Institute of Technology, 77 Massachusetts Ave, Cambridge, MA, 02139,
USA.

⁵Bristol Myers Squibb, 250 Water Street, Cambridge, MA, 02141, USA.

⁶Novartis Institutes for BioMedical Research, Inc., 250 Massachusetts
Avenue, Cambridge, MA, 02139, USA.

*Corresponding author(s). E-mail(s): ccoley@mit.edu;
Contributing authors: ztu@mit.edu; sjchoure@mit.edu;
fong410@mit.edu; jroh99@mit.edu; itail@mit.edu; kyu3@mit.edu;
jjoung@mit.edu; knathan@mit.edu; scli@mit.edu; xiaoqis@mit.edu;
linhq@mit.edu; murninm@mit.edu; jliles24@mit.edu;
Thomas.Struble@bms.com; mike.fortunato@novartis.com;
mjliu@mit.edu; whgreen@mit.edu; kfjensen@mit.edu;

[†]Current affiliation: AstraZeneca. Work done while at MIT.

S1 Details of solubility prediction and solvent screening

The screenshot shows the ASKCOS web interface for solubility prediction. The top section contains input fields for Solute (SMILES: NC(=O)OC(Cn1ncn1)c1cccc1Cl), Solvent (SMILES: C1CCOC1), and Temperature (323 K). Below these are chemical structure drawings for the solute and solvent. A row of buttons includes SUBMIT, RUN BATCH, MORE OPTIONS (highlighted with a dashed blue box and an annotation 'Click to add solubility reference data or solute thermodynamic data'), CLEAR RESULTS, and MODEL I/O DETAILS (with an annotation 'Click to select columns to display'). A 'DOWNLOAD' button is also present. A 'Select Columns' dropdown menu is shown with options: Solubility(T) [mg/mL], Solubility(298) [mg/mL], and Pred. Solute Data. Below the menu is a table of results with columns: Solute, Solvent, Temperature, Solubility (method1) [mg/mL], Solubility (method2) [mg/mL], Solubility(298) [mg/mL], Pred. Hsub298 [kcal/mol], and Pred. Cpg298 [cal/K/m]. The table contains two rows of data for the solute in THF at 298 K and 323 K. At the bottom, there is a pagination control showing 'Items per page: 10' and '1-2 of 2'.

Solute	Solvent	Temperature	Solubility (method1) [mg/mL]	Solubility (method2) [mg/mL]	Solubility(298) [mg/mL]	Pred. Hsub298 [kcal/mol]	Pred. Cpg298 [cal/K/m]
NC(=O)OC(Cn1ncn1)c1cccc1Cl	C1CCOC1	298	1.03e+2	1.03e+2	1.03e+2	37.	63.1
NC(=O)OC(Cn1ncn1)c1cccc1Cl	C1CCOC1	323	3.20e+2	3.05e+2	1.03e+2	37.	63.1

Fig. S1: Annotated screenshot of solubility prediction results in ASKCOS. The SMILES strings of the solute and the solvent are input in the **Solute** and **Solvent** panels. The desired temperature (in K) for the solubility prediction is input in the **Temperature** panel. The predicted solubility and solvation properties of solute cenobamate, defined by the SMILES string NC(=O)OC(Cn1ncn1)c1cccc1Cl, in solvent THF, defined by the SMILES string C1CCOC1, at 298 K and 323 K are displayed. The property values to display (e.g., $\log S$, Abraham parameters, etc.) can be selected from the **Select Columns** panel.

The required input for the solubility prediction model is a solute and a solvent molecule, which are specified using a SMILES string or by using the drawing functionality (Figure S1). After clicking the **SUBMIT** button, the predicted solubility in milligrams per milliliter is shown in a table at the bottom of the screen. The default prediction temperature is 298 K, but this can be adjusted by using the temperature field on the main screen. Other optional input fields are found by clicking **MORE**

OPTIONS. This opens a window where users can supply a few thermodynamic parameters to improve the model’s predictions. The model makes solubility predictions by first predicting the aqueous solubility at 298 K ($\log(S_{aq,298\text{K}})$) and then by using other thermochemical quantities to correct that solubility to the desired solvent and temperature. The known solubility of the solute in another solvent (a reference solubility) can be given to the model to use instead of the aqueous solubility prediction, which is especially useful for solutes that have low aqueous solubility. The solubility is corrected from the reference solvent to the target solvent using ML predicted solvation free energies ($\Delta G_{solv,298\text{K}}$) in both solvents. $\Delta G_{solv,298\text{K}}$ in the target solvent can also be viewed in the output table by selecting its column in the drop down menu, which is found by clicking the arrow in the **Select Columns** field.

The solubility in the target solvent at 298 K, $\log(S_{298\text{K}})$, is corrected to other temperatures using the solute’s dissolution enthalpy in the target solvent, $\Delta H_{diss,T}$, which in turn is estimated using ML predictions of more common thermochemical quantities in a thermodynamic cycle. These quantities include the solute’s sublimation enthalpy at 298 K ($\Delta H_{sub,298\text{K}}$), solid phase heat capacity at 298 K ($C_{p,s}$), gas phase heat capacity at 298 K ($C_{p,g}$), and solvation enthalpy either at 298 K ($\Delta H_{solv,298\text{K}}$) or the target temperature ($\Delta H_{solv,T}$). All of these can be viewed in the output table by selecting them in the **Select Columns** drop down menu, similar to $\Delta G_{solv,298\text{K}}$. The first three quantities are predicted via correlations that use ML predicted Abraham solute parameters [1]. These parameters can also be viewed in the output table. Known values for the first three quantities can be given to the model to use instead of the correlations from the solute parameters. This is done using the same window as the reference solubility.

The ΔH_{solv} used in the thermodynamic cycle can be predicted in two ways, which gives two temperature dependent solubility estimates. These are labeled “method1” and “method2” in the output table. The first method neglects the temperature dependence of ΔH_{solv} and uses an ML model to predict $\Delta H_{solv,298\text{K}}$, which is used as $\Delta H_{solv,T}$. This works best for temperatures below 350 K. The second method uses an ML model to predict $\Delta G_{solv,T}$, from which $\Delta H_{solv,T}$ and the temperature dependent entropy ($\Delta S_{solv,T}$) can be calculated. This model requires the critical temperature and density of the solvent, which limits this method’s use to around 100 solvents. These temperature dependent quantities can be viewed in the output table for supported solvents.

A few other optional columns can be viewed in the output table, including the solubility in moles per liter, prediction uncertainty estimates, and any error or warning messages. The uncertainty estimates are the variance of predictions made by an ensemble of models. Each model in the ensemble starts with a different random initialization and is trained on the same data. An ensemble of 10, 12, and 30 models is used for $\Delta G_{solv,298\text{K}}$, $\Delta H_{solv,298\text{K}}$, and $\log(S_{aq,298\text{K}})$ respectively. The prediction values in the output table for these quantities are the mean of the ensemble’s predictions.

It is also possible to use the solubility utility to make bulk predictions. Clicking the **RUN BATCH** button opens a window where a csv or json file containing input data can be uploaded. Clicking **MODEL I/O DETAILS** shows the format of this file. In the

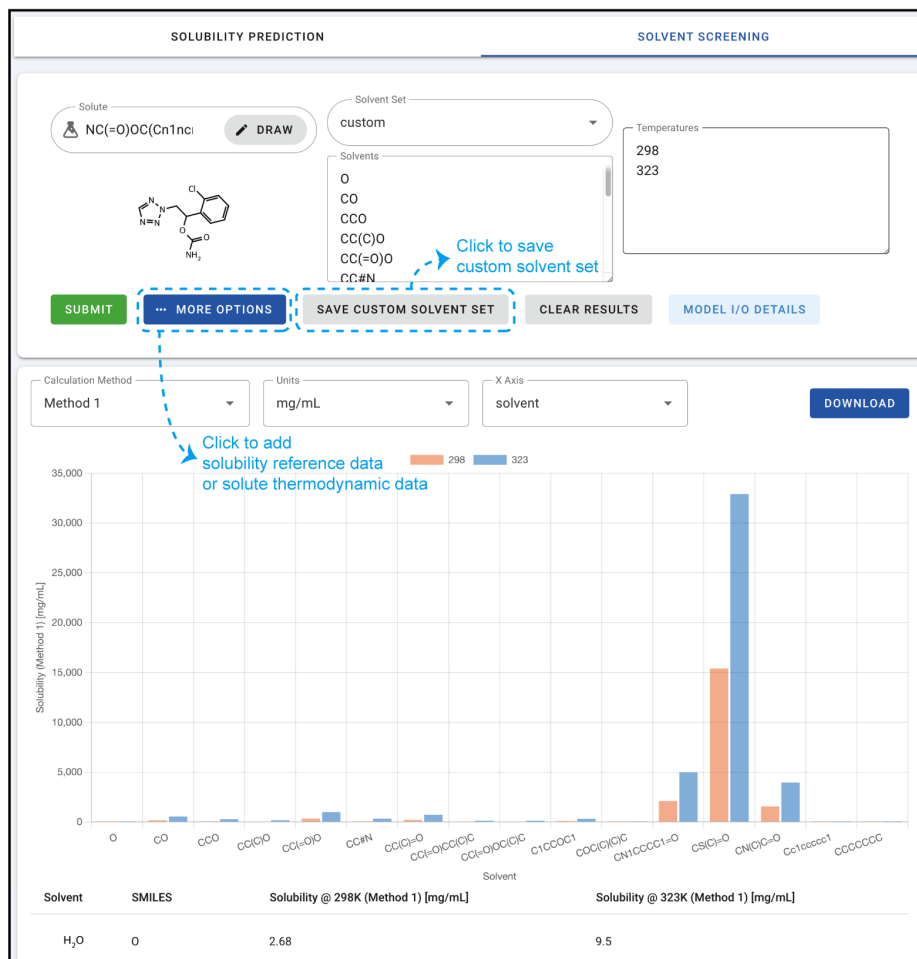


Fig. S2: Annotated screenshot of solvent screening results in ASKCOS. The SMILES string of the solute is input in the **Solute** panel. Users can select the solvents from the two predefined sets or create custom solvent sets in the **Solvent Set** panel. The desired temperatures (in K) at which to compute the solubility and solvation properties are input in the **Temperatures** panel. The predicted solubility for solute cenobamate defined by the SMILES string NC(=O)OC(Cn1ncnn1)c1ccccc1Cl in a custom set of solvents at 298 K and 323 K is displayed.

output section, there is also a **DOWNLOAD** button that allows for downloading all the results as a csv or json file.

The solvent screening utility, as shown in Figure S2, also allows for bulk solubility predictions of a single solute in a set of solvents and temperatures. Two sets of a variety of solvents are predefined, and users can also make a custom set. The output

predictions are shown in both table and chart forms. A bar chart is used if solvent is plotted on the x-axis, whereas a line chart is used if temperature is plotted on the x-axis.

S2 Details of QM descriptor prediction

ML has been widely used for chemical property predictions, but its performance and generalizability often depends on the availability of large datasets. However, acquiring large-scale datasets can be expensive and time-consuming in scenarios such as *de novo* drug and material design. Driven by the belief that QM descriptors can provide deeper physical insights, recent studies [2–4] have tried augmenting ML models with QM descriptors. It has been shown that this approach can improve the accuracy and generalizability of the model, especially when trained on smaller datasets. Thus, this method could potentially accelerate the exploration of complex reaction pathways in retrosynthesis and the prioritization of candidates in high-throughput screening.

Since calculating QM descriptors for all the molecules of interest can be very expensive, we include trained models from Li et al.’s work [3] in ASKCOS to predict 37 general QM descriptors on the fly. These models are directed message passing neural network (D-MPNN) models trained on QM descriptors computed using ω B97XD/def2-SVP//GFN2-xTB. The 37 descriptors consist of 13 atom descriptors, 4 bond descriptors, and 20 molecular descriptors. The atom descriptors include NPA charges, Parr functions, NMR shielding constants, and valence orbital occupancies. The bond descriptors consist of bond order, bond length, bonding electrons, and bond natural ionicity. The molecular descriptors encompass energy gaps, ionization potential (IP), electron affinity (EA), and dipole and quadrupole moments. For partial charge prediction, the summation is constrained by the molecule’s net charge, nucleophilic/electrophilic Parr functions sum to 1, and no constraints apply to other atom or bond descriptors.

The required input for QM descriptor prediction is either a SMILES string or the use of the drawing functionality. After clicking the **SUBMIT** button, the predicted QM descriptors will be displayed in a table at the bottom of the screen. By default, only the NPA charges and Parr functions are shown. Other predicted descriptors can be displayed in the output table by selecting them from the **Select Columns** dropdown menu. Each atom and bond property of a molecule is saved in a list, with the order defined by the RDKit [5] molecule object. To better visualize these descriptors, users can click the 3D visualization button in the table to view the predicted values individually (Figure S3) for the selected molecule. In the output section, there is also a **DOWNLOAD** button that allows you to download all the results as a CSV or JSON file.

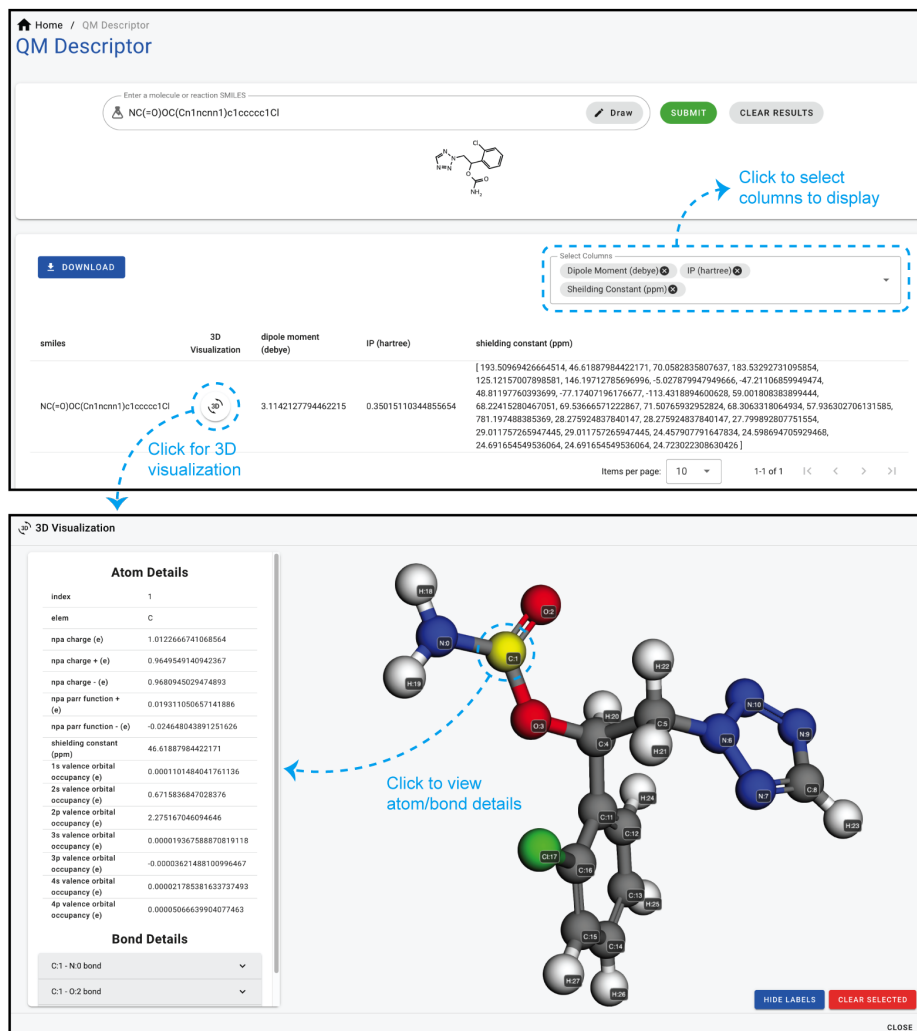


Fig. S3: Annotated screenshot of QM descriptor prediction results in ASKCOS. The predicted QM descriptors of cenobamate, defined by the SMILES string NC(=O)OC(Cn1ncnn1)c1ccccc1Cl, are displayed (top). The atom and bond descriptors can also be viewed individually with the 3D visualization page (bottom). The selected atom is highlighted in yellow, and the predicted atom properties for the selected atom is displayed on the left panel. This panel also provides predicted values for each bond connected to the selected atom.

S3 Technical details of software engineering

S3.1 Refactor into a microservice-based architecture

ASKCOS was originally developed as a pure Python 2 package beginning in 2016 to plan synthetic routes for a robotic flow chemistry platform as part of the DARPA Make-It program, with a full name of “Automated System for Knowledge-based Continuous Organic Synthesis”. Various restful endpoints for modules were defined using the Django package, and the user interface was designed mostly using Javascript and basic Jinja templates. Since then, the number of different predictive models incorporated in ASKCOS has continued to grow. Given its long development history, the monolithic nature of ASKCOS pre-2023 presented a significant challenge in terms of maintainability and extensibility. Managing ever-changing dependencies and their deprecation (e.g., Keras, Tensorflow 1) and synchronizing environments across disparate models developed years-apart by different graduate students and postdocs became unsustainable.

The main deciding factor for a major refactor of ASKCOS was the need to resolve package dependency conflicts. As an extreme example, there is no straightforward way to run a Python program in which part of the code depends on Python 2, and another part depends on Python 3. This type of dependency conflicts became inevitable as we tried to integrate newer and more powerful modules from different contributors into ASKCOS. The de facto and somewhat obvious solution was to turn prediction modules into containerized microservices, which was exactly what has been done in the 2023 refactor, with careful decoupling and re-modularization to retain existing functionalities.

S3.2 Containerized microservices for backend modules

Prediction modules as microservices have been the cornerstone of ASKCOS, but it was not until the 2023 refactor that this philosophy was formalized. We enforce modularization and consistency of these microservices by, without loss of generality, having one containerized service per prediction module. No python function call between the modules is possible and any dependency call has to be made via http requests. For example, when the Tree Builder calls any one-step expansion model, it sends an API requests (via the API gateway) to that one-step prediction service. This preserves modularization for maintainability and extensibility in the long run, at little expense as we did not observe noticeable overhead from the conversion of function calls to API calls, especially when the services are hosted on the same machine.

The main challenges when transitioning into a microservice-based architecture was decoupling and modularization, as with code refactor in general. Once the prediction modules have been fully modularized, it is easy to replace function calls with http calls, which are now made by specialized API classes in our implementation. These modules can then be wrapped as services (i.e., as REST endpoints) using packages like FastAPI [6] and TorchServe [7], and subsequently containerized with Docker. We refer the reader to the ASKCOS wiki [8] which has fully documented

the process of converting python code into containerized services with guiding examples, under the Sections `Development-Packaging python codes as services` and `Development-Containerizing backend services`.

S3.3 Centralized API gateway

Calls to backend services are routed via the *API gateway*, which is the central service in ASKCOS for API management. The API gateway is implemented with FastAPI, and consists mainly of *wrappers* for routing API calls to backend services, as well as *utils* which define endpoints for lightweight services such as drawing and authentication within the gateway itself. In order to take full advantage of FastAPI's capability to automatically generate API documentation based on endpoint definitions, these definitions have all been type hinted. In particular, each wrapper defines the expected input schema such as the name of the model used and other hyperparameters in a `dataclass`, which becomes visible from the auto-generated API documentation as will be illustrated in Section [API usage](#). The schema for the responses returned from the API gateway is defined in similar `dataclasses`. In this way, we maintain a standard pattern for all API endpoints in the gateway and make it easy to add wrappers for new prediction services: simply copy and paste the definition from an existing wrapper, and then specify the input and response schema.

We have also reworked the asynchronous workflows which are necessary for API calls from the frontend and for long-running tasks. The logic for asynchronous task management has been abstracted away and standardized as an additional endpoint, namely, `/call-async` within each wrapper. These asynchronous endpoints merely sub-call the corresponding synchronous endpoints in the same wrapper, generally with no other complication. This design keeps the effort for maintenance and extension to the minimum, while still allowing for customization if needed. The async workflows are implemented with the Celery package with RabbitMQ as the broker and Redis for storing results.

S3.4 Frontend

The ASKCOS frontend is implemented in Vue 3 [9] which offers a more flexible, maintainable, and performant architecture. Vue 3, with its composition API and enhanced reactivity system, allows developers to create scalable components and modularize code more effectively. In the context of ASKCOS, Vue 3 serves as the foundation for a dynamic user interface, where complex chemical retrosynthesis data can be seamlessly displayed and interacted with. The integration of Vuetify [10], a robust Material Design component library, ensures a visually cohesive and responsive user interface. Vuetify components are highly customizable, enabling developers to build feature-rich and intuitive user interfaces while adhering to modern web standards. This synergy of Vue 3 and Vuetify helps create an engaging user experience that can handle complex data visualization (in particular, on the IPP canvas) while maintaining accessibility across various devices.

A few other tools are used for testing, application state management and routing. To ensure the reliability and robustness of this complex system, Cypress [11] is

employed for end-to-end testing. Cypress allows for the testing of every aspect of the application from user interaction to API calls, ensuring that ASKCOS functions as intended in real-world scenarios. Pinia [12] serves as the state management tool to provide a more lightweight, modular, and type-safe approach to managing application state. This helps maintain data consistency across different components, especially in handling chemical datasets and user interactions. Furthermore, Vue Router [13] integrates seamlessly with the Vue 3 framework, managing complex navigational patterns, enabling smooth transitions between different sections of the application without reloading the page. These technologies together form the backbone of the ASKCOS web platform, ensuring high performance, reliability, and ease of use for chemists and researchers accessing the platform.

S3.5 Application monitoring and logging

Several utilities have been designed to help monitor application status while ASKCOS is running. From the **Server Status** page of the ASKCOS user interface, users can see a status summary of the services for the celery workers, the database, and prediction modules. For the celery services, the status summary shows the numbers of their pending tasks as well as of available and busy workers. For the data collections (e.g., reactions, templates, and buyable building blocks) which are stored in the Mongo database, a description of each collection and the number of documents imported during database seeding are displayed. For backend services, metadata including model names and descriptions are populated from the config file used to deploy and start ASKCOS. The service status (online vs. offline) is checked with API calls and updated upon refresh.

The logs of the API calls made from the frontend are accessible from the **Logs** page, which can be particularly helpful for debugging. The storage of such logs is made possible by Pinia which allows a state to be shared across components and pages. When a user performs an action on any ASKCOS page, if requests to the backend APIs are made, Pinia can relay information about these requests to the **Logs** page. The logs can be cleared by refreshing any page in ASKCOS.

For superusers, another useful API endpoint is `/api-logging/get`, which counts the number of API calls to each endpoint and aggregates by date. It provides succinct statistics of the most used features in ASKCOS to help superusers optimize the deployment, e.g., by increasing the number of celery workers for more frequently used modules.

S4 Advanced features

S4.1 User-friendly and customizable deployment

ASKCOS can be deployed from scratch with the following five commands:

```
1 $ mkdir ASKCOSv2
2 $ cd ASKCOSv2
3 $ git clone git@gitlab.com:mlpds_mit/askcosv2/askcos2_core.git
4 $ cd askcos2_core
5 $ make deploy
```


We refer the reader to the ASKCOS wiki [8] for the full instructions, including hardware and software requirements. The last command, `make deploy`, is the main deployment command which does the following in sequence:

1. cloning all other repositories under `ASKCOSv2/` based on the central config file in `askcos2_core`, while downloading data and model checkpoints if needed;
2. generating deployment scripts for building Docker images, starting services and stopping services based on the central config file;
3. building Docker images for all services using the generated script;
4. downloading database data and seeding into the Mongo database;
5. starting all services using using the generated script.

A single centralized config file specifies all the required configurations. It is easy to *turn off* unneeded modules for users who are more resource-constrained or have interest in only a subset of modules. Partial deployment has been documented in the ASKCOS wiki, and we include several sample config files for typical use cases (e.g., retro-only and/or backend-only).

S4.2 Model retraining and integration

A few models for one-step retrosynthesis and reaction outcome prediction have been designed for easy retraining and integration with new datasets. In particular, the retraining pipeline of the template relevance model has been streamlined. Only a list of reaction IDs and reaction SMILES are needed, and the user has the choice of either providing own train/validation/test splits, or letting the retraining engine handle the splitting. Thereafter, automated retraining and testing can be performed with the provided script (e.g., `benchmark*.sh`) after specifying the paths to the reaction files.

Integration of newly trained models into ASKCOS requires *model archiving*, followed by updating the existing deployment config, and possibly seeding additional data into the database. Archiving models into a servable `.mar` file is done using `torch-model-archiver`, while the rest of the integration pipeline is mostly bookkeeping to place model files into the correct location and make the API gateway aware of the existence of the new model. The full process of retraining and integration has been documented in the ASKCOS wiki under the Section `Deployment-Model retraining and integration`.

S4.3 User customization

ASKCOS provides options to customize many aspects of the application. The design of the deployment pipeline discussed in Section [User-friendly and customizable deployment](#) allows users to easily customize their deployment. The complete application including the frontend and the backend is deployed by default, but if ASKCOS is being integrated with other systems where only API calls are being made and the user interface is not needed, backend-only deployment may be more suitable. Similarly, while some datasets are copied over and loaded into ASKCOS at first deployment, additional datasets can be added, or others removed, at a later date by simply invoking subroutines in the `deploy.sh` script in the `askcos2_core` directory to modify the

database. Additionally, as eluded to previously in Section [User-friendly and customizable deployment](#), superusers can chose to disable certain modules depending on usage and computational resources available.

Other aspects of the environment can be customized too, and most customizable environment (env) variables are found in the `.env.example` file. The user interface can be configured to use a different logo, for example, a company logo, and have a different welcome message which may include the company name. This is very important as ASKCOS is publicly accessible, and as such, there are no protections regarding intellectual property (IP). If users want to use ASKCOS on their proprietary, sensitive data, ASKCOS should be deployed behind their company firewalls. Customizing the user interface quickly informs users they are on their internal deployment and can safely and confidently enter potentially sensitive compound data. The superuser can also change the support email addresses (by modifying the `VITE_CONTACT_EMAIL` and `VITE_SUPPORT_EMAILS` env variables) from the default MIT support groups to internal support addresses to ensure sensitive data is not shared outside of the company.

Many companies have policies that block applications requesting data outside of their firewalls to safeguard intellectual property. ASKCOS has only one external lookup, calling on the NIH name resolving service [14] to convert a compound name to its corresponding SMILES string. This service can be quickly disabled on each page by clicking the icon beside the target input area. It can permanently be disabled on deployment by setting the env variable `VITE_ENABLE_SMILES_RESOLVER` to `False`.

ASKCOS allows users to ban chemicals and reactions. These lists are specific to that user account and not global. This allows users to ban compounds or reactions that are patented, that may be toxic or harmful, or that are otherwise considered undesirable. Chemicals or reactions can be added to the ban list directly from the IPP. Banned chemicals or reactions will not appear in the predicted pathways afterwards. Users can view, add, or delete items in their ban lists on My [Banlist](#) page accessible from the left hand sidebar.

The [Buyables](#) page and backend data structures have also been enhanced to allow superusers to add important metadata such as lead time and availability to their own building blocks. These custom metadata are easily viewable during interactive planning (via the `SEARCH BUYABLES` button in the chemical node detail panel) to assist chemists' decision making. Links to the procurement website can also be added to enable easier and quicker ordering. Superusers can add buyables individually or submit a list of buyables in a file for bulk upload.

S4.4 API usage

As the frontend and the backend are fully separated in ASKCOS, it is possible to send http requests directly to the API gateway, e.g., from the command line or using some client library from any programming language. Using the APIs this way may be particularly useful for integrating some of ASKCOS' predictive modules into other workflows and for batch queries, for example, by checking whether routes with buyable building blocks can be found for a given list of (possibly many) molecules. The auto-generated API documentation mentioned in Section [Centralized API gateway](#), which is guaranteed to be up-to-date by design, can be accessed from the browser as an

interactive web page for any locally deployed ASKCOS instance. The endpoints have been organized by groups on the documentation page, and many of these endpoints have been pre-filled with working sample queries which can be easily executed with a few clicks. The interactive nature of the documentation makes it easy for experimenting with various APIs, as well as understanding their schema and performance.

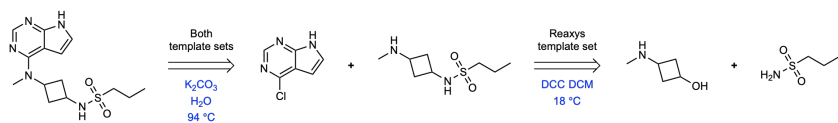
We refer the readers to the ASKCOS wiki [8] for detailed and more visual illustration of the API usage, under the Section **Advanced Usage-Using the APIs directly** where python script examples for API queries have also been provided.

S5 Illustration of applying ASKCOS to FDA-approved small molecule drugs from 2019 to 2023

We perform a case study on FDA-approved New Chemical Entities from 2019 to 2023 (New Drug Applications only) to demonstrate how a chemist may use ASKCOS to plan routes for new targets. The compilation of raw data is retrieved from the FDA website [15]. The raw chemical names of drug components are converted to SMILES using the NIH name resolving API [14]. Unresolvable names and components found in the buyable database of ASKCOS are dropped, and the remaining list of targets are deduplicated based on the canonical SMILES. The final list contains 75 unique targets in total. The raw and processed data files, as well as the scripts for data processing and for sending tree building queries are provided with instructions under https://gitlab.com/mlpds_mit/askcosv2/askcos2.core/-/tree/main/examples.

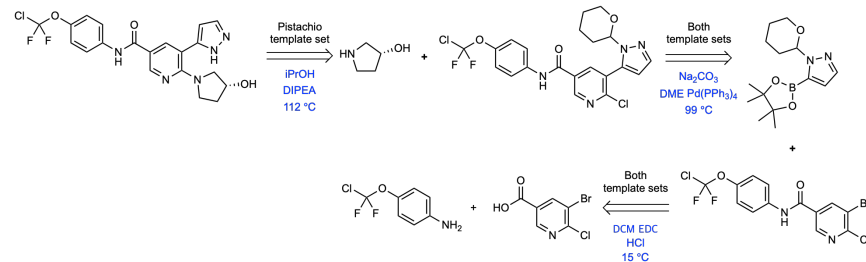
S5.1 Initial automated tree building with typical settings

As a baseline, we first query the MCTS endpoint with typical search settings. Specifically, Pistachio-trained and Reaxys-trained template relevance models are used simultaneously for one-step retrosynthetic expansion, both with a maximum of 1,000 templates and a maximum cumulative template probability of 0.999 per expansion step. The minimum threshold for the plausibility from the binary fast filter is set to 0.001. For the tree search, a maximum branching factor of 25 and a maximum search depth of 6 are used. The number of chemical nodes to be explored is capped at 5,000 and the maximum price for buyable building blocks is set at \$100/g, with no limit on expansion time. This baseline run takes about an hour to finish for all 75 targets on a typical desktop (with Intel i7-12700 CPU, 32 GB of RAM, and no GPU), corresponding to roughly one minute of wall time per molecule. After hypothetical retrosynthetic routes terminating in buyable building blocks are found for a target, routes are algorithmically sorted by length, by average plausibility based on the binary filter, and then by average template score. We include the rank-1 (by this definition) route for each of 10 sample targets in Figures S4, S5, S6, and S7. Note that these pathways are *exactly* as returned by the Tree Builder and have not been postprocessed (i.e., filtered or rescored) by any other modules in ASKCOS. We show the top-1 proposed conditions for each step using the V1 condition recommendation model. Users can change the settings to view lower-ranking predictions by the V1 model and quantitative predictions by the V2 model.



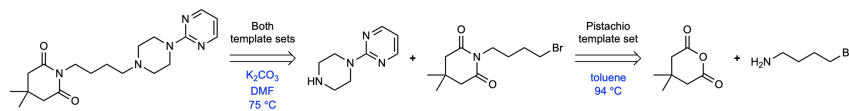
Abrocitinib:

CCCS(=O)(=O)NC1CC(N(C)c2ncnc3[nH]ccc23)C1



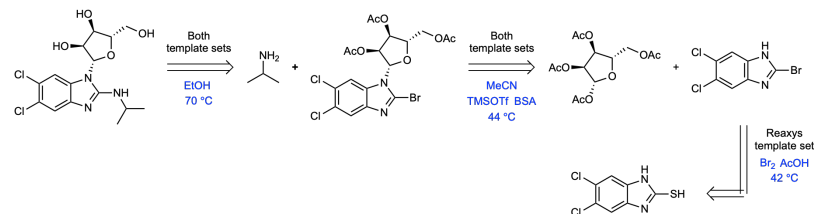
Asciminib:

O=C(Nc1ccc(OC(F)(F)Cl)cc1)c1cnc(N2CC[C@@H](O)C2)c(-c2ccn[nH]2)c1



Gepirone:

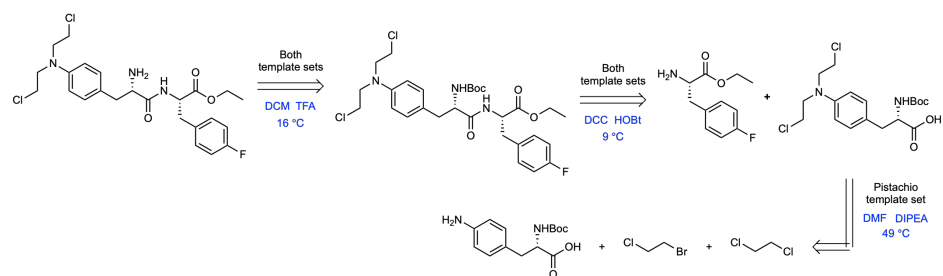
CC1(C)CC(=O)N(CCCCN2CCN(c3ncccn3)CC2)C(=O)C1



Maribavir:

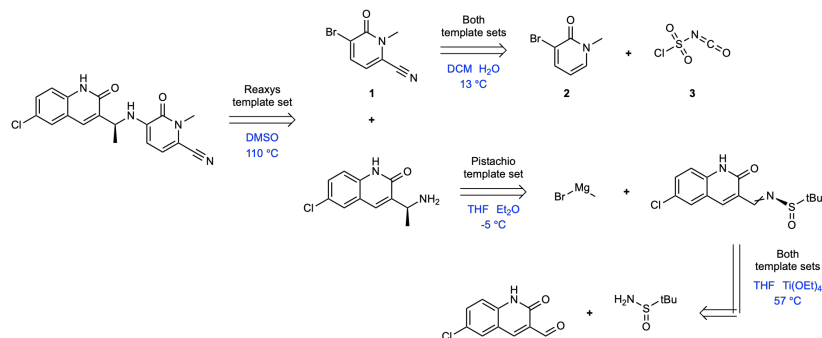
CC(C)Nc1nc2cc(Cl)c(Cl)cc2n1[C@H]1O[C@@H](CO)[C@H](O)[C@@H]1O

Fig. S4: Shortest retrosynthetic routes suggested by ASKCOS for FDA-approved small molecule drug components, part I. Top-1 recommendations by the condition recommender (V1) are shown in blue below each arrow. Ions are written in salt form. Abbreviations are defined in Table S1.



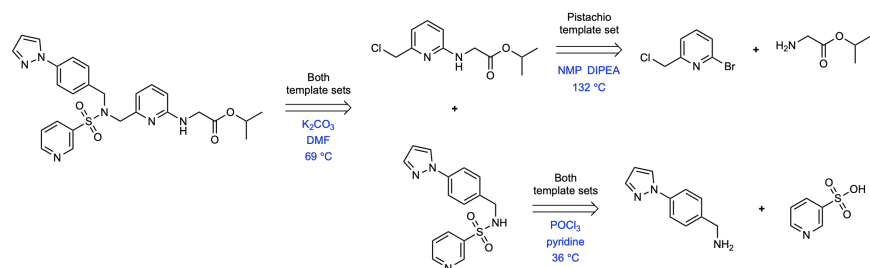
Melphalan flufenamide:

CCOC(=O)[C@H](Cc1ccc(F)cc1)NC(=O)[C@@H](N)Cc1ccc(N(CCCl)CCCl)cc1



Olutasidenib:

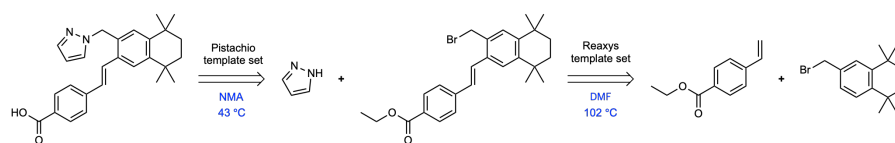
C[C@H](Nc1ccc(C#N)n(C)c1=O)c1cc2cc(Cl)ccc2[nH]c1=O



Omidenepag isopropyl:

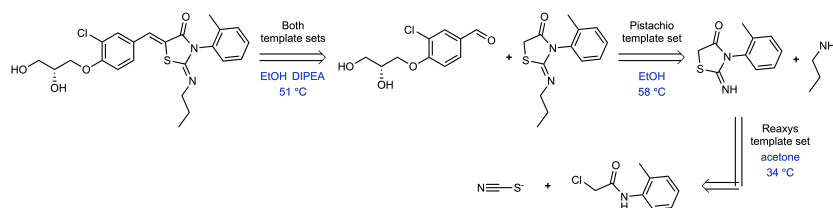
CC(C)OC(=O)CNc1cccc(CN(Cc2ccc(-n3ccc3)cc2)S(=O)(=O)c2cccnc2)n1

Fig. S5: Shortest retrosynthetic routes suggested by ASKCOS for FDA-approved small molecule drug components, part II. Top-1 recommendations by the condition recommender (V1) are shown in blue below each arrow. Ions are written in salt form. Abbreviations are defined in Table S1.



Palovarotene:

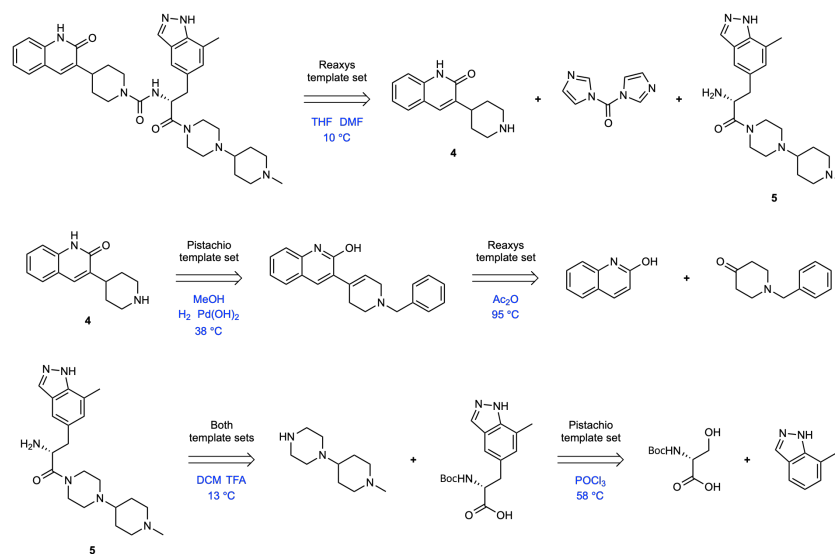
CC1(C)CCC(C)(C)c2cc(Cn3cccn3)c(/C=C/c3ccc(C(=O)O)cc3)cc21



Ponesimod:

CCCN=C1S/C(=C\c2ccc(OC[C@H](O)CO)c(Cl)c2)C(=O)N1c1cccc1C

Fig. S6: Shortest retrosynthetic routes suggested by ASKCOS for FDA-approved small molecule drug components, part III. Top-1 recommendations by the condition recommender (V1) are shown in blue below each arrow. Abbreviations are defined in Table S1.



Zavegepant:

Cc1cc(C[C@@H](NC(=O)N2CCC(c3cc4ccccc4[nH]c3=O)CC2)C(=O)N2CCN(C3CCN(C)CC3)CC2)cc2cn[nH]c12

Fig. S7: Shortest retrosynthetic routes suggested by ASKCOS for FDA-approved small molecule drug components, part IV. Top-1 recommendations by the condition recommender (V1) are shown in blue below each arrow. Ions are written in salt form. Abbreviations are defined in Table S1.

Table S1: List of abbreviations used in Figures S4, S5, S6, and S7

Abbreviation	Full name
BSA	benzenesulfonamide
DCC	1,3-dicyclohexylcarbodiimide
DCE	1,2-dichloroethane
DCM	dichloromethane
DEAD	diethylazodicarboxylate
DIPEA	<i>N,N</i> -diisopropylethylamine
DMAP	<i>N,N</i> -dimethyl-4-aminopyridine
DME	ethylene glycol dimethyl ether
DMF	<i>N,N</i> -dimethylformamide
DMSO	dimethyl sulfoxide
DPP	diphenyl phosphate
EDC	1-ethyl-3-(3-dimethylaminopropyl)carbodiimide
HATU	hexafluorophosphate azabenzotriazole tetramethyl uronium
HMPA	hexamethylphosphoric triamide
HOBt	hydroxybenzotriazole
LDA	lithium diisopropylamine
NMA	<i>N</i> -methylacetamide
NMP	<i>N</i> -methyl-2-pyrrolidone
TFA	trifluoroacetic acid
THF	tetrahydrofuran
TMSOTf	trimethylsilyl trifluoromethanesulfonate

S5.2 Step-wise verification and further analysis *within* ASKCOS

Successfully proposing routes for a target means that the recursive retrosynthetic tree search was able to identify hypothetical pathways that terminate in buyable starting materials, but does not necessarily mean that those hypothetical pathways are chemically plausible. Here, we discuss how ASKCOS facilitates cross-referencing with literature when needed, and how we can analyze suggestions more thoroughly.

At least a few steps among the routes in Figures S4, S5, S6, and S7 are counter-intuitive or seemingly implausible. For example, in the route for olutasidenib (Figure S5), cyanopyridone **1** is prepared from the corresponding less substituted pyridone **2** in one step using chlorosulfonyl isocyanate (CSI) **3**. Since any step proposed by a template relevance model can be traced back to the associated template(s) and literature precedents, we can easily cross-reference the origin of the suggestion in the graphical interface by clicking the reaction node, clicking the template, and clicking the reference link sequentially as described in the Results Section. In this case, the precedent substrates for this cyanation [16–18] are all pyrroles. Visualization of the template shows that it captures only the requirement of an aromatic N-heterocycle without accounting for the ring size or other substituent effects. While this is typical for algorithmically-extracted templates, these types of errors can be identified and manually corrected. We exported the route into the Interactive Path Planner (IPP), removed this problematic step, and re-expanded manually with the Reaxys template set to find an alternative using trimethylsilyl cyanide. Checking the new template

and some associated references [19–21] provides stronger evidence that cyanation with trimethylsilyl cyanide might afford **1** from **2**. It is worth noting that we can always manually replace this final step after exporting the route into the IPP. In a more general setting, if we are unsatisfied with some particular intermediate step(s), we can either re-expand the whole sub-tree, or pick another route with the unpromising step appearing more upstream in the synthesis so that we can “fail early”.

Similar to the reaction steps, the validity and/or optimality of the top proposed reaction condition is not guaranteed. In particular, for the last step of the route for abrocitinib (Figure S4), DCC is suggested as part of the top conditions. As a reagent generally for amide coupling, it does not typically directly activate alcohols (into better leaving groups for nucleophilic attack). The V2 condition recommendation model instead proposes Mitsunobu conditions for this step (Table S2) with sufficient literature precedent to suggest plausibility [22]. For the last step of the route for asciminib (Figure S4), a mild, base-free set of conditions is initially proposed for the proposed amide coupling with EDC. These conditions may be improved by addition of a base, and as in the rank 2 and 3 predictions by the V1 condition recommender, DMAP and HOBt (Table S3), which are commonly used in combination with EDC. For the second step of the route for maribavir (Figure S4), bis(trimethylsilyl)acetamide in the rank 3 condition by V1 may be a better base to use than benzenesulfonamide (BSA) in the rank 1 & 2 recommended conditions by V1 (Table S4). The final retrosynthetic step for zavegepant (Figure S7) would also warrant further investigation of potential protecting group chemistry and/or specialized conditions to achieve the desired site selectivity. Modifications of reactions conditions based on user expertise are common; in practice, we interpret recommended conditions as starting points for empirical screening and further optimization, rather than a claim that the reaction will/must proceed exactly as proposed.

Table S2: Recommended conditions for the last step in the proposed retrosynthesis of abrocitinib. The top 3 recommendations of the V2 condition recommender using fingerprint (fp) and graph representations of molecules are shown.

V1	Rank 1	DCM, DCC at 18°C	
	Rank 2	DCM, DCC, DMAP at 17°C	
	Rank 3	THF, DCC at 24°C	
V2 (fp)	Rank 1	THF (110.71 equiv.), PPh ₃ (0.57 equiv.), DEAD (1.52 equiv.) at 22°C with reactants 2A (1 equiv.), 2B (1.06 equiv.)	
	Rank 2	H ₂ O (38.37), THF (70.34 equiv.), PPh ₃ (0.71 equiv.), DEAD (1.39 equiv.) at 24°C with reactants 2A (1 equiv.), 2B (1.11 equiv.)	
	Rank 3	THF (68.31 equiv.) at 36°C with reactants 2A (1 equiv.), 2B (1.34 equiv.)	
V2 (graph)	Rank 1	pyridine (8.61 equiv.) at 8°C with reactants 2A (1 equiv.), 2B (1.15 equiv.)	
	Rank 2	MeCN (75.16 equiv.) at 19°C with reactants 2A (1 equiv.), 2B (1.22 equiv.)	
	Rank 3	DCM (12.96 equiv.) at 6°C with reactants 2A (1 equiv.), 2B (1.22 equiv.)	

Table S3: Recommended conditions for the last step in the proposed retrosynthesis of asciminib. The top 3 recommendations of the V2 condition recommender using fingerprint (fp) and graph representations of molecules are shown.

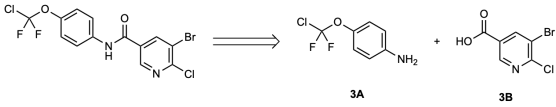
			
V1	Rank 1	DCM, EDC, HCl at 15°C	
	Rank 2	DCM, DMAP, EDC, HCl at 15°C	
	Rank 3	DCM, HOBT, EDC, HCl at 15°C	
V2 (fp)	Rank 1	DMF (29.75 equiv.), H ₂ O (2034.89 equiv.), HATU (1.64 equiv.), DIPEA (3.30 equiv.) at 37°C with reactants 3A (1.06 equiv.), 3B (1 equiv.)	
	Rank 2	DMF (21.93 equiv.), HATU (1.68 equiv.), DIPEA (3.10 equiv.) at 40°C with reactants 3A (1.06 equiv.), 3B (1 equiv.)	
	Rank 3	DMF (18.34 equiv.), AcOEt (174.46 equiv.), HATU (1.69 equiv.), DIPEA (2.91 equiv.) at 39°C with reactants 3A (1.06 equiv.), 3B (1 equiv.)	
V2 (graph)	Rank 1	H ₂ O (162.66 equiv.), MeCN (111.64 equiv.), Na ₂ CO ₃ (2.65 equiv.) at 44°C with reactants 3A (1.16 equiv.), 3B (1 equiv.)	
	Rank 2	MeCN (76.23 equiv.) at 43°C with reactants 3A (1.38 equiv.), 3B (1 equiv.)	
	Rank 3	MeCN (108.14 equiv.), Et ₃ N (4.39 equiv.) at 44°C with reactants 3A (1.16 equiv.), 3B (1 equiv.)	

Table S4: Recommended conditions for the second step in the proposed retrosynthesis of maribavir. The top 3 recommendations of the V2 condition recommender using fingerprint (fp) and graph representations of molecules are shown.

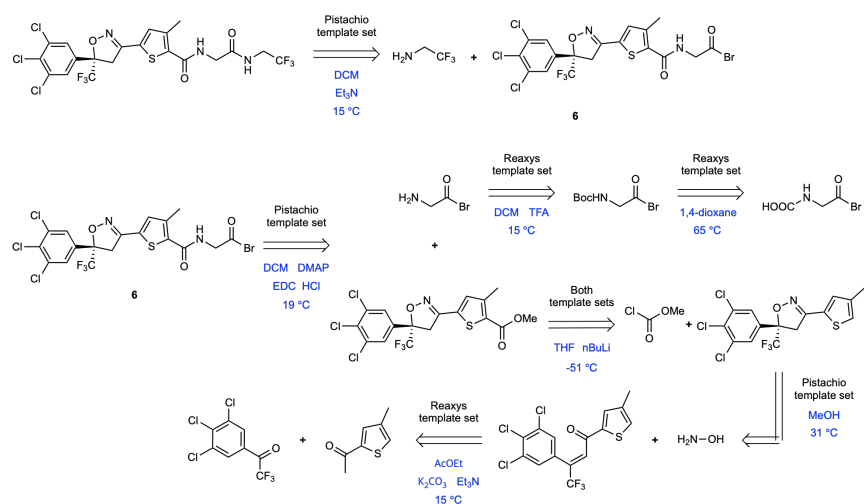
V1	Rank 1	MeCN, TMSOTf, BSA	at 44°C
	Rank 2	TMSOTf, BSA	at 63°C
	Rank 3	MeCN, TMSOTf, bis(trimethylsilyl)acetamide	at 45°C
V2 (fp)	Rank 1	MeCN (166.38 equiv.),	at 48°C with reactants 4A (1 equiv.), 4B (1 equiv.)
	Rank 2	AcOEt (176.53 equiv.), MeCN (159.74 equiv.)	at 51°C with reactants 4A (1 equiv.), 4B (1 equiv.)
	Rank 3	H ₂ O (471.41 equiv.), AcOEt (148.80 equiv.), MeCN (143.76 equiv.)	at 53°C with reactants 4A (1 equiv.), 4B (1 equiv.)
V2 (graph)	Rank 1	THF (50.86 equiv.), <i>n</i> -pentane (13.81 equiv.), <i>t</i> -BuLi (1.43 equiv.)	at -59°C with reactants 4A (1 equiv.), 4B (1 equiv.)
	Rank 2	O ₂ (3.56 equiv.)	at 62°C with reactants 4A (1 equiv.), 4B (1 equiv.)
	Rank 3	DCM (51.48 equiv.), dimethyl sulfide (3.71 equiv.)	at -6°C with reactants 4A (1 equiv.), 4B (1 equiv.)

S5.3 Re-running the automated Tree Builder with different search settings

There are a number of known failure modes that may explain why routes are not found for all 75 targets under default search settings. For example, the tree search process may get stuck in some local optimum where it only explore paths following a particular disconnection it incorrectly thought to be highly promising; the template sets used (Pistachio and Reaxys) may not have templates corresponding to rarer (and arguably more interesting) transformations. One potential solution is straightforward: simply re-run the Tree Builder jobs with different settings, which only takes machine time. Like in the initial run, we queue up a large number of Tree Builder jobs via python scripts (available in the example folder mentioned at the start of this Section) and allow them to run to completion at the background. Results are automatically stored in the database for later inspection through the graphical interface.

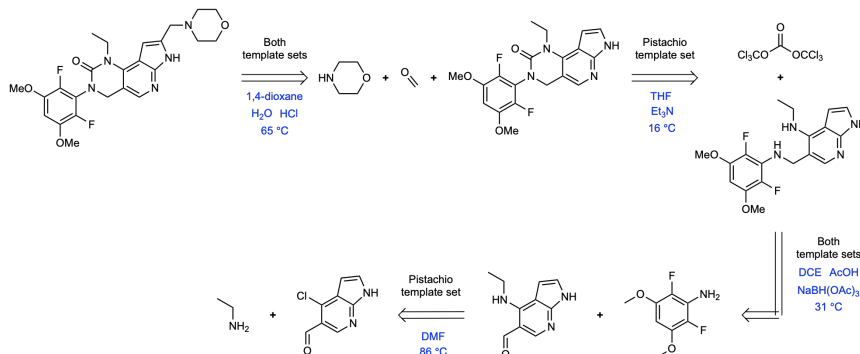
We experimented three re-runs for demonstration purposes, where each run took about an hour to finish for all 75 targets, similar to the baseline. First, we change the maximum number of templates per expansion step to encourage the exploration of more diverse transformations. In an unconstrained setting (e.g., no limit on the number of reactions or chemicals explored) or during interactive planning, the maximum number of templates should be *increased* so that each expansion may individually explore more transformations. In a constrained setting like the initial run where we cap the maximum number of chemicals explored to 5,000, however, it may be desirable to *decrease* the maximum number of templates per expansion so that each may be explored more thoroughly. As an extreme numerical example, if the maximum number of templates per expansion is set to 5,000, the tree search may immediately reach the limit of 5,000 chemicals after a single expansion of the target and will terminate before it is able to consider any pathways of depth 2 or greater. In contrast, if the maximum number of templates per expansion is set to 10, the limit of 5,000 chemicals may be reached much more slowly, allowing for more thorough exploration of all of the 10 templates proposed for the target.

We re-ran with the maximum number of templates reduced to 100 while keeping the other settings exactly the same as in the baseline. Hypothetical routes could then be found for more targets. The shortest routes for lotilaner and pemigatinib are depicted in Figure S8. In both cases, the pathways returned by ASKCOS involve heterocycle forming steps. For lotilaner, the isoxazoline ring is prepared from an α,β -unsaturated ketone with hydroxylamine; for pemigatinib, the cyclic urea in the fused ring system is prepared by an intramolecular cyclization with triphosgene. To re-iterate the typical disclaimer for template-based retrosynthetic proposals, these steps have some supporting precedents but may or may not be achievable experimentally (e.g., due to the presumed stereoselectivity of the former). We can perform similar verification and further analysis as detailed in Section [Step-wise verification and further analysis within ASKCOS](#), and it is up to the user’s discretion on whether to accept or reject proposed routes or steps.



Lotilaner:

Cc1cc(C2=NO[C@@](c3cc(Cl)c(Cl)c(Cl)c3)(C(F)(F)F)C2)sc1C(=O)NCC(=O)NCC(F)(F)F

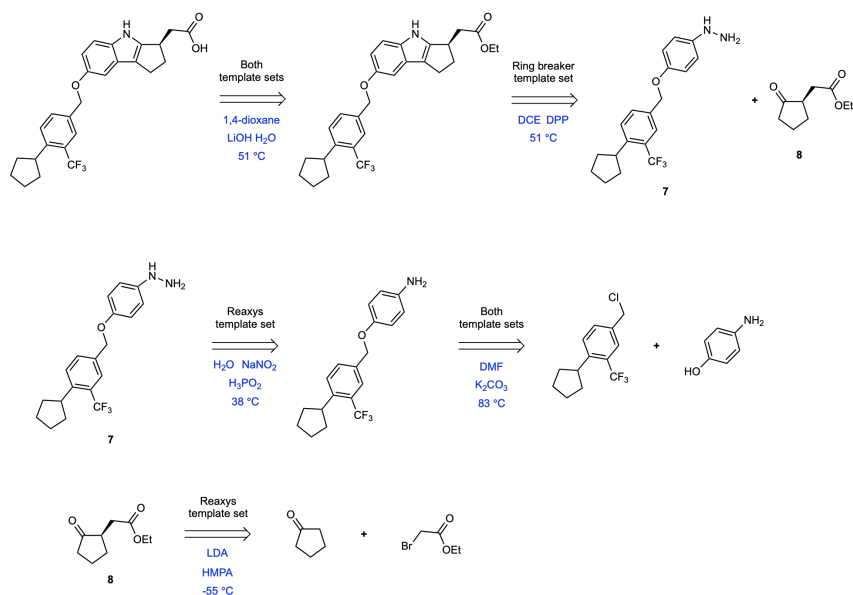


Pemigatinib:

CCN1C(=O)N(c2c(F)c(OC)cc(OC)c2F)C2cnc3[nH]c(CN4CCOCC4)cc3c21

Fig. S8: Shortest retrosynthetic routes suggested by ASKCOS for lotilaner and pemigatinib when re-running with smaller numbers of templates per expansion step. Top-1 recommendations by the condition recommender (V1) are shown in blue below each arrow. Ions are written in salt form. Abbreviations are defined in Table S1.

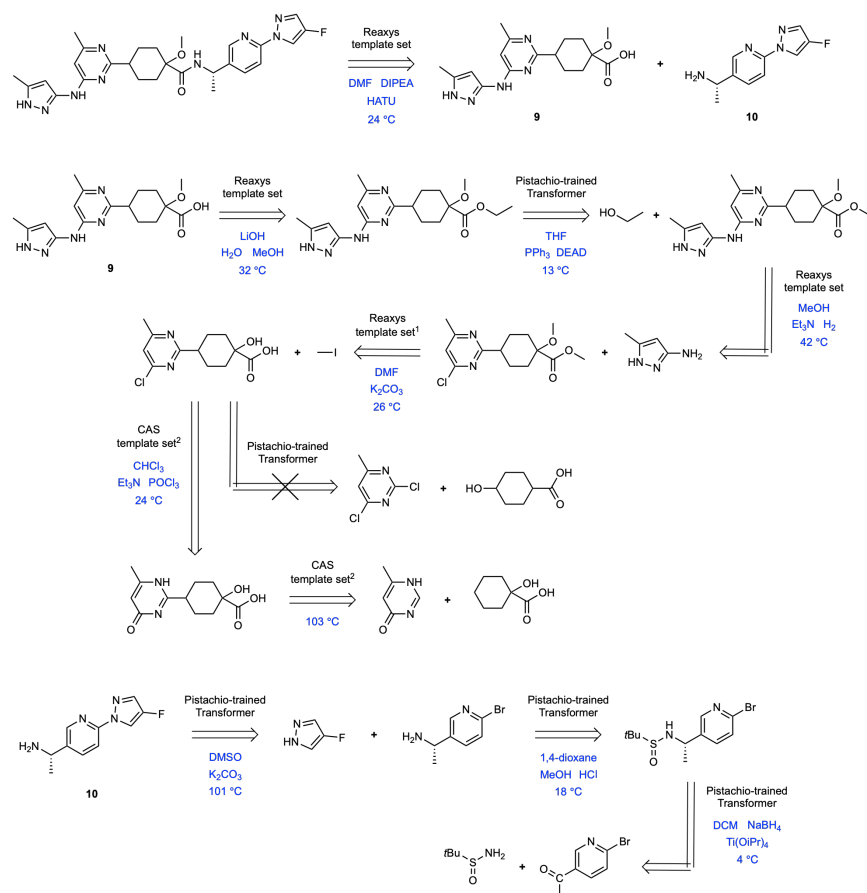
For the second re-run, we included the ring breaker [23] template set in addition to Pistachio and Reaxys, while keeping all the other settings exactly the same as the baseline. While the ring breaker template set is derived from the Pistachio dataset [24] and overlaps with the Pistachio template set, it prioritizes (retrosynthetically) ring breaking transformation. As shown in Figure S9, in the shortest proposed route for etrasimod which was newly solved from this re-run, the 6-5-5 fused ring system is prepared via a Fischer indole synthesis between the phenylhydrazine and a chiral cyclic ketone. Templates corresponding to Fischer indole synthesis exist in the default template sets but may be better represented and prioritized by the ring breaker template relevance model.



Etrasimod:

O=C(O)C[C@H]1CCc2c1[nH]c1ccc(OCc3ccc(C4CCCC4)c(C(F)(F)F)c3)cc21

Fig. S9: Shortest retrosynthetic route suggested by ASKCOS for etrasimod when re-running with the addition of the ring breaker template set. Top-1 recommendations by the condition recommender (V1) are shown in blue below each arrow. Ions are written in salt form. Abbreviations are defined in Table S1.



Pralsetinib:

COC1(C(=O)N[C@@H](C)c2ccc(-n3cc(F)cn3)nc2)CCC(c2nc(C)cc(Nc3cc(C)[nH]n3)n2)CC1

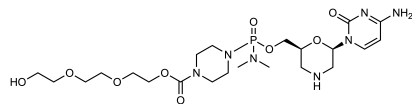
Fig. S10: Shortest retrosynthetic route suggested by ASKCOS for pralsetinib when re-running with the addition of a template-free strategy. Top-1 recommendations by the condition recommender (V1) are shown in blue below each arrow. Ions are written in salt form. ¹The original shortest route contains an alchemical step here (not shown) proposed by the Transformer model using a reactant with an extra nitrogen in the ring that has no way of disappearing in the product; excluding this reactant with one-click filters out all routes with this nonsensical step, and the replacement step from the new shortest route after filtering is shown. ²These two steps are obtained from manual expansion in the IPP upon rejecting the Transformer-proposed step marked with a crossed out arrow. Abbreviations are defined in Table S1.

For the last re-run, we replaced the Pistachio-trained template relevance model with the Pistachio-trained Transformer, thereby combining a template-based expansion strategy with a template-free one; all the other settings remained exactly the same as the baseline. Template-free models are generally less constrained in formulation and *may* generate more creative suggestions at the expense of chemical validity, helping the multi-step search converge towards buyable starting materials.

The shortest synthesis route for newly solved prasetinib is shown in Figure S10 after filtering out a blatantly wrong intermediate with an extra nitrogen in the ring (not shown for succinctness). The inclusion of the Transformer model facilitates the solution of key intermediates **9** and **10**. For **9**, the Transformer model proposed a transesterification before making subsequent simplifying disconnections; while circuitous and seemingly unnecessary from a chemistry standpoint, from an algorithmic standpoint, the template-based model might have happened to generate better recommendations for the methyl ester than it could for the ethyl ester directly. The protection chemistry used in this route could be modified manually as a postprocessing step. There is another unusual suggestion from the Transformer in this route, which is marked with a cross. We chose to discard this step and re-expand manually in the IPP, but recommendations from the Pistachio and Reaxys template models did not propose to disconnect the two rings. We therefore made use of a proprietary template relevance model trained on data from the CAS Content Collection [25], which then successfully proposed a step based on analogy to metal-free cross coupling (e.g., as reported in [26], though missing the necessary NaN_3 and an external oxidant). Cross-referencing against literature precedents reveals that carboxylic acids seem to be out of the substrate scope; the proposal is perhaps more creative than feasible, but the underlying transformation can still be a starting point for a chemist to develop further. Notably, the interactive planning does not have to stop here; we can continue to modify, e.g., by (re-)expanding with different settings until we are satisfied with the route.

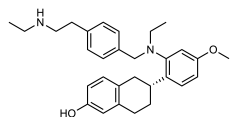
S5.4 Contextualizing the limitation of automatic planning

Even with a few different settings for the automated tree search, some targets still failed to yield retrosynthetic routes. Five such targets are shown in Figure S11. Each may have a different explanation for *why* the Tree Builder fails—missing chemistry, missing building blocks, or insufficient search time—which is not a particular focus for this case study. Interactive planning, as discussed in the Results Section, may be more helpful and insightful, as it offers the user full control over which intermediate to expand and thus the direction of retrosynthesis at a higher level. Moreover, since only one expansion step is performed per click, the user can explore different and more aggressive expansion settings (e.g., more templates, and/or more one-step strategies simultaneously, and/or more tolerant definition of building blocks). Last but not least, the interactive planning mode in ASKCOS is designed to cater even to the most general scenario of synthesis planning, by allowing the user to delete steps, manually add steps (by inputting SMILES as precursors), as well as to add step-wise notes and share routes with others.



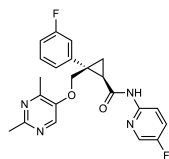
Casimersen:

CN(C)P(=O)(OC[C@@H]1CNC[C@H](n2ccc(N)nc2=O)O1)N1CCN(C(=O)OCCOCCOCCO)CC1



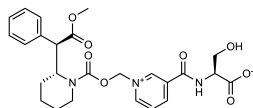
Elacestrant:

CCNCCc1ccc(CN(CC)c2cc(OC)ccc2[C@@H]2CCc3cc(O)ccc3C2)cc1



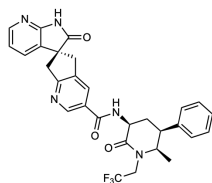
Lemborexant:

Cc1ncc(OC[C@@]2(c3cccc(F)c3)C[C@H]2C(=O)Nc2ccc(F)cn2)c(C)n1



Sordexmethylphenidate:

COC(=O)[C@H](c1cccc1)[C@H]1CCCCN1C(=O)OC[n+]1cccc(C(=O)N[C@@H](CO)C(=O)[O-])c1



Ubrogrepant:

C[C@@H]1[C@H](c2cccc2)C[C@H](NC(=O)c2cnc3c(c2)C[C@@]2(C3)C(=O)Nc3ncccc32)C(=O)N1CC(F)(F)F

Fig. S11: Sample targets for which ASKCOS fails to propose any routes with buyable building blocks within the defined search criteria from all four automatic runs.

References

- [1] Chung, Y., H. Vermeire, F., Wu, H., J. Walker, P., Abraham, M.H., Green, W.H.: Group Contribution and Machine Learning Approaches to Predict Abraham Solute Parameters, Solvation Free Energy, and Solvation Enthalpy. *Journal of Chemical Information and Modeling* **62**(3), 433–446 (2022) <https://doi.org/10.1021/acs.jcim.1c01103> . Publisher: American Chemical Society
- [2] Guan, Y., Coley, C.W., Wu, H., Ranasinghe, D., Heid, E., Struble, T.J., Pattanaik, L., Green, W.H., Jensen, K.F.: Regio-selectivity prediction with a machine-learned reaction representation and on-the-fly quantum mechanical descriptors. *Chemical Science* **12**(6), 2198–2208 (2021) <https://doi.org/10.1039/D0SC04823B> . Publisher: The Royal Society of Chemistry
- [3] Li, S.-C., Wu, H., Menon, A., Spiekermann, K.A., Li, Y.-P., Green, W.H.: When Do Quantum Mechanical Descriptors Help Graph Neural Networks to Predict Chemical Properties? *Journal of the American Chemical Society* **146**(33), 23103–23120 (2024) <https://doi.org/10.1021/jacs.4c04670> . Publisher: American Chemical Society
- [4] Stuyver, T., Coley, C.W.: Quantum chemistry-augmented neural networks for reactivity prediction: Performance, generalizability, and explainability. *The Journal of Chemical Physics* **156**(8) (2022)
- [5] RDKit: Open-Source Cheminformatics Software. <https://www.rdkit.org/> Accessed 2024-07-25
- [6] FastAPI. <https://fastapi.tiangolo.com/> Accessed 2024-07-25
- [7] TorchServe. <https://pytorch.org/serve/> Accessed 2024-07-25
- [8] The ASKCOS wiki. <https://gitlab.com/mlpds.mit/askcosv2/askcos-docs/-/wikis/home> Accessed 2024-07-09
- [9] Vue.js. <https://vuejs.org/> Accessed 2024-09-22
- [10] Vuetify. <https://vuetifyjs.com/> Accessed 2024-09-22
- [11] Cypress. <https://www.cypress.io/> Accessed 2024-09-22
- [12] Pinia: The intuitive store for Vue.js. <https://pinia.vuejs.org/> Accessed 2024-09-22
- [13] VueRouter. <https://router.vuejs.org/> Accessed 2024-09-22
- [14] Chemical Identifier Resolver from NIH. <https://cactus.nci.nih.gov/chemical/structure> Accessed 2024-11-04

- [15] Compilation of CDER New Molecular Entity (NME) Drug and New Biologic Approvals. <https://www.fda.gov/drugs/drug-approvals-and-databases/compilation-cder-new-molecular-entity-nme-drug-and-new-biologic-approvals/> Accessed 2024-10-14
- [16] Anderson, H.J., Loader, C.E., Xu, R.X., Lê, N., Gogan, N.J., McDonald, R., Edwards, L.G.: Pyrrole chemistry. XXVIII. Substitution reactions of 1-(phenylsulfonyl)pyrrole and some derivatives. *Canadian Journal of Chemistry* **63**(4), 896–902 (1985) <https://doi.org/10.1139/v85-149> . Publisher: NRC Research Press
- [17] Elliott, L.D., Berry, M., Orr-Ewing, A.J., Booker-Milburn, K.I.: The Intramolecular Photometathesis of Pyrroles. *Journal of the American Chemical Society* **129**(11), 3078–3079 (2007) <https://doi.org/10.1021/ja070254l> . Publisher: American Chemical Society
- [18] Koovits, P.J., Knowles, J.P., Booker-Milburn, K.I.: Conformationally Driven Two- and Three-Photon Cascade Processes in the Stereoselective Photorearrangement of Pyrroles. *Organic Letters* **18**(21), 5608–5611 (2016) <https://doi.org/10.1021/acs.orglett.6b02829> . Publisher: American Chemical Society
- [19] Boogaard, A.T., Pandit, U.K., Koome, G.-J.: Ring D modifications of ellipticine. Part 1. New ellipticine derivatives from 1-cyano-6-methyllellipticine. *Tetrahedron* **50**(8), 2551–2560 (1994) [https://doi.org/10.1016/S0040-4020\(01\)86971-4](https://doi.org/10.1016/S0040-4020(01)86971-4)
- [20] Ornstein, P.L., Schoepp, D.D., Arnold, M.B., Leander, J.D., Lodge, D., Paschal, J.W., Elzey, T.: 4-(Tetrazolylalkyl)piperidine-2-carboxylic acids. Potent and selective N-methyl-D-aspartic acid receptor antagonists with a short duration of action. *Journal of Medicinal Chemistry* **34**(1), 90–97 (1991) <https://doi.org/10.1021/jm00105a016> . Publisher: American Chemical Society
- [21] Price, D.J., Drewry, D.H., Schaller, L.T., Thompson, B.D., Reid, P.R., Maloney, P.R., Liang, X., Banker, P., Buckholz, R.G., Selley, P.K., McDonald, O.B., Smith, J.L., Shearer, T.W., Cox, R.F., Williams, S.P., Reid, R.A., Tacconi, S., Faggioni, F., Piubelli, C., Sartori, I., Tessari, M., Wang, T.Y.: An orally available, brain-penetrant CAMKK2 inhibitor reduces food intake in rodent model. *Bioorganic & Medicinal Chemistry Letters* **28**(10), 1958–1963 (2018) <https://doi.org/10.1016/j.bmcl.2018.03.034>
- [22] Henry, J.R., Marcin, L.R., McIntosh, M.C., Scola, P.M., Harris Jr, G.D., Weinreb, S.M.: Mitsunobu reactions of n-alkyl and n-acyl sulfonamides-an efficient route to protected amines. *Tetrahedron letters* **30**(42), 5709–5712 (1989)
- [23] Thakkar, A., Selmi, N., Reymond, J.-L., Engkvist, O., Bjerrum, E.J.: “Ring Breaker”: Neural Network Driven Synthesis Prediction of the Ring System Chemical Space. *Journal of Medicinal Chemistry* **63**(16), 8791–8808 (2020) <https://doi.org/10.1021/acs.jmedchem.9b01919> . Publisher: American Chemical Society

- [24] The Pistachio dataset. <https://www.nextmovesoftware.com/pistachio.html>
Accessed 2024-07-08
- [25] The CAS reactions collection. <https://www.cas.org/cas-data/cas-reactions>
Accessed 2024-07-08
- [26] Antonchick, A.P., Burgmann, L.: Direct Selective Oxidative Cross-Coupling of Simple Alkanes with Heteroarenes. *Angewandte Chemie International Edition* **52**(11), 3267–3271 (2013) <https://doi.org/10.1002/anie.201209584> . eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/anie.201209584>