

# Direct Optimization across Computer-Generated Reaction Networks Balances Materials Use and Feasibility of Synthesis Plans for Molecule Libraries

Hanyu Gao, Jean Pauphilet, Thomas J. Struble, Connor W. Coley, and Klavs F. Jensen\*



Cite This: <https://dx.doi.org/10.1021/acs.jcim.0c01032>



Read Online

ACCESS |



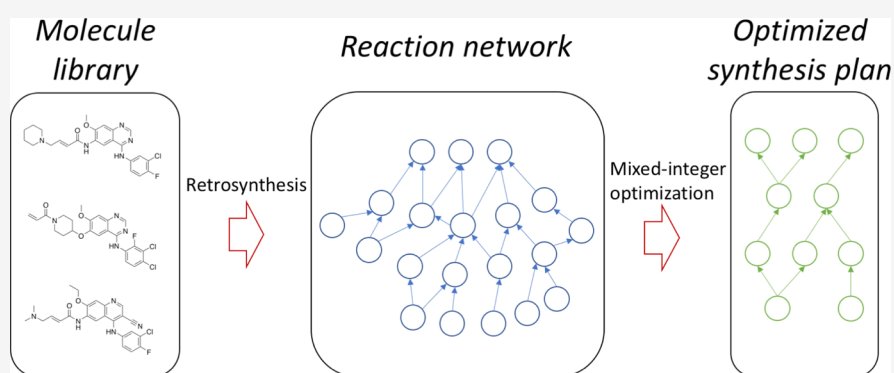
Metrics & More



Article Recommendations



Supporting Information



**ABSTRACT:** The synthesis of thousands of candidate compounds in drug discovery and development offers opportunities for computer-aided synthesis planning to simplify the synthesis of molecule libraries by leveraging common starting materials and reaction conditions. We develop an optimization-based method to analyze large organic chemical reaction networks and design overlapping synthesis plans for entire molecule libraries so as to minimize the overall number of unique chemical compounds needed as either starting materials or reaction conditions. We consider multiple objectives, including the number of starting materials, the number of catalysts/solvents/reagents, and the likelihood of success of the overall syntheses plan, to select an optimal reaction network to access the target molecules. The library synthesis planning task was formulated as a network flow optimization problem, and we design an efficient decomposition scheme that reduces solution time by a factor of 5 and scales to instance with 48 target molecules and nearly 8000 intermediate reactions within hours. In four case studies of pharmaceutical compounds, the approach reduces the number of starting materials and catalysts/solvents/reagents needed by 32.2 and 66.0% on average and up to 63.2 and 80.0% in the best cases. The code implementation can be found at [https://github.com/Coughy1991/Molecule\\_library\\_synthesis](https://github.com/Coughy1991/Molecule_library_synthesis).

## INTRODUCTION

For each new successful pharmaceutical compound, hundreds to thousands of small molecules are typically designed, synthesized, and tested. Synthesizing each molecule separately is an expensive process, both in time and resources. To improve the efficiency of drug discovery, the concept of molecule libraries has been widely adopted in medicinal chemistry.<sup>1–3</sup> A molecule library in the hit-to-lead or lead optimization stage is a collection of compounds that are intended for the same function but exhibit some diversity in structures and will optimize desired properties. Molecule libraries can be compiled in different ways. A library may be a collection of compounds with similar functionality manually extracted from the literature;<sup>4,5</sup> a library may be designed by enumerating possible side groups as decorations of a common core scaffold;<sup>6–8</sup> or a library may be designed using an *in silico* tool for the generation of drug-like compound libraries.<sup>9–13</sup> Accessing all the molecules in a library would enable rapid

exploration of structure–activity relationships during the hit-to-lead or lead optimization phases.

While structure/property-based enumeration and generative models can design molecule libraries comprising many molecules, these molecules might not be easy or even possible to synthesize.<sup>14</sup> Even when they are, finding the most efficient way to synthesize them is not straightforward for at least two reasons. First, for each molecule, there can be multiple possible synthetic pathways, where each pathway involves multiple reaction steps with many possible choices of reaction conditions (e.g., if we consider 3-step pathways with 10

**Received:** September 1, 2020

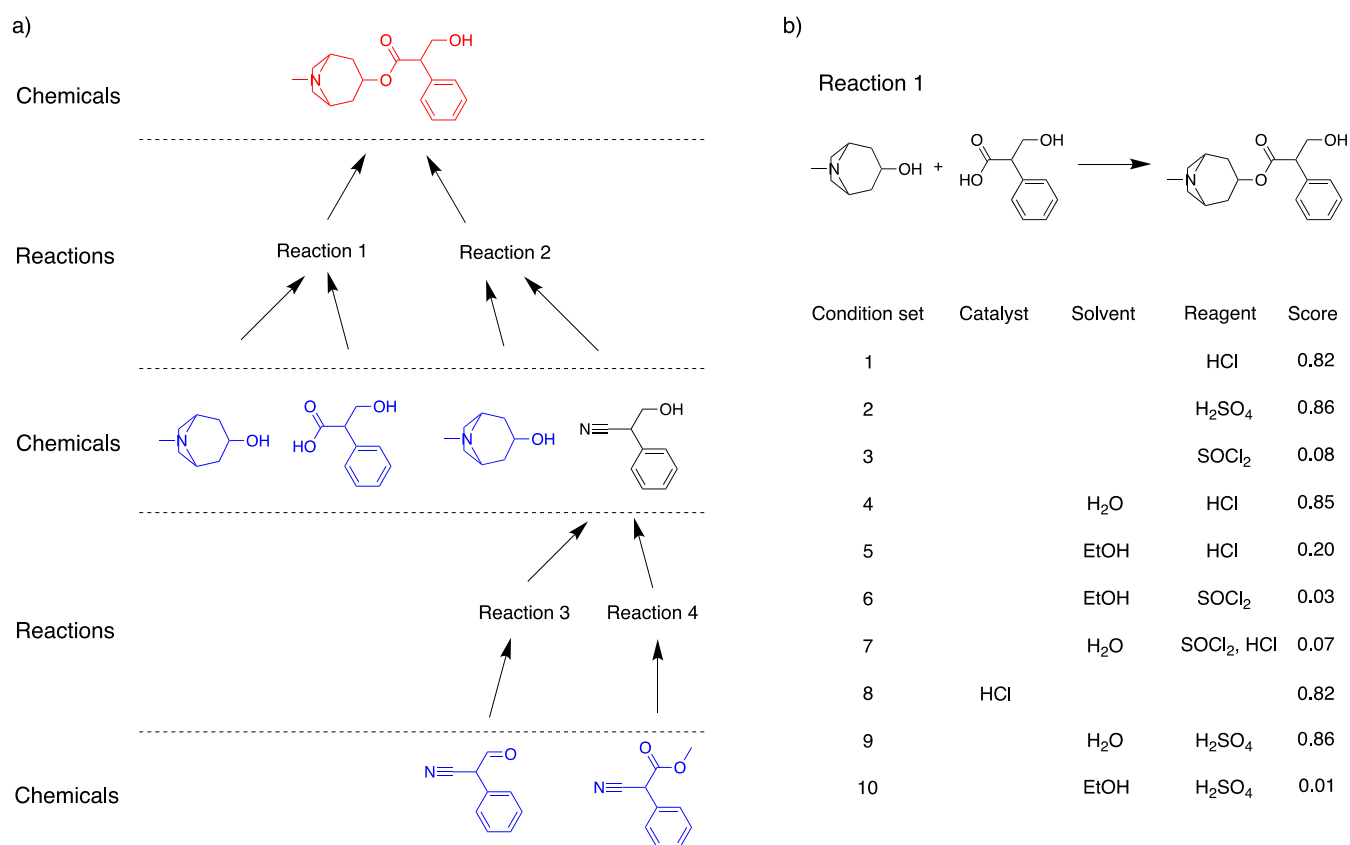


ACS Publications

© XXXX American Chemical Society

A

<https://dx.doi.org/10.1021/acs.jcim.0c01032>  
J. Chem. Inf. Model. XXXX, XXX, XXX–XXX



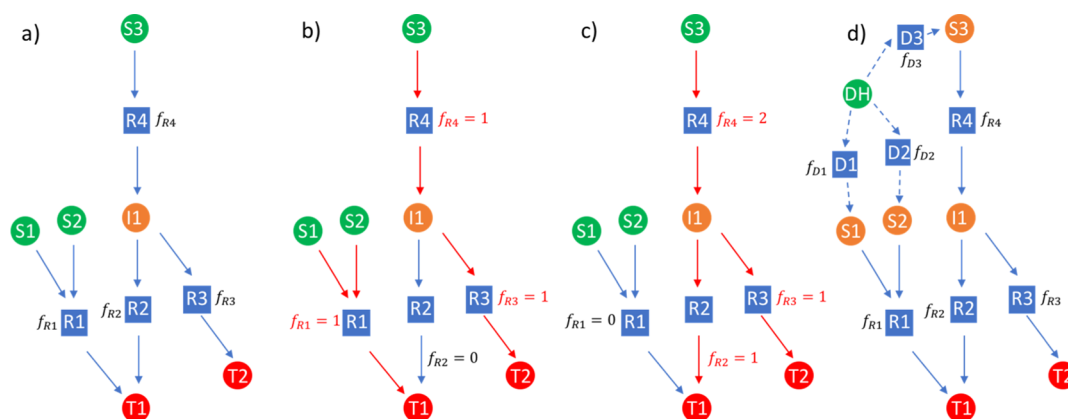
**Figure 1.** (a) Illustration of the two type of nodes—chemicals and reactions and their connectivity in the reaction network (red molecule represents target molecules and blue molecules represent starting materials); (b) each reaction has 10 sets of predicted conditions and reaction evaluation scores associated with each set of conditions as illustrated for reaction 1. Note that condition sets might be degenerate because of the same chemical predicted with different roles (e.g., 1 and 8), and common solvents might be missing in the conditions (e.g., water for 1 and 2). These issues are considered in the Discussion section.

different reactions for each step and 10 possible choices of reagents for each reaction, the number of options is  $10^4$  for one molecule). Second, selecting the most efficient pathway for each molecule separately does not leverage structural similarities between molecules. For instance, it is desirable to share starting materials and to use reactions with similar conditions. Given the combinatorial nature of the problem and the increased number of molecules and reactions to consider, the problem easily grows beyond the capability of manual enumeration. Computer-assisted synthesis planning has progressed quickly in recent years, including retrosynthetic search,<sup>15–19</sup> reaction condition recommendation,<sup>20,21</sup> and reaction outcome prediction.<sup>22,23</sup> These developments have enabled fast construction of chemical reaction networks connecting molecular targets to commercially available starting materials in numerous possible ways. However, for synthesis planning of an entire molecule library, this network comprises hundreds of thousands of reactions, rendering it difficult to fully explore all the possibilities and optimize synthesis plans.

The few studies that have tackled synthesis planning for multiple targets truncate the information and the search space in different ways. Molga *et al.*<sup>24</sup> developed a method to search for pathways for multiple targets on the same growing network. In order to promote the use of common intermediates and similar reactions, they modified the retrosynthetic search—penalizing the search when additional reaction types are encountered. While they demonstrated this approach on analogue of compounds with variation of some functional

groups, the penalization of additional reaction types potentially limits the scope of the reactions that can be explored when constructing the reaction network. In previous work from our group,<sup>25</sup> we selected a limited number of pathways for each molecule and then solved an optimization problem to select a combination of pathways for the entire library, accounting for the number of unique chemicals and the perceived likelihood of success for the syntheses. This explicit, though incomplete, enumeration of the pathways for each molecule reduced the size of the subsequent optimization problem.<sup>25</sup> However, the approach did not take full advantage of the entire reaction network and could not consider multiple different reaction conditions (for the same reaction) because of the increased combinatorial complexity.

Herein, we present an alternative approach to solve the problem directly by formulating and solving a mixed-integer linear programming (MILP) problem on the entire reaction network, instead of enumerating the trees explicitly. This approach avoids information loss because of the incomplete enumeration of reaction pathways and increases the flexibility of exploring different sets of reaction conditions. This additional flexibility, however, comes at the expense of increased combinatorial complexity that we address in two steps. First, we formulate the reaction pathway selection problem as a network flow problem that can be directly processed in computational solvers. Second, we design a decomposition strategy to accelerate the numerical convergence of our algorithm and improve the solutions found,



**Figure 2.** Graphical interpretation of the flow variables on the reaction network. (a) Example reaction network, where there are 6 chemicals, T1 and T2 are targets (red circles), S1, S2, and S3 are starting materials (green circles), and I1 is an intermediate (orange circle). R1–R4 are reactions (blue squares). The  $f$ s are the flow variables defined over these four reactions. (b,c) Two sets of possible pathways to access the two targets. In (b), R1 is used to access T1. R4 and R3 are used to access T2. In (c), R4 and R2 are used to access T1, and R4 and R3 are used to access T2. The directed edges (arrows) are highlighted in red. (d) Introduction of a dummy head node (DH) and dummy reactions (D1–3). Note that by doing this, the starting materials (S1–3) can be treated the same way as intermediates (orange circles), and the choice of starting materials is equivalent to the choice of dummy reactions.

especially for large-size problems. We test the approach on four molecule libraries to demonstrate its effectiveness.

## METHODS

**Retrosynthesis.** As a first step, we performed retrosynthetic analysis for multiple targets by using the ASKCOS platform, as implemented in Coley *et al.*,<sup>16</sup> to search for possible retrosynthetic pathways. Specifically, reaction templates were recursively applied to a target molecule to break it down into starting materials. The priority of a template application was given by a classifier trained to predict the most relevant reaction templates for a given product molecule. An upper confidence bound tree search balances exploitation and exploration. The terminal nodes in the tree search were defined as chemicals that are small enough (no more than 10 carbon, 3 nitrogen, and 5 oxygen atoms). The search time was limited to 60 s for every target. The search for different targets happened on the same graph to facilitate the exploitation of previously explored reactions. While these settings are used in this work, there are other options which the users can tune when using the ASKCOS platform (*e.g.*, different stop criteria, a longer search time) for improved performance, which can be case-dependent. Because the reaction evaluation model does not take stereochemistry into consideration,<sup>22</sup> for consistency, we did not account for chirality in retrosynthesis analysis.

The result of the retrosynthetic analysis is a directed graph where two types of nodes, chemicals and reactions, are connected alternately (Figure 1); the child nodes of a chemical are possible reactions for producing the chemical, and the child nodes of a reaction are the reactants of that reaction. A condition recommendation model predicted the 10 most suitable sets of reaction conditions for every reaction,<sup>20</sup> and then, a reaction evaluation model estimated the likelihood of success of the desired reaction under each set of reaction conditions.<sup>22</sup>

Note that the reaction evaluation model was developed using the USPTO database, which is publicly available, while other models were developed using proprietary Reaxys data, which has a wider coverage of the chemical reaction space. While it would be desirable to unify the data sources and

compare their effect on the solution of the retrosynthesis analysis, it is beyond the scope of this work.

**Optimization.** After obtaining the chemical-reaction graph through retrosynthetic analysis and reaction evaluation, the task was to select an optimal set of reactions to access all the targets in the molecule library from the set of starting materials. A multiobjective MILP problem was formulated to minimize the resources needed for the syntheses of all the target molecules as well as the possibility of failure of the synthesis plan. A numerical strategy based on the Benders decomposition method was developed to efficiently solve the MILP problem. The optimal synthesis plan could be obtained and compared with planning synthesis for each target molecule separately.

**Objective Function.** We considered a combination of three concurrent objectives. The first objective was to minimize the number of starting materials used in all the syntheses, so as to simplify inventory and supply chain management. The second objective was to minimize the chemicals used as reaction conditions (catalysts, solvents, and reagents—referred to as “C/S/R” hereafter). This objective promotes similar type of reactions that are likely to happen under similar conditions. The third objective was to minimize the probability at which the synthesis plan might fail. From the reaction evaluation model, we associated each reaction with a penalty, indicating whether the reaction had low probability of success. By minimizing the overall penalty, we encouraged the selection of fewer reactions and the selection of the more plausible ones. The mathematical formulation of the objective function for multimolecule synthesis planning is as follows

$$\min_{s_i, c_k, o_{mn}} \lambda_1 \sum_{i \in S} s_i + \lambda_2 \sum_{k \in K} c_k + \lambda_3 \sum_{m \in R} \sum_{n \in \{1, 2, \dots, 10\}} Q_{mn} o_{mn} \quad (1)$$

where  $\lambda_1$ ,  $\lambda_2$  and  $\lambda_3$  are positive weights to trade-off the different objectives: (1) The total number of starting materials  $\sum_{i \in S} s_i$ , where  $s_i$  are binary variables indicating whether starting material  $i$  is selected and  $S$  is the set of all possible starting materials. (2) The total number of C/S/R chemicals  $\sum_{k \in K} c_k$ , where  $c_k$  are binary variables indicating whether a chemical  $k$  (catalyst, solvent or reagent) is used in any reaction of the plan

and  $K$  denotes the set of all possible C/S/R. (3) The overall penalty  $\sum_{m \in R, n=1, \dots, 10} Q_{mn} o_{mn}$ , where  $o_{mn}$  are binary variables indicating whether reaction  $m$  chooses its  $n$ th option of reaction conditions,  $R$  is the set of all reactions, and  $Q_{mn}$  are precalculated parameters that represent the penalty associated with choosing option  $n$  for reaction  $m$ . Here, we defined them as

$$Q_{mn} = \min\left(\frac{1}{\text{score}_{mn}}, 20\right) \quad (2)$$

where  $\text{score}_{mn}$  is the reaction evaluation score between 0.0 and 1.0 obtained from the reaction evaluation model. With this definition, we apply a higher penalty to reactions with lower scores, which can be understood as a coarse estimate of the probability of experimental success for that reaction. The values of  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$  can be changed by users of the model based on their optimization goals. For example, to minimize the risk of finding unsuccessful reactions, a user could increase  $\lambda_3$ .

**Constraints.** In order to find the optimal synthesis plan, we first needed to ensure that we searched among valid synthesis plans. We defined a valid synthesis plan as a subset of reactions from the network that connected all the targets back to some commercially available starting materials, with one set of conditions chosen for each reaction. To mathematically describe a valid synthesis plan, we formulated the selection of reactions as a network flow optimization problem. Intuitively, starting materials could be thought of as the source of the flow, and target molecules were the final sink for the flow. Reactions were treated as flow channels that could be switched on or off. A valid synthesis plan thus required that there was flow reaching all target molecules from starting materials. Rigorously, we defined the flow through a reaction as the number of times this reaction was used in synthesizing all targets and denoted it  $f_m$ . Figure 2a presents a stylized reaction network with 2 targets and 4 potential reactions. Figure 2b,c represents two different pathways to synthesize the targets, with the corresponding flow variables. In Figure 2b, reaction R1 is used to synthesize target T1, and reactions R4 and R3 are used to synthesize T2. R2 is not used. Therefore, the flow variables  $f_{R1}$ ,  $f_{R4}$ , and  $f_{R3}$  all take the value 1 and  $f_{R2}$  equals zero. In Figure 2c, reactions R4 and R2 are used to synthesize T1, and reactions R4 and R3 are used to synthesize T2. In this case, R4 is used twice, so  $f_{R4} = 2$ ,  $f_{R2} = f_{R3} = 1$ , and  $f_{R1} = 0$ . We introduced a dummy head node so that starting materials could be modeled in the same way as intermediates (Figure 2d).

Below are a few sets of constraints that ensure that the decision variables impacting the objective, namely, the set of starting materials  $s_i$ , of C/S/R chemicals  $c_k$ , and reaction-condition  $o_{mn}$  are not chosen arbitrarily but rather correspond to a valid synthesis plan, where all the targets are produced in some way.

**Flow Constraints.** Each chemical can either be synthesized or produced by a reaction, we denote by  $\text{CH}_i$  the set of reactions that produce chemical  $i$  (child nodes) and  $\text{PR}_i$  the set of reactions that consume it (parent nodes). An admissible set of reactions must produce all the targets, that is, for all targets  $i$ , we must have  $\sum_{m \in \text{CH}_i} f_m - \sum_{m' \in \text{PR}_i} f_{m'} = 1$ . Similarly, for all the intermediates, the inflow must equal the outflow. We treated starting materials in the same way as intermediates, after the introduction of a dummy head node and dummy reactions

(DH and D1–3 in Figure 2d). Consequently, the selection of starting materials was equivalent to the selection of the dummy reactions. These flow constraints can be concisely written as follows

$$\sum_{m \in \text{CH}_i} f_m - \sum_{m' \in \text{PR}_i} f_{m'} = b_i, \quad \forall i \in S \cup I \cup T, \text{ with } b_i = \begin{cases} 1, & \text{if } i \in T \\ 0, & \text{if } i \in S \cup I \end{cases} \quad (3)$$

where  $S$ ,  $I$ , and  $T$  denote the set of starting material, intermediates, and targets respectively. We also specified that flow variables are non-negative continuous variables. Yet, because net inflows  $b_i$  are binary, the flow amount can only take integer values.

**Reaction Selection Constraints.** As stated above, the flow variables indicate how many times a reaction is used by different targets, but in the final objective function, we only want to know whether the reaction is used by any target.

We modeled this logic through the constraint

$$C\phi_m \geq f_m, \quad \forall m \in R \cup D \quad (4)$$

where  $R$  is the set of reactions and  $D$  is the set of dummy reactions. This set has a one-to-one mapping to the set of starting materials  $S$ .  $\phi_m$  is a new binary variable encoding whether or not to select reaction  $m$ .  $C$  is a sufficiently large constant. In this case,  $f_m$  cannot be greater than the total number of targets, so  $C$  can be set equal to the total number of targets. Note that the binary variables for selecting dummy reactions  $\phi_i$ ,  $i \in D$  are equivalent to the variables for selecting starting materials,  $s_i$ ,  $i \in S$ .

**Condition Set Selection Constraints (Sum-to-One Constraints).** For every reaction, we predicted ten sets of conditions, but in the end, only one set of conditions could be chosen, which was imposed through the following constraints

$$\sum_{n \in \{1, 2, \dots, 10\}} o_{mn} = \phi_m, \quad \forall m \in R \quad (5)$$

**Chemical Selection Constraints for Reaction Conditions.** Based on the sets of the conditions we chose, the individual chemicals (catalysts, solvents, and reagents) involved can be determined, as captured through the constraints

$$\begin{aligned} \forall m \in R, \quad n \in \{1, 2, \dots, 10\} \\ \forall k \in K, \\ \text{if } k \text{ is used in the } n\text{th option of conditions for reaction } m \\ c_k \geq o_{mn} \end{aligned} \quad (6)$$

**Final Formulation.** All in all, the final mixed-integer optimization problem is stated as follows



$$\begin{aligned}
& \min_{\phi_m, c_k, o_{mn}} \lambda_1 \sum_{i \in D} \phi_i + \lambda_2 \sum_{k \in K} c_k + \lambda_3 \sum_{m \in R} \sum_{n \in \{1, 2, \dots, 10\}} Q_{mn} o_{mn} \\
& \text{s. t. } f_m \geq 0, \\
& \phi_m, c_k, o_{mn} \text{ binary,} \\
& \sum_{m \in \text{CH}_k} f_m - \sum_{m' \in \text{PR}_k} f_{m'} = b_i, \quad \forall i \in S \cup I \cup T, \\
& C\phi_m \geq f_m, \quad \forall m \in R \cup D, \\
& \sum_{n \in \{1, 2, \dots, 10\}} o_{mn} = \phi_m, \quad \forall m \in R, \\
& c_k \geq o_{mn}, \quad \forall m, n, k \text{ such that } k \text{ is in the } n \\
& \quad \text{th option for reaction}
\end{aligned} \tag{7}$$

**Solution Method.** The optimization problem formulated above is an MILP that can be recognized as a variant of the network design problem—a notoriously hard family of MILP problems.<sup>26</sup> It can be directly fed to an optimization solver. We chose Gurobi<sup>27</sup> in this work, which has been demonstrated to perform the best on a majority of test problems.<sup>28</sup> The problem can be solved quickly (within 100 s) for small-size molecule libraries (2–7 targets) and reaction networks (less than 1000 reactions). However, as the size of the molecule library and the size of the reaction network increase, solving the problem to optimality becomes computationally prohibitive. Even obtaining a reasonable optimality gap, which measures the difference between the current best solution and a lower bound of the optimal solution (which is typically obtained from solving the linear relaxation of an integer optimization problem), can be challenging. In the worst case, the optimization solver will explore all potential solutions which grow exponentially in the number of targets and the number of steps per pathway.

Hence, for the large-scale scenarios, we designed a decomposition strategy to accelerate the convergence of the solver and obtain better quality solutions. The structure of the problem makes it suitable to use a Benders decomposition strategy,<sup>29</sup> as it can be divided into two stages that can be solved iteratively. We refer to Rei *et al.*<sup>30</sup> for a comprehensive review of Benders decomposition.

**Benders Decomposition.** The synthesis planning problem can be decomposed into two stages: in the first stage, the reactions of the synthesis plans are selected, while the second stage focuses on finding the best set of conditions for these reactions. Given a set of reactions  $f_m$ ,  $m \in \text{DUR}$ , the best set of conditions can be found by solving the second-stage problem (or sub-problem)

$$\begin{aligned}
F(f_m, \phi_m) &= \min_{c_k, o_{mn}} \lambda_2 \sum_{k \in K} c_k + \lambda_3 \sum_{m \in R} \sum_{n \in \{1, 2, \dots, 10\}} Q_{mn} o_{mn}, \\
& \text{s. t. } c_k, o_{mn} \text{ binary,} \\
& \sum_{n \in \{1, 2, \dots, 10\}} o_{mn} = \phi_m, \quad \forall m \in R, \\
& c_k \geq o_{mn}, \quad \forall m, n, k \text{ such that } k \text{ is in the } n \\
& \quad \text{th option for reaction}
\end{aligned} \tag{8}$$

With this notation, the overall problem (or master problem) can be summarized as follows

$$\begin{aligned}
& \min_{\phi_m, f_m} \lambda_1 \sum_{i \in D} \phi_i + F(f_m) \\
& \text{s. t. } f_m \geq 0, \quad \phi_m \text{ binary,} \\
& \sum_{m \in \text{CH}_k} f_m - \sum_{m' \in \text{PR}_k} f_{m'} = b_i, \quad \forall i \in S \cup I \cup T, \\
& C\phi_m \geq f_m, \quad \forall m \in R \cup D
\end{aligned} \tag{9}$$

We implemented Benders decomposition in the following way: In the subproblem, we relaxed the constraint that  $c_k, o_{mn}$  are binary and considered them as continuous variable between 0 and 1 instead. As a result, one can show that the function  $F(f_m, \phi_m)$  is convex. In particular, strong duality applies, and  $F(f_m, \phi_m)$  can be formulated as a maximization problem, that is,

$$\begin{aligned}
F(f_m, \phi_m) &= \max_{u_m, x_k, y_{mn}, v_{mnk}} - \sum_{m \in R} \hat{\phi}_m u_m - \sum_{k \in K} x_k \\
& \quad - \sum_{m \in R} \sum_{n \in \{1, 2, \dots, 10\}} y_{mn} \\
& \text{s. t. } x_k + \lambda_2 - \sum_m \sum_n v_{mnk} \geq 0 \quad \forall k \in K \\
& y_{mn} + \lambda_3 Q_{mn} + u_m + \sum_k v_{mnk} \geq 0 \\
& \forall m \in R, n \in \{1, 2, \dots, 10\}
\end{aligned} \tag{10}$$

In Benders decomposition, the master problem is solved by replacing  $F(f_m, \phi_m)$  by a piece-wise linear lower approximation, namely, we solve

$$\begin{aligned}
& \min_{\phi_m, f_m} \lambda_1 \sum_{i \in D} \phi_i + z \\
& \text{s. t. } f_m \geq 0, \\
& \phi_m \text{ binary,} \\
& \sum_{m \in \text{CH}_k} f_m - \sum_{m' \in \text{PR}_k} f_{m'} = b_i, \quad \forall i \in S \cup I \cup T, \\
& C\phi_m \geq f_m, \quad \forall m \in R \cup D, \\
& z \geq [\text{piece-wise lower approximation of } F]
\end{aligned} \tag{11}$$

Once we have solved the master problem, we use the subproblem to provide feedback to the master problem. Indeed, for the returned solution, we solve  $F(f_m, \phi_m)$  as a maximization problem and obtain a new linear outer-approximation of  $F$

$$z \geq - \sum_{m \in R} \hat{u}_m \phi_m - \sum_{k \in K} \hat{x}_k - \sum_{m \in R} \sum_{n \in \{1, 2, \dots, 10\}} \hat{y}_{mn} \tag{12}$$

These two problems, the master and the sub problem, are both much easier to solve than the original problem. For continuous linear optimization, this decomposition strategy is guaranteed to converge to the same optimal solution as the original solution and usually demonstrates to be much faster. In this problem, because there are binary variables in the original subproblem and relaxed this constraint, the decomposition is not guaranteed to fully close the optimality gap, but

we found that it indeed converged to improved solutions for large-scale problems.

**Overall Numerical Strategy.** In practice, we found that an off-the-shelf solver was able to find reasonable heuristic solutions quickly and then slowly reduced the optimality gap (by improving the solution or the lower bound). On the other hand, the decomposition method was much faster at improving the lower bound but needed a good initial feasible solution for the master problem to be effective. Therefore, we implemented this decomposition method in the following procedure, with the Gurobi solver:

- 1 Solve the original problem directly using Gurobi;
- 2 Terminate the solution process after 100 s without improvement in the optimal solution (if not solved to optimality). Record the current best solution;
- 3 Initiate the decomposition method with the current best solution, generating optimality cuts at any feasible solution found;
- 4 Terminate the solution if (a) optimality gap is less than 1%, (b) solution time is over an hour and optimality gap is less than 20%, or (c) solution time is greater than 24 h.

The logic behind the termination criteria is that if the problem proves to be difficult (*i.e.*, it has been solved for an hour without converging), we loosen the termination condition to 20% optimality gap. Also, note that in this implementation, the decomposition method would be automatically triggered if there were signs that the original problem could be difficult to solve (unable to solve within 100 s). Therefore, the model user does not have to decide whether or not to use the decomposition method.

The code implementation can be found at [https://github.com/Coughy1991/Molecule\\_library\\_synthesis](https://github.com/Coughy1991/Molecule_library_synthesis). All computations were performed on a dual Intel Xeon(R) CPU E6-2690@2.9 GHz processors, and the computational times reported are the real time elapsed (wall time).

## RESULTS

This approach was applied to drug-like chemical libraries to optimize the synthesis plans for all molecules in the libraries. We demonstrated the method on four molecule libraries: (1) tamatinib and fostamatinib; (2) a library of indole analogs; (3) a library of molecules that are similar to dacomitinib; and (4) a library of derivatives of *k*-opioid agonist ICI-199441. The four case studies included 2, 7, 7, and 48 targets, respectively. The resulting reaction networks for all case studies are summarized in Table 1. The number of reactions ranges from around 500 to more than 7000.

**Table 1. Summary of the Retrosynthesis Analysis Results for the Four Case Studies**

molecule library	num. of targets	num. of starting materials	num. of Intermediates	num. of reactions	num. of C/S/R <sup>a</sup>
tamatinib	2	195	224	568	368
tryptophans and indoles	7	381	418	1029	628
dacomitinib	7	534	1099	2397	522
ICI-199441	48	1033	2805	7151	918

<sup>a</sup>The number of catalysts/solvents/reagents.

We evaluated a synthesis plan based on a weighted combination three concurrent objectives: minimizing the number of unique starting materials, the number of C/S/R needed, and the probability at which the synthesis plan might fail. For each library, we compared the solution obtained by choosing the “best” pathway for each target individually, that is, the separate planning solution, with combined synthesis plans obtained by solving the optimization formulation, and for different weights on the three objectives (1:1:1, 10:1:0.1, and 0.1:1:10, representing the number of starting materials: number of C/S/R: total reaction penalty. See definition of the objectives in the [Methods](#) section).

Table 2 summarizes results from each retrosynthesis analysis. On average, our optimized synthesis plan reduces the number of starting materials and C/S/R by 32.2 and 66.0%, respectively. In terms of computational burden, our optimization formulation can be solved within seconds for the smallest libraries and within hours for the largest instances. In particular, the proposed decomposition algorithm significantly improves scalability to larger-size libraries and reduces computational time by a factor of 5–10 compared to commercial solvers (Figure S1). This demonstrates that Benders decomposition drastically improves scalability of the optimization approach. The case studies are discussed in the following sections based on the results of combined synthesis planning with 1:1:1 weighting. Results with other weighting factors (10:1:0.1 and 0.1:1:10) are shown in the [Supporting Information](#) (Figures S2–S5).

### Case 1. Syntheses of Tamatinib and Fostamatinib.

We investigated the combined synthesis planning for two molecules, tamatinib and fostamatinib (red molecules shown in Figure 3). These two molecules are very similar in structure (Tanimoto similarity of radius 2 Morgan fingerprint = 0.816), and tamatinib can be viewed as a substructure of fostamatinib. After running retrosynthesis for these two molecules, the reaction network included 568 reactions, 195 starting materials, 224 intermediates, and 368 catalysts/solvents/reagents. In the separate synthesis planning, the pathway chosen for these two targets do not have much overlap. In the combined synthesis planning, however, the synthesis of fostamatinib takes advantage of the synthesis of tamatinib and then further functionalizes tamatinib to obtain fostamatinib. Also, there are two S<sub>N</sub>2 reactions that share the same solvents and reagents (Figure 3).

### Case 2. Synthesis of a Library of Indole Analogues.

In this case, we tested the model on a library of indoleamine 2,3-dioxygenase 1 inhibitors adapted from Figure 3 of ref 4. One of the eight compounds presented in the original reference was so small that it qualified as a starting material, so we excluded it and performed the analysis on the remaining seven molecules (Figure 4a). Stereochemistry was not included in the retrosynthesis analysis as described in the [Methods](#) section. There was a larger structural variability across molecules, but they did show some overlapping substructures. For this library, we constructed a reaction network of 1029 reactions, 381 starting materials, 418 intermediates, and 628 catalysts/solvents/reagents.

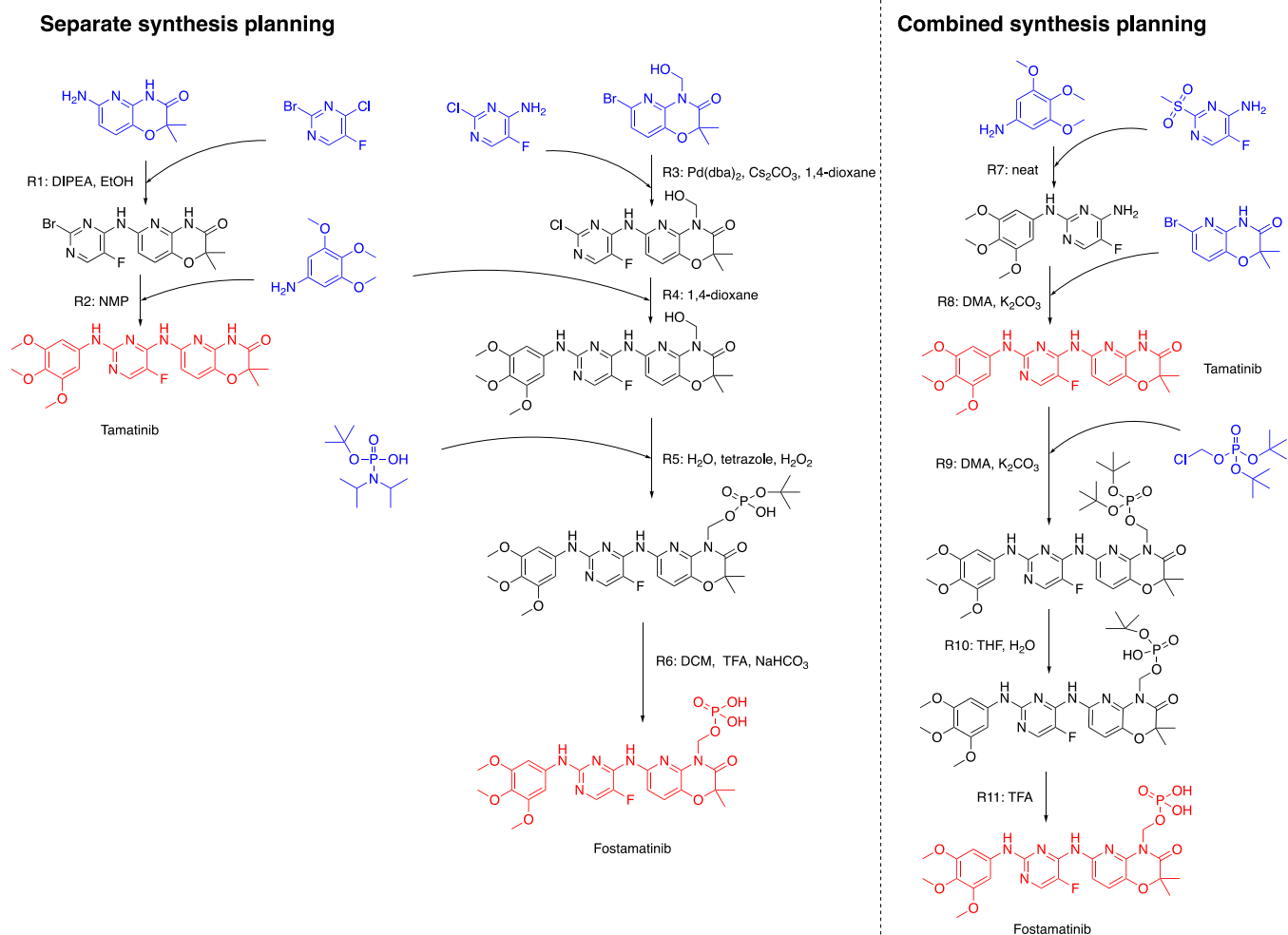
The pathways resulting from separate and combined synthesis planning for all seven targets are shown in Figure S3 in the [Supporting Information](#). Here for better visualization, we presented three molecules that had significant overlap when running combined synthesis planning in Figure 4b. These three syntheses shared the same intermediate (the product of R5),

Table 2. Summary of the Optimization Results for the Four Case Studies

molecule library	number of targets	separate synthesis planning		combined synthesis planning								
		/		1:1:1 <sup>c</sup>			10:1:0.1			0.1:1:10		
		SM <sup>a</sup>	C/S/R <sup>b</sup>	SM	C/S/R	time/s	SM	C/S/R	time/s	SM	C/S/R	time/s
tamatinib	2	6	12	4	5	1	4	3	2	4	7	1
tryptophans and indoles	7	19	25	16	12	4	12	11	23	17	18	2
dacomitinib	7	18	30	15	6	18	14	9	354 <sup>d</sup>	16	6	5
ICI-199441	48	38	42	14	11	463 <sup>d</sup>	13	9	12,361 <sup>e</sup>	21	12	79

<sup>a</sup>The number of starting materials. <sup>b</sup>The number of catalysts/solvents/reagents. <sup>c</sup>The weighting factors for the three different objectives.

<sup>d</sup>Decomposition used. <sup>e</sup>Decomposition used; solved to 20% optimality gap.



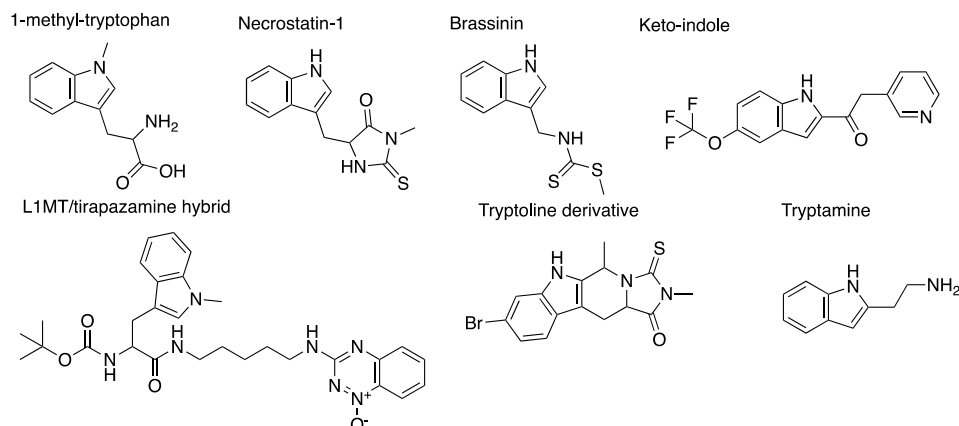
**Figure 3.** Comparison of the synthetic routes found for the two targets in the separate and combined synthesis planning. Molecules highlighted in blue are starting materials; those in red are target molecules. Text next to the arrows takes the format “reaction index: reaction conditions”. Chemistry steps possibly needing further investigation: S<sub>N</sub>Ar reactions potentially have selectivity issues (R1, R2); R7: reaction is possible but similar literature precedence used a strong base.<sup>31</sup> R9: there might be side reactivities of other amines with the phosphate ester; R10 and R11 represent two deprotection steps of the same group so they might be combined in one step.

which can facilitate the development of these syntheses. This intermediate being shared by multiple targets also indicates its potential to be diversified to other molecules. The reaction conditions are also greatly simplified compared to the separate synthesis planning. Moreover, most reactions use common solvents and reagents.

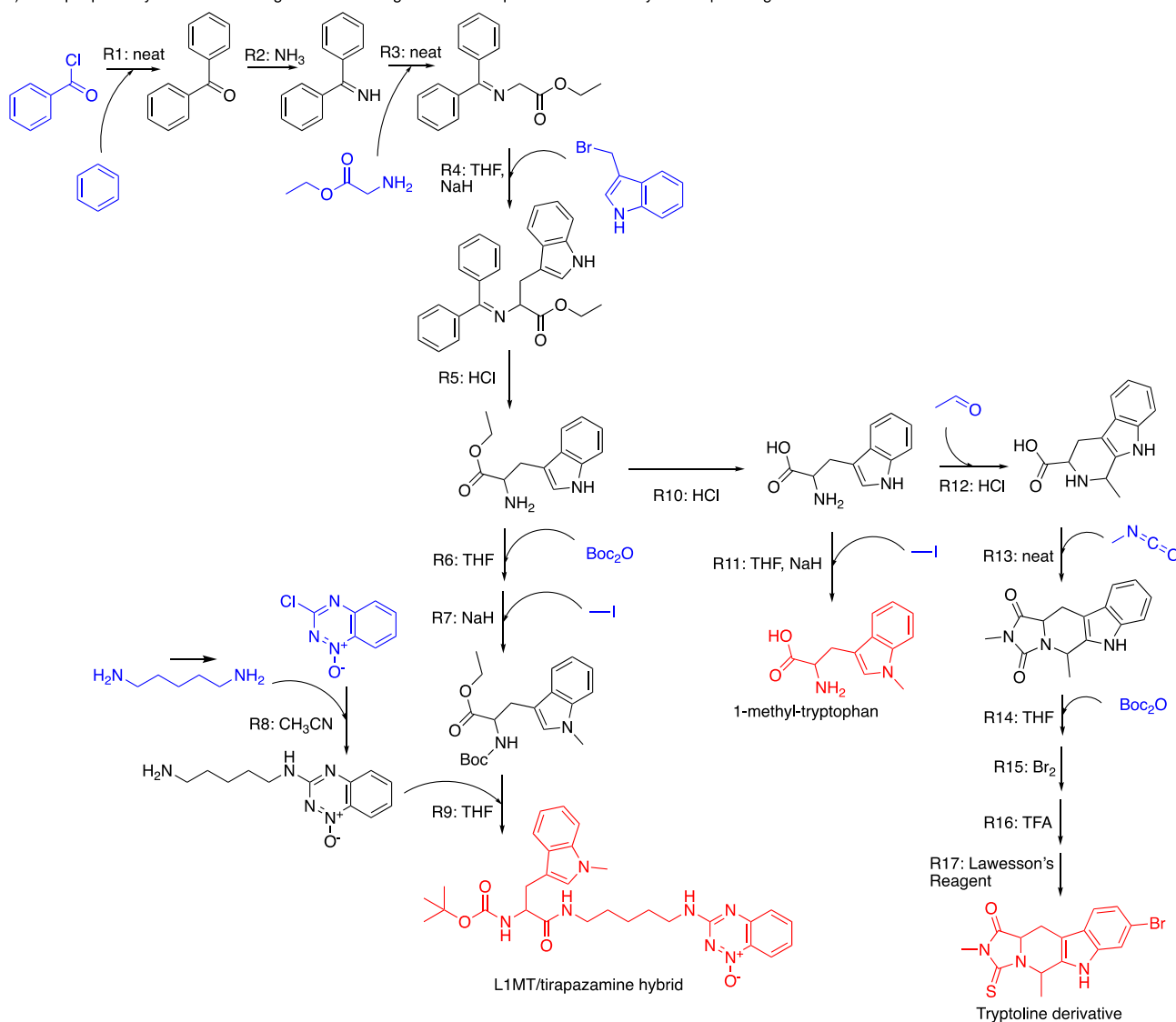
**Case 3. A Library of Molecules That Are Similar to Dacomitinib.** Dacomitinib is a drug that was approved in 2018. We performed a similarity search in the Drugbank

database (<https://www.drugbank.ca/>) to identify other molecules with a similarity score of 0.7 or higher and obtained a library of seven molecules (red molecules in Figure 5). While it has the same number of target molecules as the previous case study, the retrosynthesis analysis resulted in a reaction network of 2397 reactions, 534 starting materials, 1099 intermediates and 533 catalysts/solvents/reagents, which was significantly larger than the previous case study. In this case, when the weighting factors are 10:1:0.1, the problem becomes

a) all target molecules in this library

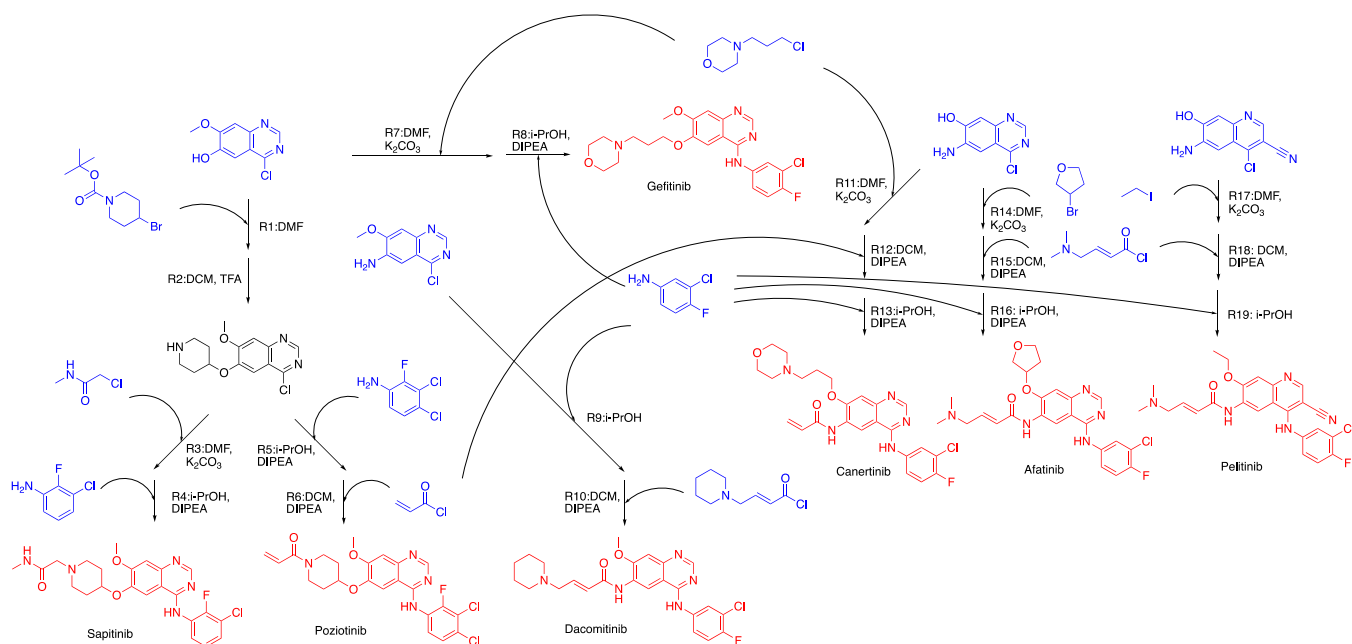


b) example pathways for selected targets that have significant overlap in the combined synthesis planning



**Figure 4.** (a) Structures of tryptophan and indole analogs used in this case study. (b) Examples of three molecules that share common intermediates in the combined synthesis planning. Molecules highlighted in blue are starting materials; those in red are target molecules. Only key intermediates are shown. Text next to the arrows takes the format “reaction index: reaction conditions”. Full pathways are shown in Figure S3. Chemistry steps possibly needing further investigation: comparing R6 + R7 and R11, they are methylating the same amine while one uses protection and the other do not, so it is likely that R10 and R11 would need to be swapped in this sequence. R15 and R17 might have site selectivity issues because there are multiple sites where the same reaction can happen.





**Figure 5.** Optimal reaction network selected in the combined synthesis planning for the seven targets in case 3. Molecules highlighted in blue are starting materials; those in red are target molecules. Only key intermediates are shown. Text next to the arrows takes the format “reaction index: reaction conditions”. Full pathways are shown in Figure S4 in the [Supporting Information](#). Chemistry steps possibly needing further investigation: in R17, there might be selectivity issues because ethyl iodide might also react with the free aniline.

challenging to solve directly. It took 2011 s to reach 1% optimality gap. With the decomposition strategy applied, computational time to reach the same accuracy reduces by a factor of 5, down to 354 s.

Figure 5 shows the pathways for these seven targets from combined synthesis planning. It can be seen that the pathways are highly interconnected, with many starting materials being shared by two or more targets. In this case, the reaction conditions were simplified to only six distinct chemicals compared to 30 in the separate synthesis planning.

**Case 4. A Library of Derivatives of *k*-Opioid Agonist ICI-199441.** In this case, we explored a molecule library with a larger size. It was a case study also investigated by Molga *et al.*<sup>24</sup> All molecules are enumerated derivatives of compound ICI-199441, allowing variations on four different sites (Figure 6a). A total of 48 targets were present in this library. In the report by Molga *et al.*,<sup>24</sup> their algorithm were able to identify the pathways for the top-five most accessible targets, which is likely the result of biasing the search toward common reactions. Here, we attempted to identify the maximum overlap for all the targets. In our work, the reaction network constructed for this library included 7151 reactions, 1033 starting materials, 2805 intermediates, and 918 catalysts/solvents/reagents. For this problem size, the MILP was much more challenging to solve. For weighting factors of 1:1:1, the solver took 793 s to reach 1% optimality gap. For weighting factors of 10:1:0.1, the optimality gaps reached after 24 h using the raw commercial solver Gurobi were 22.3%. The decomposition strategy, however, solves the problem with weights 1:1:1 to 1% optimality gap in 463 s and the instance with weighting factors 10:1:0.1 to 20% optimality gap in 12,361 s (nearly 3.5 h). These results together with what was observed in case 3 demonstrated the effectiveness of the decomposition method in improving the computational performance of the optimization algorithm on large-scale instances.

Figure 6b shows examples of molecules whose synthetic routes identified by the combined synthesis planning significantly overlap. They are the four molecules with the four different variations in  $R_2$  group and the same  $R_1$ ,  $R_3$ , and  $R_4$  groups. The reaction transformations in the synthetic routes for the four targets are very similar, except for two steps in the first synthesis (reactions R2 and R3 are different from R7, R11, and R15). For the same type of reactions, usually the same condition is shared among them (e.g., reactions R4, R8, R12, and R16).

## DISCUSSION

**Scalability.** There are three major components of the computational cost. For retrosynthesis, the computational time scales linearly with the number of targets in the library. If we use the 60 s expansion time for each target, it would be possible to analyze hundreds of targets within hours.

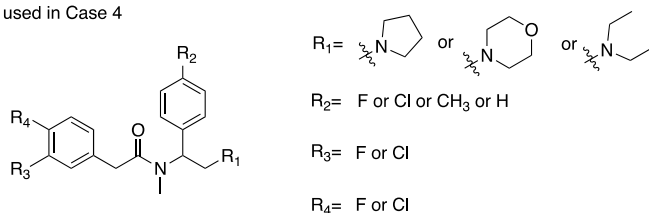
For condition prediction and reaction evaluation, while it is relatively slow per reaction (10 s on an NVIDIA GeForce GTX 1080 GPU), this step has the potential of being parallelized, which could reduce the computational cost given sufficient computing resources.

For optimization, as demonstrated in case 4, it would be challenging to solve large problems to optimality, but within a reasonable period of time, a good heuristic solution can usually be found. Afterwards, the majority of time is spent on improving the lower bound. Therefore, the user can specify the allowed computation time based on the time sensitivity of the task and obtain good solutions to libraries with a few hundred molecules.

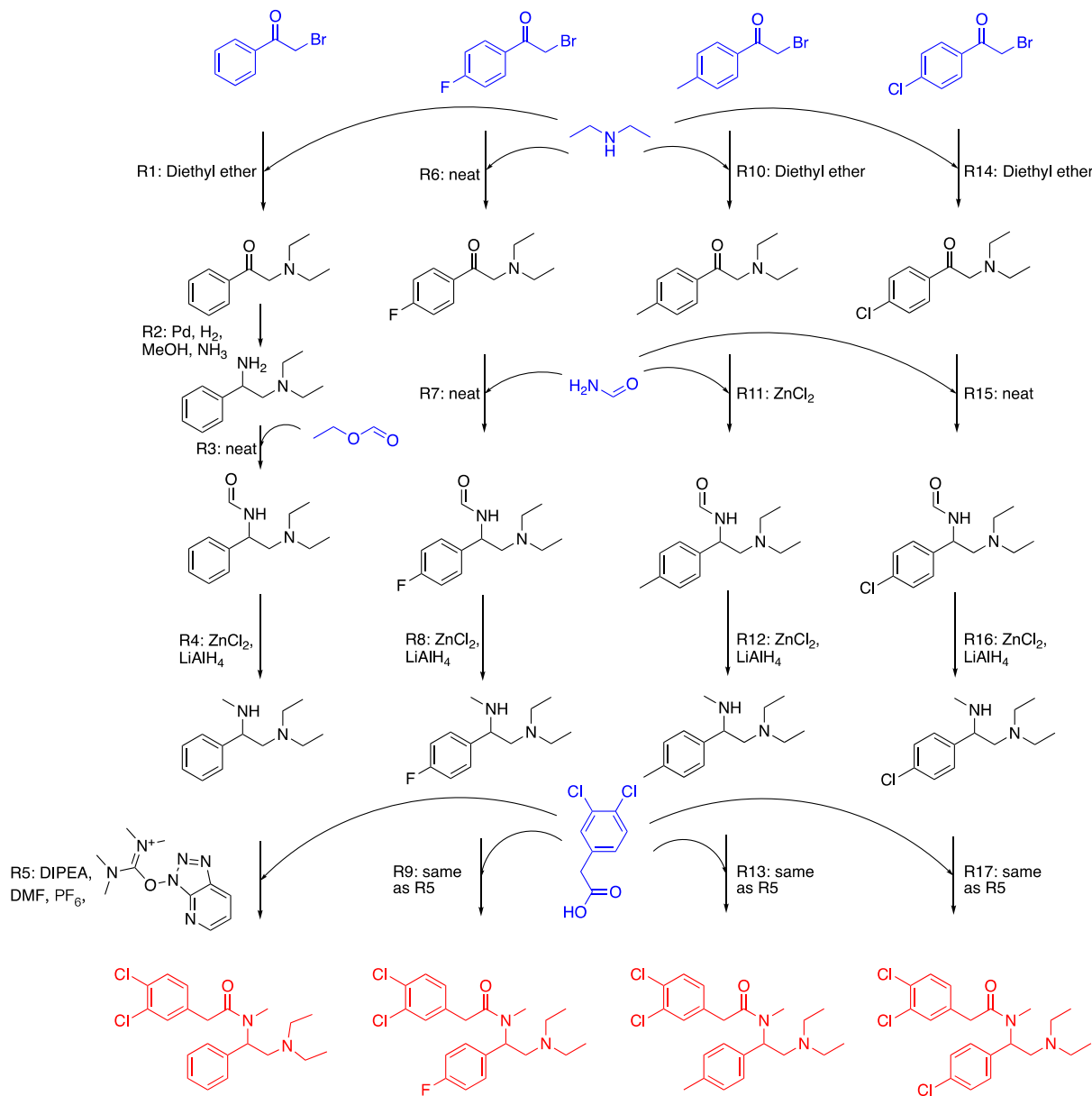
In all, the proposed method should be able to solve molecule library synthesis planning tasks with a few hundred molecules within one or two days.

**Limitations.** We are aware that a subset of suggestions given by the retrosynthesis model might not be chemically feasible because of different aspects of model limitations. In

a) analogs of ICI-199441 used in Case 4



b) example pathways for selected targets that have significant overlap in the combined synthesis planning



**Figure 6.** (a) Analogues of ICI-199441 used in case 4. (b) Examples of targets that have significant overlap in the combined synthesis planning. Molecules highlighted in blue are starting materials; those in red are target molecules. Text next to the arrows takes the format “reaction index: reaction conditions”. Full pathways are shown in Figure S5 in the [Supporting Information](#).

some reactions, necessary reaction conditions are missing (typically solvent, *e.g.*, Figure 3, R7 and Figure 6b, R6). This omission can be relatively easily spotted and corrected by chemists. In addition, missing solvents are likely to be frequently used chemicals, which are often omitted in data records. Therefore, adding these solvents does not contribute much to the overall number of chemicals. Nevertheless, it is ultimately desirable to improve on this aspect by the

community creating better data sets for training the reaction condition recommendation model.

The current reaction prediction model treats molecules as a 2-D graphs which neglect chirality information, so the chiral transformations proposed by the retrosynthesis analysis cannot be adequately evaluated. Therefore, we do not account for stereochemistry in this work. With the advances in reaction prediction capabilities, it would be desirable to take stereo-

chemistry into account, through both better data curation and method development. First, it is important to create high quality data sets that include mainly chiral reactions along with their achiral counterparts (presumably under different reaction conditions) to be able to train and compare models on predicting stereochemistry. For the methods, chiral molecular fingerprints and text-based methods that are directly trained on SMILES strings of the molecules could potentially capture chiral information implicitly, but the effectiveness has not been validated on a large scale because of the lack of data availability. Graph-based representations can also be extended to include chiral information through asymmetric message passing for atoms with different chirality or using calculated atom descriptors based 3-D structures of the molecules.

As a general trend, as more weight is put on minimizing the number of starting materials, more low-score reactions are included in the final pathway. We provide notes in the [Supporting Information](#) for these low-score reactions on whether it reflects some limitations of the retrosynthesis model or it might be a valid reaction based on similar literature precedence. Meanwhile, it is observed that many of such low-score reaction exist across combined synthesis planning and separate synthesis planning, indicating that these low-score reactions are not a consequence of optimizing pathway selection. Therefore, with the improvement of retrosynthesis analysis, the optimization framework developed here will likely provide a similar level of benefit in resource minimization to realize a library of molecules.

## CONCLUSIONS

We combined retrosynthesis analysis and a mixed-integer optimization algorithm to plan syntheses of multiple molecules in a molecule library. We considered multiple objectives, including the number of starting materials, the number of catalysts/solvents/reagents, and the likelihood of success of the overall syntheses plan to select an optimal reaction network to access the target molecules. Instead of preenumerating a fixed number of pathways, we directly formulated the optimization problem on the reaction network which avoided information loss through incomplete tree enumeration. For each reaction, we allowed for selection among 10 different sets of reaction conditions, greatly enhancing the flexibility of selecting reactions with similar conditions. The framework was demonstrated on four case studies, with the size of the library ranging from 2 to 48 targets. Solving the optimization problem effectively reduced the number of starting materials and catalysts/solvents/reagents by promoting the sharing of chemicals between the syntheses of different targets. Benders decomposition was used to accelerate the optimization and improved computational performance on large-size problems. Overall, this framework can serve as a general tool for planning efficient syntheses for molecular libraries, which can simplify chemical inventory and supply chain management, facilitate reaction development, and thus reduce cost of discovery and development.

## ASSOCIATED CONTENT

### Supporting Information

The Supporting Information is available free of charge at <https://pubs.acs.org/doi/10.1021/acs.jcim.0c01032>.

Comparison of computational time, full pathways for Case 1, full pathways for Case 2, full pathways for Case 3, and full pathways for Case 4 ([PDF](#))

SMILES strings for molecules used in the case studies ([XLSX](#))

## AUTHOR INFORMATION

### Corresponding Author

Klavs F. Jensen – Department of Chemical Engineering, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, United States; Email: [kfjensen@mit.edu](mailto:kfjensen@mit.edu)

### Authors

Hanyu Gao – Department of Chemical Engineering, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, United States; [orcid.org/0000-0002-6346-0739](https://orcid.org/0000-0002-6346-0739)

Jean Pauphilet – London Business School, London NW1 4SA, U.K.; [orcid.org/0000-0001-6352-0984](https://orcid.org/0000-0001-6352-0984)

Thomas J. Struble – Department of Chemical Engineering, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, United States; [orcid.org/0000-0003-1695-2367](https://orcid.org/0000-0003-1695-2367)

Connor W. Coley – Department of Chemical Engineering, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, United States; [orcid.org/0000-0002-8271-8723](https://orcid.org/0000-0002-8271-8723)

Complete contact information is available at:

<https://pubs.acs.org/10.1021/acs.jcim.0c01032>

### Author Contributions

All authors have given approval to the final version of the manuscript.

### Funding

This work was supported by the Machine Learning for Pharmaceutical Discovery and Synthesis Consortium and the DARPA Make-It program under contract ARO W911NF-16-2-0023.

### Notes

The authors declare no competing financial interest.

## ACKNOWLEDGMENTS

We thank Elsevier for the access to Reaxys for developing the retrosynthesis platform. We also thank Gurobi Optimization, LLC for the free academic license to use the Gurobi solver.

## ABBREVIATIONS

DIPEA, *N,N*-diisopropylethylamine; NMP, *N*-methyl-2-pyrrolidone; TFA, trifluoroacetic acid; DCM, dichloromethane; DMA, dimethylacetamide; THF, tetrahydrofuran; dba, dibenzylideneacetone; Lawesson's reagent, 2,4-bis(4-methoxyphenyl)-1,3,2,4-dithiadiphosphetane-2,4-disulfide; DMF, dimethylformamide

## REFERENCES

- (1) Kleiner, R. E.; Dumelin, C. E.; Liu, D. R. Small-Molecule Discovery from DNA-Encoded Chemical Libraries. *Chem. Soc. Rev.* **2011**, *40*, 5707–5717.
- (2) Clark, M. A.; Acharya, R. A.; Arico-Muendel, C. C.; Belyanskaya, S. L.; Benjamin, D. R.; Carlson, N. R.; Centrella, P. A.; Chiu, C. H.; Creaser, S. P.; Cuzzo, J. W.; Davie, C. P.; Ding, Y.; Franklin, G. J.

- Franzen, K. D.; Gefter, M. L.; Hale, S. P.; Hansen, N. J. V.; Israel, D. I.; Jiang, J.; Kavarana, M. J.; Kelley, M. S.; Kollmann, C. S.; Li, F.; Lind, K.; Mataruse, S.; Medeiros, P. F.; Messer, J. A.; Myers, P.; O'Keefe, H.; Oliff, M. C.; Rise, C. E.; Satz, A. L.; Skinner, S. R.; Svendsen, J. L.; Tang, L.; van Vloten, K.; Wagner, R. W.; Yao, G.; Zhao, B.; Morgan, B. A. Design, Synthesis and Selection of DNA-Encoded Small-Molecule Libraries. *Nat. Chem. Biol.* **2009**, *5*, 647–654.
- (3) Dandapani, S.; Rosse, G.; Southall, N.; Salvino, J. M.; Thomas, C. J. Selecting, Acquiring, and Using Small Molecule Libraries for High-throughput Screening. *Curr. Protoc. Chem. Biol.* **2012**, *4*, 177–191.
- (4) Röhrig, U. F.; Majjigapu, S. R.; Vogel, P.; Zoete, V.; Michielin, O. Challenges in the Discovery of Indoleamine 2,3-Dioxygenase 1 (IDO1) Inhibitors. *J. Med. Chem.* **2015**, *58*, 9421–9437.
- (5) Xin, B.-T.; Huber, E. M.; De Bruin, G.; Heinemeyer, W.; Maurits, E.; Espinal, C.; Du, Y.; Janssens, M.; Weyburne, E. S.; Kisselev, A. F.; Florea, B. I.; Driessen, C.; Van Der Marel, G. A.; Groll, M.; Overkleeft, H. S. Structure-Based Design of Inhibitors Selective for Human Proteasome B2c or B2i Subunits. *J. Med. Chem.* **2019**, *62*, 1626–1642.
- (6) Kannan Sivaraman, K.; Paiardini, A.; Sieńczyk, M.; Ruggeri, C.; Oellig, C. A.; Dalton, J. P.; Scammells, P. J.; Drag, M.; McGowan, S. Synthesis and Structure-Activity Relationships of Phosphonic Arginine Mimetics as Inhibitors of the M1 and M17 Aminopeptidases from *Plasmodium falciparum*. *J. Med. Chem.* **2013**, *56*, 5213–5217.
- (7) Fleeman, R.; Lavoie, T. M.; Santos, R. G.; Morales, A.; Nefzi, A.; Welmaker, G. S.; Medina-Franco, J. L.; Giulianotti, M. A.; Houghten, R. A.; Shaw, L. N. Combinatorial Libraries as a Tool for the Discovery of Novel, Broad-Spectrum Antibacterial Agents Targeting the ESKAPE Pathogens. *J. Med. Chem.* **2015**, *58*, 3340–3355.
- (8) Keseru, G. M.; Erlanson, D. A.; Ferenczy, G. G.; Hann, M. M.; Murray, C. W.; Pickett, S. D. Design Principles for Fragment Libraries: Maximizing the Value of Learnings from Pharma Fragment-Based Drug Discovery (FBDD) Programs for Use in Academia. *J. Med. Chem.* **2016**, *59*, 8189–8206.
- (9) Gómez-Bombarelli, R.; Wei, J. N.; Duvenaud, D.; Hernández-Lobato, J. M.; Sánchez-Lengeling, B.; Sheberla, D.; Aguilera-Iparraguirre, J.; Hirzel, T. D.; Adams, R. P.; Aspuru-Guzik, A. Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules. *ACS Cent. Sci.* **2018**, *4*, 268–276.
- (10) Blaschke, T.; Olivecrona, M.; Engkvist, O.; Bajorath, J.; Chen, H. Application of Generative Autoencoder in de Novo Molecular Design. *Mol. Inf.* **2018**, *37*, 1700123.
- (11) Segler, M. H. S.; Kogej, T.; Tyrchan, C.; Waller, M. P. Generating Focused Molecule Libraries for Drug Discovery with Recurrent Neural Networks. *ACS Cent. Sci.* **2018**, *4*, 120–131.
- (12) Gupta, A.; Müller, A. T.; Huisman, B. J. H.; Fuchs, J. A.; Schneider, P.; Schneider, G. Generative Recurrent Networks for de Novo Drug Design. *Mol. Inf.* **2018**, *37*, 1700111.
- (13) Sattarov, B.; Baskin, I. I.; Horvath, D.; Marcou, G.; Bjerrum, E. J.; Varnek, A. De Novo Molecular Design by Combining Deep Autoencoder Recurrent Neural Networks with Generative Topographic Mapping. *J. Chem. Inf. Model.* **2019**, *59*, 1182–1196.
- (14) Gao, W.; Coley, C. W. The Synthesizability of Molecules Proposed by Generative Models. *J. Chem. Inf. Model.* **2020**, DOI: 10.1021/acs.jcim.0c00174.
- (15) Segler, M. H. S.; Preuss, M.; Waller, M. P. Planning Chemical Syntheses with Deep Neural Networks and Symbolic AI. *Nature* **2018**, *555*, 604.
- (16) 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* **2019**, *365*, No. eaax1566.
- (17) Szymkuć, S.; Gajewska, E. P.; Klucznik, T.; Molga, K.; Dittwald, P.; Startek, M.; Bajczyk, M.; Grzybowski, B. A. Computer-Assisted Synthetic Planning: The End of the Beginning. *Angew. Chem., Int. Ed.* **2016**, *55*, S904–S937.
- (18) Schwaller, P.; Petraglia, R.; Zullo, V.; Nair, V. H.; Haeuselmann, R. A.; Pisoni, R.; Bekas, C.; Iuliano, A.; Laino, T. Predicting Retrosynthetic Pathways Using a Combined Linguistic Model and Hyper-Graph Exploration Strategy. **2019**, arXiv Prepr. arXiv1910.08036.
- (19) Segler, M. H. S.; Waller, M. P. Neural-Symbolic Machine Learning for Retrosynthesis and Reaction Prediction. *Chem.—Eur. J.* **2017**, *23*, S966–S971.
- (20) 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 Cent. Sci.* **2018**, *4*, 1465–1476.
- (21) Walker, E.; Kammeraad, J.; Goetz, J.; Robo, M. T.; Tewari, A.; Zimmerman, P. M. Learning To Predict Reaction Conditions: Relationships between Solvent, Molecular Structure, and Catalyst. *J. Chem. Inf. Model.* **2019**, *59*, 3645–3654.
- (22) 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. *Chem. Sci.* **2019**, *10*, 370–377.
- (23) Schwaller, P.; Gaudin, T.; Lányi, D.; Bekas, C.; Laino, T. “Found in Translation”: Predicting Outcomes of Complex Organic Chemistry Reactions Using Neural Sequence-to-Sequence Models. *Chem. Sci.* **2018**, *9*, 6091–6098.
- (24) Molga, K.; Dittwald, P.; Grzybowski, B. A. Computational Design of Syntheses Leading to Compound Libraries or Isotopically Labelled Targets. *Chem. Sci.* **2019**, *10*, 9219–9232.
- (25) Gao, H.; Coley, C. W.; Struble, T. J.; Li, L.; Qian, Y.; Green, W. H.; Jensen, K. F. Combining Retrosynthesis and Mixed-Integer Optimization for Minimizing the Chemical Inventory Needed to Realize a WHO Essential Medicines List. *React. Chem. Eng.* **2020**, *5*, 367–376.
- (26) Schrijver, A. *Combinatorial Optimization: Polyhedra and Efficiency*; Springer Science & Business Media, 2003; Vol. 24.
- (27) Gurobi Optimization, LLC. *Gurobi Optimizer Reference Manual*, 2020.
- (28) Jablonský, J. Benchmarks for Current Linear and Mixed Integer Optimization Solvers. *Acta Univ. Agric. Silv. Mendelianae Brun.* **2015**, *63*, 1923–1928.
- (29) Benders, J. F. Partitioning Procedures for Solving Mixed-Variables Programming Problems. *Numer. Math.* **1962**, *4*, 238–252.
- (30) Rahmaniani, R.; Crainic, T. G.; Gendreau, M.; Rei, W. The Benders Decomposition Algorithm: A Literature Review. *Eur. J. Oper. Res.* **2017**, *259*, 801–817.
- (31) Moon, Y.-C.; Baiazitov, R.; Du, W.; Lee, C.-S.; Hwang, S.; Almstead, N. G. Chemoselective Reactions of 4, 6-Dichloro-2-(Methylsulfonyl) Pyrimidine and Related Electrophiles with Amines. *Synthesis* **2013**, *45*, 1764–1784.