

Efficient Enumeration of Branched Novel Biochemical Pathways Using a Probabilistic Technique

Zhiqing Xu^{*,†} and Radhakrishnan Mahadevan^{*,†,‡,¶}

**Department of Chemical Engineering and Applied Chemistry, University of Toronto,
Toronto, Ontario, M5S3E5, Canada*

*†Institute of Biomedical Engineering, University of Toronto, Toronto, Ontario, M5S 3G9,
Canada*

¶Current address: 200 College Street, M5S3E5, Toronto, ON, Canada

E-mail: zhiqing.xu@mail.utoronto.ca; krishna.mahadevan@utoronto.ca

Phone: (1)647-391-5399; (1)416-946-0996

Abstract

The advancement in the field of synthetic biology has allowed metabolic engineers to construct *de novo* pathways in engineered cells. The promiscuous activities of many enzymes which diversify the natural metabolic pathway network bring lots of possibilities for producing desired non-natural compounds. Although various graph-based computational tools were developed to predict novel pathways and have largely expanded the range of chemicals that can be produced biologically, the chemical repertoire reachable in microorganisms can be further broadened. The main challenge remains in dealing with the combinatorial growth of nodes and branches in the putative reaction network generated when searching the branched pathways. In this paper, a branched novel pathway computation tool that leverages probabilistic methods and machine learning,

Anneal Path, is presented. It targets chemical diversity and explores a larger theoretically possible chemical space. The pathway network has been modeled through a loop-less directed hypergraph, where a list of starting compounds and a target compound are connected with a series of hyper-edges, each representing a multi-molecular reaction.

This pathway identification program efficiently enumerates chemically possible branched novel biochemical pathways on the basis of (1) multi-molecular reaction hyper-edge rules which identify enzyme's promiscuity on multiple substrates and/or products, (2) simulated annealing-based pruning algorithm that efficiently handle the huge hypergraphic reaction network generated by hyper-edge rules and (3) a decision-tree-based pathway length predictor which accurately estimates the minimum number of enzymatic reaction steps required between two chemical structures. We also show that this *Anneal Path* can identify a large set of branched pathways more efficiently than purely similarity based search methods. Source code for producing Anneal Path and pathway length predictor are available on GitHub at: <https://github.com/LMSE/Anneal-Path> and <https://github.com/LMSE/PathwayLenPred>, respectively.

Introduction

Engineering microbial production of many valuable non-natural chemicals is becoming prevalent in industry.^{1,2} The key components for the production of non-natural chemicals are (1) prediction of new pathways and (2) enzyme discovery for the novel pathways.³ Numerous successes in metabolic engineering have motivated the development of computational tools to predict promising novel pathways for biosynthesis of target compounds.^{4,5,6} *In-silico* pathway prediction has recently become the preliminary step in the development of novel biosynthetic pathways due to extreme difficulty in non-computational pathway design approaches.⁷ A generalized workflow for *de novo* pathway design *in-silico* mainly involves the following two critical steps: (1) the identification of all potential reaction routes from a list

of starting materials to the desired target compound, on the basis of structural transformation from substrates to products, where each neighbouring compound is connected by promiscuous enzymatic rules^{8,9}(2) the pathway prioritization based on various criteria, including thermodynamic feasibility, number of non-natural reaction steps, enzyme docking and cellular resources cost etc.^{10,11,12} Different scoring systems have been designed to rank the pathways predicted from the identification step; most of them use a combination of the filtering or scoring attributes and make good evaluations, for example Yim and Haselbeck (2011) has accurately selected a novel 1, 4-butanediol pathway from over ten thousand 4-to-6-steps pathways (identified by the SimPheny Biopathway Predictor) to be tested in vitro; the ranking is based on theoretical yield, pathway length and thermodynamic feasibility, etc.¹³ Recent research has also provided great solution to retrosynthetic pathway step suggestion based on training a transformer model.¹⁴ Other pathway finding tools (i.e., BNICE,¹⁰ ATLAS,¹⁵ GEM-path,¹⁶ ReactPred¹⁷) have used different molecular encoding systems, reaction simulation algorithms and different in-house transformation rules. In addition, algorithms that automatically parse biochemical reactions in database and generate reaction rules have been developed in RetroRules,^{18,19,20} which allows the program to predict hypothetical reactions based on entire biochemical reaction databases.

The growing computational resources available for in-silico pathway design and improved algorithms have allowed researchers to analyse a greater number of potentially novel pathways.²¹ However, current methods for generating novel pathways with non-natural intermediates either use simple directed edges to represent reactions (simulating only transformations from one compound to another)^{13,22} or identify multiple reactants and products of each reaction but allow only one substrate non-native to the enzyme in each promiscuous enzymatic reaction.¹⁹ There are also algorithms that recover branched pathways through merging linear pathways (such as LPAT²⁵). These existing methods avoid to deal with promiscuity on multiple substrates in reaction while expanding branched pathway networks in order to prevent the putative reaction network from being too large. And this is where a large number

of reasonable reactions are left out and causes loss in branched novel pathways. On the other hand, stoichiometry-based methods, for example, OptStoic²³ and NovoStoic²⁴ use the incidence matrix of the pathway network (namely the stoichiometric matrix) to approach conversions involving multiple substrates (or products). However, neither of these tools finds novel compounds and they are not used to generate novel pathways with non-natural chemicals. Stoichiometry-based method itself has limitations that it performs reaction prediction and path searching in two separate stages which is equivalent to identifying pathways in an extended reaction database. This approach will result in the loss of a considerable number of branched novel pathways as compared with reaction rule-based approach which expands the reaction network by iteratively applying transformation rules to the compounds.

In order to identify promiscuity on multiple substrates of multi-molecular reactions as expanding branched novel pathways, we use hyper-edge rules for simulating multiple substrates in enzymatic reactions. Hyper-edge rules are transformation rules that contains multiple substructures for retaining all possibilities of non-natural reactions. See Figure S1(A) for a comparison between four types of edges in a hypergraph correspond to four types of reactions: (1) $A \rightarrow B$ simple reaction, (2) condensation reaction, (3) decomposition reaction and (4) a multi-molecular reaction, which contains multiple reactants and products. Reactions identified by hyper-edge rules can later be used to reconstruct a hyper-graphical pathway network which includes all structurally-possible intermediate compounds. To our knowledge, there is no current computational method that uses hyper-edge reaction rules (or multi-molecular reaction rules) to simulate reactions with multiple novel (or non-natural) substrates and construct a hypergraphic pathway network that takes into account all chemically plausible novel reactions. While having the advantage of enumerating all novel biochemical routes, hypergraph-based reaction network construction would require a huge computational cost due to the explosion of network size. Many different approaches have been used in current pathway prediction tools to remove unpromising nodes (chemicals) when building the prediction network. Apart from the maximum pathway length and

upper bound on the size of molecules, reaction feasibility scores are sometimes assessed to screen potential substrates and control the network complexity.^{10,19} These approaches have targeted the chemical feasibility of constructing pathways *in vivo*, which is more often considered when ranking the result pathways.¹¹ On the other hand, screening compounds through estimating the chemical similarity to the target compound aims to select intermediates that can be easily transformed (within a few steps of hypothetical biochemical reactions) to the desired target.^{26,27} SimIndex uses an openbabel built-in “FP4” similarity metric which helps to explore the hypothetical chemical space to the fullest extent possible while reducing the computation time.^{26,28} Such relatively more systematic similarity-based pruning method was used widely in various *de novo* pathway identification algorithms.^{29,30} Apart from “FP4”, a variety of chemical fingerprints that captures distinct molecular structure information, for example, Molecular ACCess System (MACCS), Extended Connectivity Fingerprints (ECFPs) and atom-pair descriptor^{31,32,33} as well as other similarity measuring methods, for example, Maximum Common Substructure (MCS)³⁴ have been developed for compounds virtual screening with different purposes and they have similar performances when being used for pathway prediction.

In this work, we present *Anneal Path*, an algorithm for enumerating branched novel pathways based on a probabilistic method. We have improved the current similarity-based pruning method (i.e., SimIndex²⁶) through applying simulated annealing, a probabilistic optimization algorithm for solving global extremes in combinatorial problems. We rationalize that the use of such a method might be valuable for efficiently identifying relatively longer (with 8-10 reaction steps or even more) branched pathways, which involves a combinatorial explosion in the number of reactions and pathways. We formulate the pathway prediction as an optimization problem and use this algorithm to identify pathways to a broad range of targets. We have also validated our algorithm by using it to predict a set of 20 biologically relevant pathways demonstrating the biological significance of our approach. On top of the chemical structural similarity estimation for pruning the reaction network, we have developed

a decision-tree-based distance measurement system to determine the number of reaction steps required to transform one chemical structure to another, which allows better estimation than the structural similarity scores. We have used this as a pre-screening in *Anneal Path* to set the maximum depth of the reaction network.

Approach

Overview

Here, our objective is to develop a workflow for using hypergraph based reaction rules for enhancing the diversity of the predicted metabolic pathways. Our workflow involves the use of a chemistry framework to represent metabolites, hypergraph based reaction rule representation to generate the combinatorial expansion, simulated annealing based search algorithm and pruning to search the space for a pathway from a starting set of compounds to a target compound, and machine learning based reaction distance measurement to limit the number of steps in the combinatorial expansion.

Chemistry Framework

Anneal Path is scripted in python. The reaction simulator is based on *RDKit* which incorporates SMILES/SMARTS encoding of compounds and transformations with efficient built-in reaction modeling algorithm.^{35,36} Reactions can be predicted according to both user-defined in-house sets of transformation rules and different packages of rules (written in SMARTS) available online (i.e., RetroRules¹⁸). Since cofactors (e.g., NADH or ATP) or small molecules like water can be ignored when expanding the reaction network, they are not encoded in the transformation rules for reaction prediction/simulation but go with the rule ID's and can be eventually shown in the final results. Different structural similarity metrics (those calculated based on atom-pair,³³ topological torsion,³⁷ ECFP,³² MACCS,²⁸ etc.), have also been incorporated based on corresponding fingerprints programmed in *RDKit*.³⁶ Those different

similarity metrics are tested for their performance in compounds screening and are then used to benchmark the new algorithm. While a set of generalized in-house reaction rules have been used for the work in this paper, *Anneal Path* is also compatible with automatically generated RetroRules or other user-customized reaction rules on certain conditions.^{18,19}

Hypergraph-Based Network Representation

Anneal Path adopts a directed hypergraph for representing the putative reaction network which enumerates all the possible pathways. The use of hyper-edge reaction rules allow multiple non-native compounds in the predicted reactions and therefore produces a larger and more complex reaction network, as well as a combinatorial growth of nodes and branches (as compared with exponential growth in a simple graph representation).

A directed hypergraph can be written as a collection $H = \{C, R\}$, where C is a set of compounds represented by nodes or vertices of the graph and R is a set of all reactions represented by directed hyper-edges. Each reaction hyper-edge is an ordered pair $R_n = \{S_n, P_n\}$; S and P are two subsets of reactant(s) and product(s), respectively.³⁸

Hyper-edge reaction rules model the promiscuity on multiple substrates of reactions and allow multiple non-native compounds in the predicted reactions. Figure 1 shows an example of one such enzymatic reaction and corresponding rule. Figure 2 shows the comparison between simple-graph-based and a hypergraph-based pathway representation, where greyed-out nodes and edges are omitted in a simple-graph-based search. Such branched pathways can only be recovered in a hyper-graphic reaction network.

A comparison of exponential complexity of a simple-graph-based network and combinatorial complexity of a hypergraphic reaction network is discussed in detail in supporting information. The significant difference in the number of nodes and pathways between the two networks shown in Figure S1 is due to the fact that in each iteration step of hypergraph-based network expansion, all combinations of compounds that contain substructures that followed the reaction rules would add a new reaction hyper-edge as well as node(s) to the

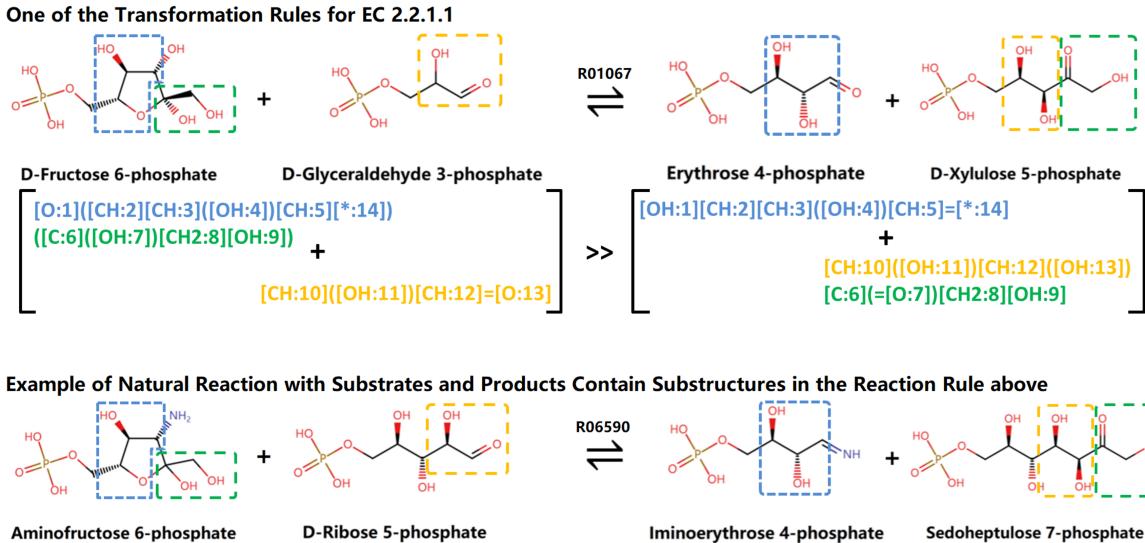


Figure 1: The hyper-edge transformation rule (that models the promiscuity on both substrates) for EC 2.2.1.1 with reaction centers highlighted. The corresponding SMARTS string that encodes the atom mapping of the reaction center are shown below. The natural reaction at the bottom shows example of promiscuity on both substrates which implies the significance of using hyper-edge rules to allow multiple non-native substrates/products in one reaction step.

system. This makes the number of nodes and branches undergo a combinatorial growth instead of an exponential growth we see in a simple-graph search.

As discussed previously, a hyper-graphic reaction network can help to generate all possibilities of biochemical transformation routes to a target compound. Figure 3A shows a small section (around 10%) of such putative reaction network centered around Fructose 6-phosphate (F6P) (in real case example), which is composed of mostly reactions with multiple reactants and/or products. All the branched pathways are highlighted in blue lines while linear pathways are highlighted in red lines. Examples of complete branched novel pathways generated from F6P to DAHP are shown in Figure 3B, where the natural pathway found in glycolysis and tyrosine metabolism are highlighted. Each reaction in the natural pathway is marked by colors corresponding to their edge types. This figure clearly illustrates the combinatorial complexity evident in hyper-graph based networks, typical in metabolite networks.

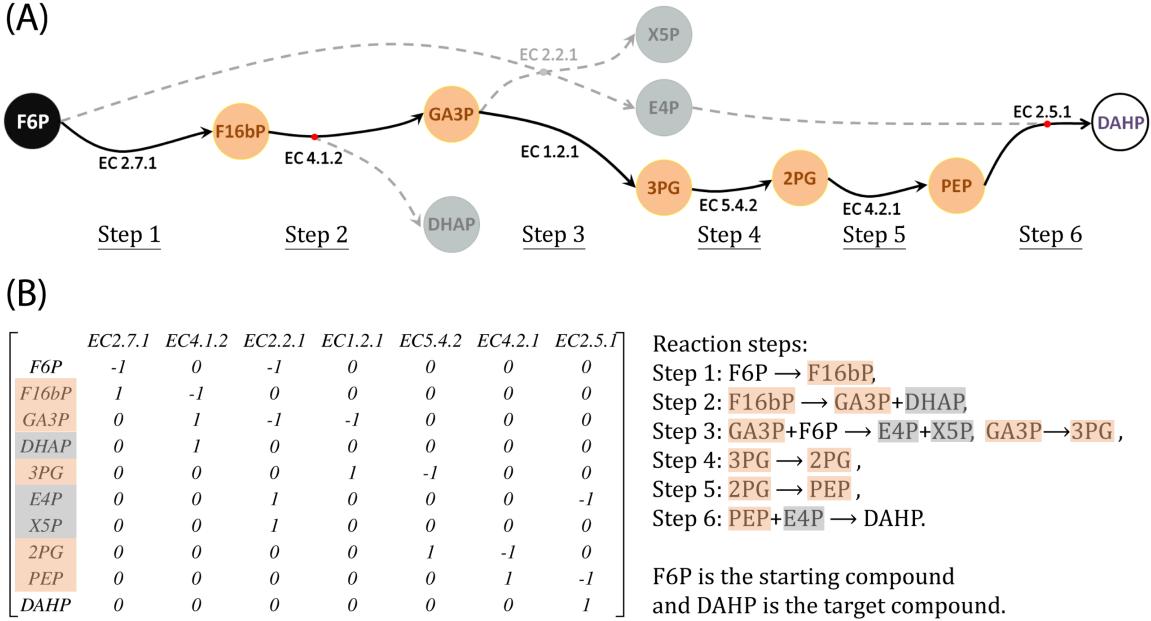


Figure 2: (A) A comparison between simple-graph-based pathway representation and a hypergraph-based pathway representation of a six step pathways from F6P to DAHP found in glycolysis and aromatic amino acid biosynthesis pathways. In a simple-graph representation, the greyed-out nodes and edges shown in the figure are omitted. (B) The matrix shows a corresponding incidence matrix (stoichiometry matrix) of the hypergraph-based pathway representation. The abbreviations for metabolites in the figure were as follows: Fructose 6-phosphate (F6P), Fructose 1,6-bisphosphate (F16bP), Glyceraldehyde 3-phosphate (GA3P), Dihydroxyacetone phosphate (DHAP), D-Erythrose 4-phosphate (E4P), D-Xylulose 5-phosphate (X5P), 3-Phospho-D-glycerate (3PG), 2-Phospho-D-glycerate (2PG), Phosphoenolpyruvate (PEP), 2-Dehydro-3-deoxy-D-arabino-heptonate 7-phosphate (DAHP).

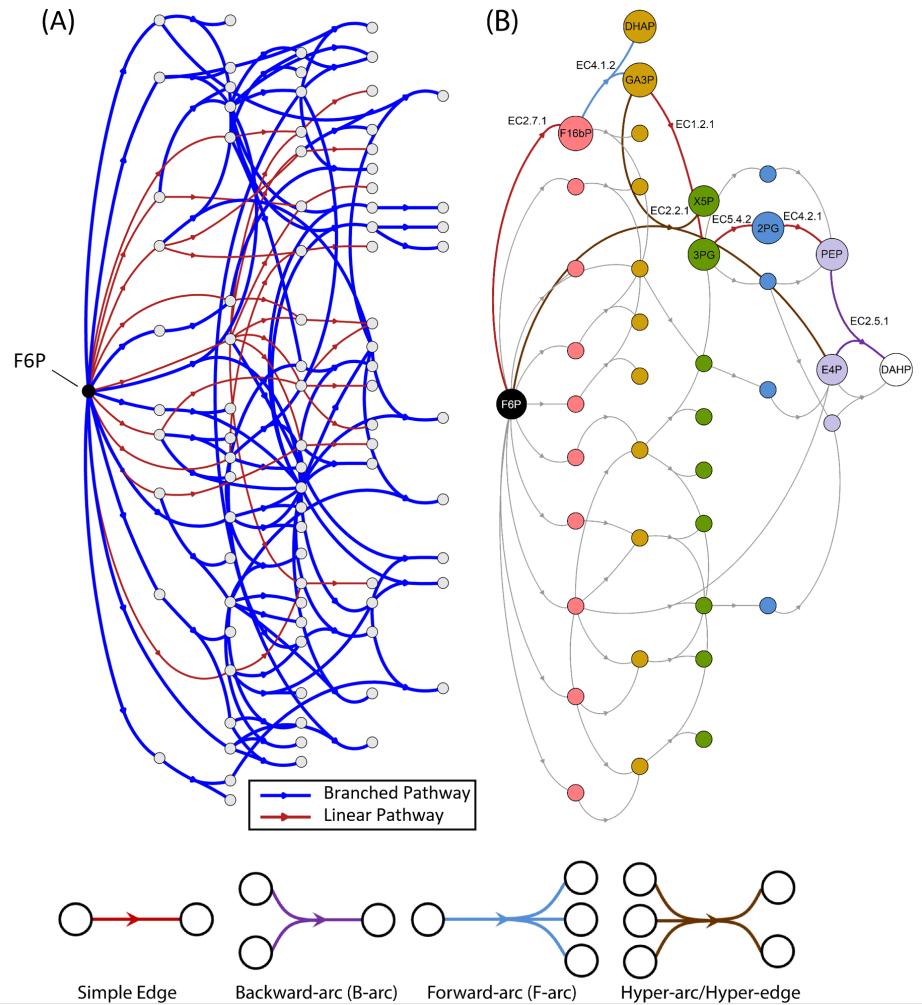


Figure 3: (A) A small portion of hypergraph-based putative reaction network generated centered around starting compound F6P. All the branched pathways are highlighted in blue lines while linear pathways are highlighted in red lines. (B) A small portion of novel branched pathways (from source compound F6P to target compound DAHP, generated by reaction rule based pathway finding algorithms), the natural pathway found in glycolysis and tyrosine metabolism is highlighted. Reactions in the natural pathway are colored according to hyper-arc types.

Search Algorithm

The retrosynthesis process is most commonly used in pathway finding algorithms. It iteratively substitutes the products with substrates based on predicted reactions, until all compounds are replaced by the starting materials. As such, a reaction network centered around the target compound is generated to link the target to the sources. This retrosynthetic technique works well for designing linear pathways or branched synthetic pathways for relatively large molecules.¹⁹ However, in a hypergraph context, the retrosynthesis approach has some limitations. Typically, retrosynthesis applies only simple edge reaction rules ($\mathbf{a} \rightarrow \mathbf{b}$) and B-arc reaction rules ($\mathbf{a} + \mathbf{b} \rightarrow \mathbf{c}$, where \mathbf{a} , \mathbf{b} and \mathbf{c} are substructures). There are two scenarios of B-arc rules used in retrosynthesis. One is the corresponding reaction has only one product. The second is the reaction rule models promiscuity of only one out of multiple products. The advantage of the B-arc rules (i.e., $\mathbf{a} + \mathbf{b} \rightarrow \mathbf{c}$) is as follows: when a known compound containing substructure \mathbf{c} is being backwardly substituted, two unique compounds containing the substructures \mathbf{a} and \mathbf{b} , respectively can be easily identified based on the rules. However, if the hyper-edge rules (i.e., $\mathbf{a} + \mathbf{b} \rightarrow \mathbf{c} + \mathbf{d}$) are used to model promiscuity of multiple products in one reaction (in retrosynthesis), there can be an infinite number of potential compounds that contain the substructure \mathbf{d} (based on the reaction rule, any compounds containing substructure \mathbf{d} would generate a different set of two compound structures containing \mathbf{a} and \mathbf{b} , respectively). As a result, for backward search, using hyper-edge rules to model promiscuity of multiple products in one reaction is not practical.

On the other hand, the hyper-edge rules can be used in the forward search direction as the compounds that can serve as the substrates in these reaction rules are finite. This finite set includes compounds that are given as starting material (e.g. intermediates such as pyruvate) and those that are generated through the application of the reaction rules from the previous steps. Reaction network constructed in this way literally enumerates all chemically plausible possibilities. Hence, there is a key difference between forward-direction reaction expansion, (starting from the source compounds) and reverse-direction retrosynthesis, when

using a hypergraphic reaction network representation.

Anneal Path uses a bidirectional breadth-first search starting from both the sources and the target to generate two reaction networks at the same time. This approach has been used previously in linear pathway identification tool presented in *Pathfinder* by Noor et al.³⁹ The search method used has reduced the number of reaction expansion steps by half and allows a more efficient search for relatively longer pathways. Figure 4 shows an example of such a double direction search. Two hyper-graphic reaction networks are generated centered around the source compound and the target compound, respectively. Nodes generated by one expansion step are lined up horizontally. The starting-compounds-side network is composed of all types of hyper-arcs while the target side one contains only B-arcs.³⁸ Bridge compounds connecting the two hypergraphs are then identified before completed pathways (from reactant to product) are recovered through hypergraph path finding algorithm.

Pruning

As mentioned in previous sections, structural similarity estimation has been used to improve the pathway finding efficiency. In a retrosynthesis context, similarity based compound screening is performed after each reaction expansion step. The structures of all intermediate compounds that have been generated are compared with the source compound(s) specified by the user. The similarity scores are then used for compound selection to further expand the reaction network. This idea is based on the assumption that those compounds with similar chemical structure can be converted from one to another within relatively fewer number of enzymatic reaction steps. In most cases, a global increasing trend can be observed in similarity scores of all intermediates along a pathway, from the starting compound to the target, as shown in Pertusi et al.²⁶

In order to validate the performance of the current similarity metrics in predicting the number of reaction steps required for biochemical transformations between any two compounds, the correlation between the similarity score and reaction distance is examined. Here,

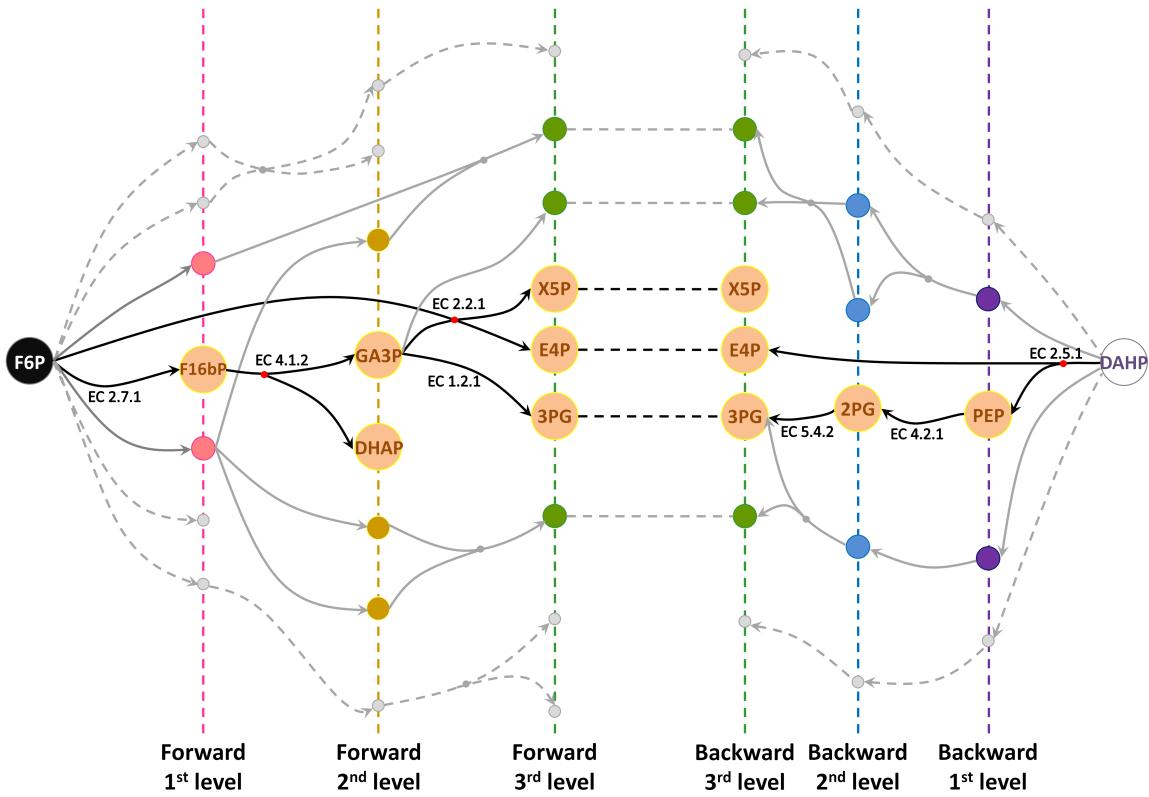


Figure 4: A double direction search for six-step pathways from F6P to DAHP. Two reaction networks are generated centered around source compound F6P and target DAHP, which is shown by black and white nodes in the figure, respectively. After three steps of reaction expansion from both sides, bridge compounds are identified which links the two reaction networks. Compounds found from one reaction expansion step are aligned.

the reaction distance between two compounds is defined to be the minimum number of reaction steps required to transform one compound to the other. Reaction data obtained from MetaNetX reaction database are used to construct reaction trees and pathway networks.⁴⁰ Subsequently, the distance data are generated based on the constructed pathway network. In Figure S2 (see Supporting Information), the distributions of similarity scores estimated based on four different similarity metrics (MCS, Atom-pairs, MACCS and ECFP) are plotted against the real reaction distances.^{29,32,33,34} A negative correlation between the similarity scores and the real reaction distances can be observed from the plot, which confirms the idea of using similarity score to select compounds that are relatively “closer” to the target. However, the R square (coefficient of determination) value is not high (the highest being the one for ECFP, which equals to 0.39). This implies the weak ability to predict the reaction distance with a calculated similarity score, motivating the development of a machine learning based pathway length prediction to manage the combinatorial explosion in the number of predicted pathways.

The efficacy of similarity-based pruning is further reduced in a hypergraphic reaction network. As transformation rules of reactions with multiple substrates and/or products are adopted, the similarity-based pruning method struggles with predicting branched pathways. We attribute this to the fact that for reactions with multiple reactants and/or products, the substrate(s) and product(s) might not have similar structures. For example, PEP and DAHP in the condensation reaction of the sample pathway shown in Figure 2, have a structural similarity as low as 0.47 (compared with 0.69 between X5P and DAHP). Figure S3 (in the Supporting Information) shows the distributions of similarity scores of all pairs of reactant and product in MetaNetX⁴⁰ reaction database grouped by simple reactions ($A \rightarrow B$, red), condensation reactions (purple), decomposition reactions (blue) and multi-reactant-multi-product reactions (brown), with all the co-factors being ignored. We observe completely different distributions between simple reactions and reactions with multiple reactants and/or multiple products. (p-value of Wilcoxon signed rank test $< 10^{-3}$), which suggests the poten-

tial problem of using this pruning method for a reaction network generation using hypergraph based reaction rules. Thus, in a hypergraphic putative reaction network, similarity pruning could possibly eliminate intermediate compounds that are relatively dissimilar to the target, which would lead to significant loss of resulting pathways.

In *Anneal Path*, we improve the pruning method using an optimization approach. The similarity-based compound screening is considered to be an optimization of the similarity score, since the similarity score along any expected resulting pathways always approaches the maximum value (≤ 1) as the starting compounds are transformed towards the target. A temporary drop of the similarity score can be seen as a decrease of the objective function after a local maximum has been reached. Current pruning algorithms, which mostly select compounds based on their score improvements as compared with those of their parent node(s), can be easily trapped by local maxima. In order to find those pathways that do not have a perfectly increasing trend in the similarity from the sources to the target, either a lot more compounds have to be taken into account during putative network construction, or a different compounds selecting schema has to be used. SimIndex,²⁶ for example, has introduced a tolerance factor on the parent node(s) similarity to allow more compounds to be expanded at each reaction steps. While the tolerance factor helps to recover a number of pathways, the increased fluctuations in the similarity scores in a hypergraphic reaction network would require a very low tolerance or that most compounds for each reaction step be conserved for more search results.

Here, we propose a more computationally efficient approach to this problem by using a global optimization algorithm, simulated annealing.⁴¹ Simulated annealing can be viewed as an adaption of the random walk Metropolis-Hastings sampler for optimizing functions. It takes a great number of random "moves" to make changes to the objective function, and accepts/rejects each random move on the basis of a probability, which is calculated using, (1) numerical change to the objective function caused by the "move" and (2) the annealing temperature.⁴² It starts with a search throughout the entire solution space. A decreasing anneal

temperature gradually narrows down the search space while the probabilistic-based schema still allows exploration outside the search space with a relatively lower chance. The broad and “random” search at the beginning and the probabilistic-based acceptance of “moves” that decreases the objective function together help to avoid local extreme in the interest of seeking the global maximum. The algorithm has been proven very useful for solving complicated optimization problems in high dimensions (e.g. traveling salesman problem).

In our case of compound screening, the objective function to be optimized is implicitly defined as the sum of similarity scores of all tail nodes in the temporary putative network after all accepted random "moves" are completed within the current reaction expansion step. Thus, rejecting a reaction corresponds to a child node being identical to its parent node, which brings zero changes to the objective function. On the other hand, accepting one reaction is simply the replacement of one parent node with its corresponding child node(s). The probability of accepting a reaction is therefore,

$$\alpha_n(rxn_x) = \min\left\{\exp\left[\frac{g(\text{accept}) - g(\text{reject})}{T_n}\right], 1\right\} = \min\left\{\exp\left[\frac{\Delta S_{rxn_x}}{T_n}, 1\right]\right\}$$

where,

- $\{T_n\}$ is a sequence of decreasing temperature parameters (namely cooling schedule),
- $g(\text{move})$ is the implicitly defined objective function,
- $g(\text{accept}) - g(\text{reject})$ is estimated by ΔS_{rxn_x} ,
- ΔS_{rxn_x} is the difference in similarity scores of reactant(s) and product(s).

Figure 5 is a flowchart showing the *Anneal Path* search mechanism. Similarities are estimated using multiple metrics (i.e., ECFP, MACCS and atom-pairs) for comparison, since their performances in reflecting reaction distances shown in Figure S2 are equally poor. In addition to the probabilistic-based screening, we force the program to select compounds

from two other groups (1) compounds with top similarity scores and (2) compounds with top similarity score improvements, to allow control of randomness in the system. The size of each bin is left as an input parameter to help adjust the reaction network size.

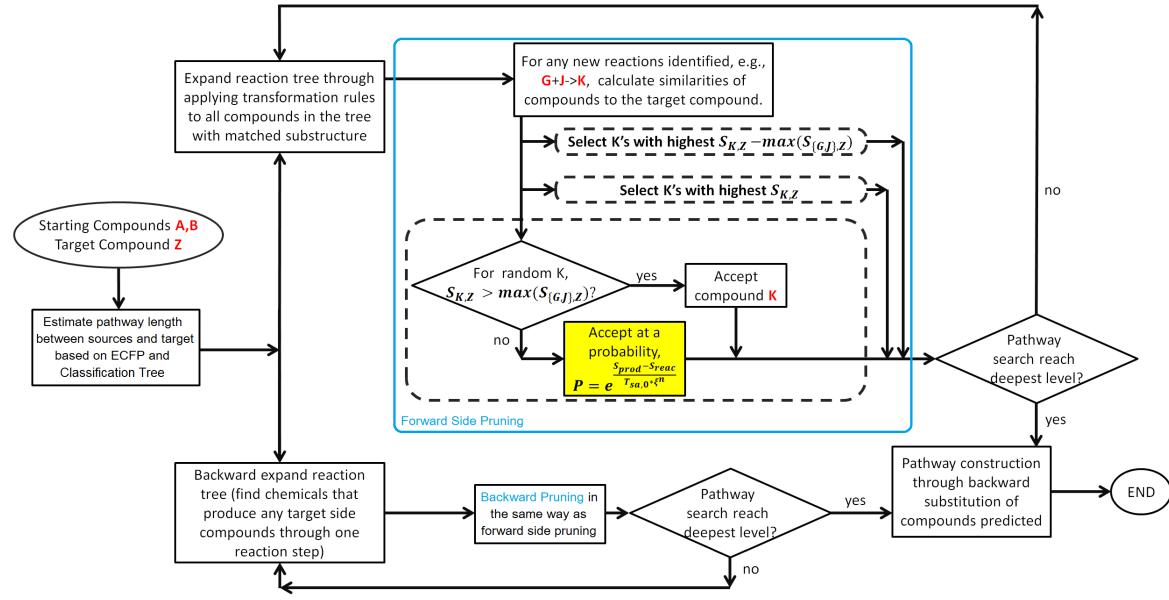


Figure 5: *Anneal Path* algorithm shown in the flowchart diagram. Simulated annealing algorithm has been used to improve the structural similarity based pruning for better computation efficiency, where $T_{sa,0}$ and ξ represent the initial temperature and cooling coefficient (parameters of simulated annealing algorithm), respectively. $S_{K,Z}$ represents the similarity score of compound K and Z , and $\max\{S_{G,J},Z\}$ indicates the highest score among all reactants in this reaction (in this case G and J). $\{T_n\} = \{T_{sa,0} \cdot \xi_n\}$ is a sequence of decreasing temperature parameters called cooling schedule of simulated annealing.

Machine Learning Based Reaction Distance Measurement

As discussed in the *Pruning* section, current similarity metrics have limited ability to estimate the number of reaction steps between pairs of chemical structures. Here, we present a reaction distance prediction model based on decision tree classification which predicts the number of biochemical transformations required (or shortest pathway length) between any two chemical structures. Data preprocessing is as follows. Raw data is obtained from three MetaNetX database files available online, *chem_prop.tsv*, *reac_prop.tsv* and *chem_ref.tsv*. The *chem_prop.tsv* file contains all chemicals and their database id's with physical/chemical

properties information. The *reac_prop.tsv* file contains all the reaction information. And *chem_ref.tsv* is a file that contains a list of different types of text representation and encoding of molecules including the SMILES string, InChI (IUPAC International Chemical Identifier). Reaction networks are then constructed through applying a ten step reaction expansion on each compound in the MetaNetX database.⁴⁰

Generative models trained for chemical structures design have used autoencoder to convert different types of discrete molecular representations to continuous vector representations. The method has been applied to different molecular representations, including chemical fingerprints (i.e., extended connectivity fingerprint, ECFP), SMILES strings etc.⁴³ The paper *Automatic Chemical Design Using a Data – Driven Continuous Representation of Molecules*⁴³ has reported that a variational autoencoder (VAE) has proven useful for converting SMILES strings to continuous vectors and achieves better performance in structure design. The SMILES-based molecular variational autoencoder has been used in this work to convert the chemical structures into numerical data. An already trained VAE is used directly to convert SMILES strings to a 292-dimensional vector space (hereinafter referred to as VAE-encoded data). The dataset follows a multivariate Gaussian distribution. Therefore this continuous representation of pairs of molecules and corresponding distances (X_1 , X_2 and Y) are used as the inputs to distance metric learning models, as shown in Figure S4 in the Supporting Information.⁴⁴

Although chemical encoding using ECFP has been reported to be outperformed by the SMILES for training a VAE,⁴³ it is straightforward to take in molecular substructures directly to train the structural similarity metric. In our work, we have used a different approach to encode the data using ECFP as an alternative to the VAE based approach. All chemicals in SMILES are converted to ECFPs before each distinctive identifier is one-hot encoded. The list of ECFPs can then be converted to a vector with each dimension representing the number of corresponding identifiers contained in the structure (hereinafter referred to as ECFP-encoded data). Figure S5 shows the ECFPs identifiers of the *oxazolidin – 4 – one*

structure and a sample ECFP-encoded data point. Differences in the two ECFP-encoded vectors are taken for each data point, and (ΔX and Y) are used as inputs to the machine learning models (other than distance metric learning models).

A list of classification models (i.e., random forest classifier, decision tree classifier, support vector classifier, ridge classifier, gradient boosting classifier, simple neural networks, large margin nearest neighbor, etc.) have been trained using the processed data. We have used only non-parametric models for the ECFP-encoded data and both parametric and non-parametric models for VAE-encoded data. This is because the ECFP representations are count-encodings and would not satisfy the distribution assumption of most parametric machine learning models, while VAE-encoded data are normally distributed. Both the ECFP-encoded data and the VAE-encoded data were divided into training set, cross-validation set (for tuning the hyperparameters) and test set (for examining the performance of predicting the reaction distance). Hyperparameters are adjusted through ROC-AUC analysis based on one-vs.-rest classification.

Machine Learning models have been further trained on synthetic data generated by Anneal Path. Through applying the list of generalized reaction rules on random starting compounds for a few iterations, a much larger and complicated putative reaction network can be constructed than the one contains the MetaNetX reactions alone. While the pathway data from MetaNetX database is used for validating the new distance metric developed, the prediction performance on the synthetic dataset is also significantly important as it reflects the capability of determining the maximum number of search steps when Anneal Path is used to find novel pathways.

Results and discussion

Improved Network Pruning Using Simulated Annealing

The computational complexity of constructing a hypergraphic putative reaction network is examined based on the increases in the number of reactions generated for each step of reaction expansion. We illustrate the combinatorial complexity of such a reaction network through comparing with a simple putative network generated by simple-edge reaction rules alone (rules that take only one substrate and output one product). Figure 6 shows a comparison between the growth of the number of edges in the two types of forward side reaction network generated centred around pyruvate (brute-force search). Computation of expanding this branched pathway network was carried out on AMD Opteron(tm) 8431 Processor and took around 120 hours. The huge difference in the growth of nodes and edges in the network illustrate the combinatorial complexity and the significantly enhanced diversity of the hypergraphic reaction network.

The performance of *Anneal Path* in enumerating branched pathways is tested with identifying novel amino acids biosynthesis/degradation pathways. Sources and targets are compound pairs found in KEGG/ATLAS maps connected by different pathways with 6-12 reaction steps.^{15,47} We demonstrate the improved computational efficiency and show the prediction of a larger number of pathways through comparing this algorithm with the similarity-based compound screening schema. Two pruning algorithms are manipulated such that the same number of intermediates is selected for further search in each reaction expansion step. As mentioned in the *Pruning* section, the performance of the current similarity metrics have been tested and compared. This reflect their abilities to select compounds that are more likely to be precursors of the target based on number of branches in the reaction network. However, the best-performing similarity metrics are quite unpredictable for different cases and it is not clear which one is most suitable for similarity-based pruning. While ECFP³² helps to find relatively more pathways in most cases, it is not favorable for identifying some

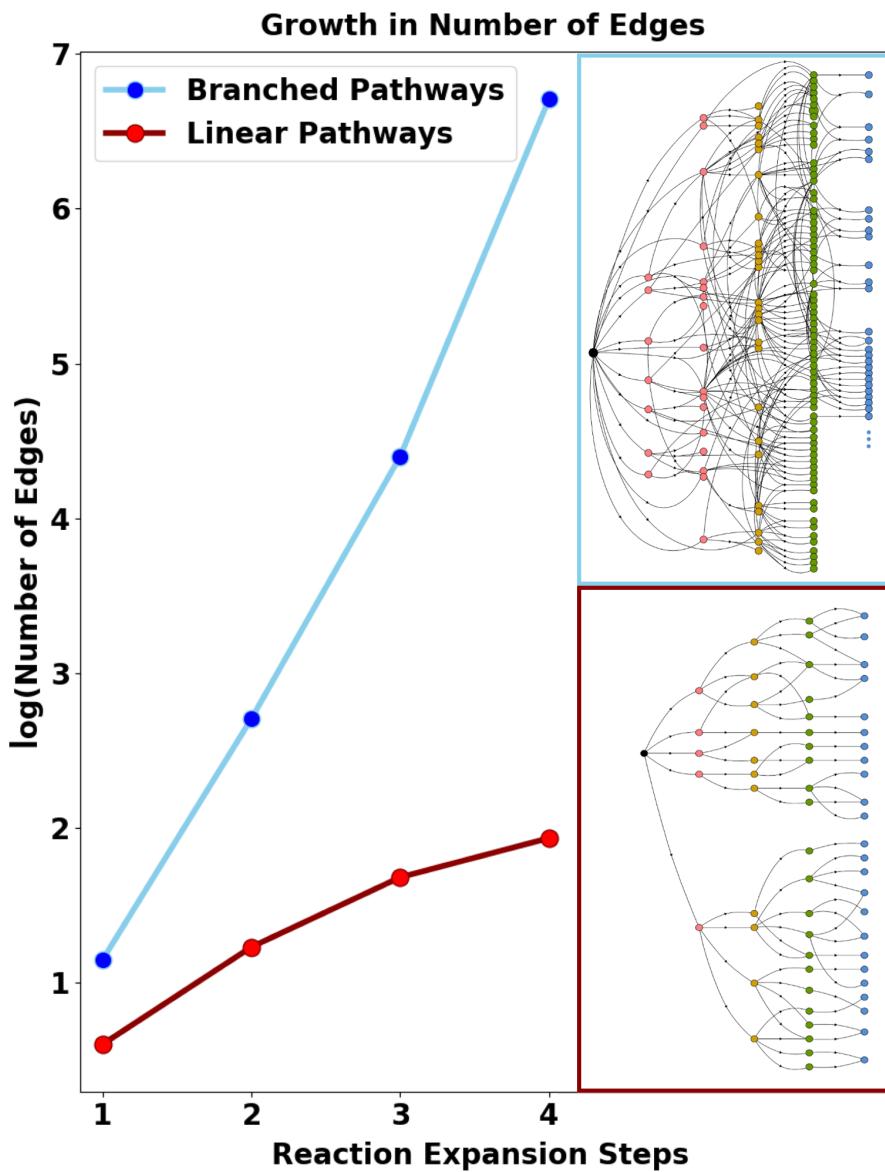


Figure 6: A comparison between the number of edges searched at each level for reaction networks generated centered around pyruvate on the basis of (a) hyper-edge transformation rules (upper curve highlighted in blue) and (b) simple-edge reaction rules (lower curve highlighted in red). On the right of each curve, a small portion of corresponding reaction network are visualized where compounds searched at each level are aligned.

pathways.

Based on tests and comparisons of different similarity metrics, we have chosen ECFP to be used in both similarity-based pruning and *Anneal Path* pruning. Figure 7A shows a comparison between the numbers of pathways generated based on similarity pruning (shown by red markers) and *Anneal Path* pruning (100 trials for each pathway prediction task) which is shown by the box plots. Among all the pathway prediction trials we have performed, the count of Anneal Path predicting more pathways than similarity-based pruning approach was on average, 71.5%. Figure 7B shows the relative efficiency of pathway enumeration through the use of two different pruning methods (red line here is the best-performed similarity metric).

In the benchmarking tests, we have used an implicitly defined cost function and a fixed set of primitive parameters (cooling schedule) as the number of compounds to be checked is unpredictable. Without much tuning of the parameters in simulated annealing algorithm, the results have a wide interval in some cases (see Figure 7).^{41,42} Nevertheless, the new pruning algorithm still generates an increased number of pathways and has better computation efficiency in most of the tests (see Figure 7).

Machine Learning Prediction of Reaction Distance for Improved Pathway Prediction in *Anneal Path*

Among all the machine learning models that has been trained on VAE-encoded data and ECFP-encoded MetaNetX data, decision tree classifier trained on the ECFP-encoded data has achieved the best performance of predicting the reaction distance of compound pairs in the test set . It has achieved a Spearman' correlation of 0.703 (as compared with an highest correlation of 0.363 between similarity scores and actual reaction distances, with different similarity metrics tested). The accuracy of reaction distance prediction is 60% (as compared with an highest accuracy of 0.26 estimated by different similarity metrics). Table 1 shows performances of the new distance metric model that has been trained (decision

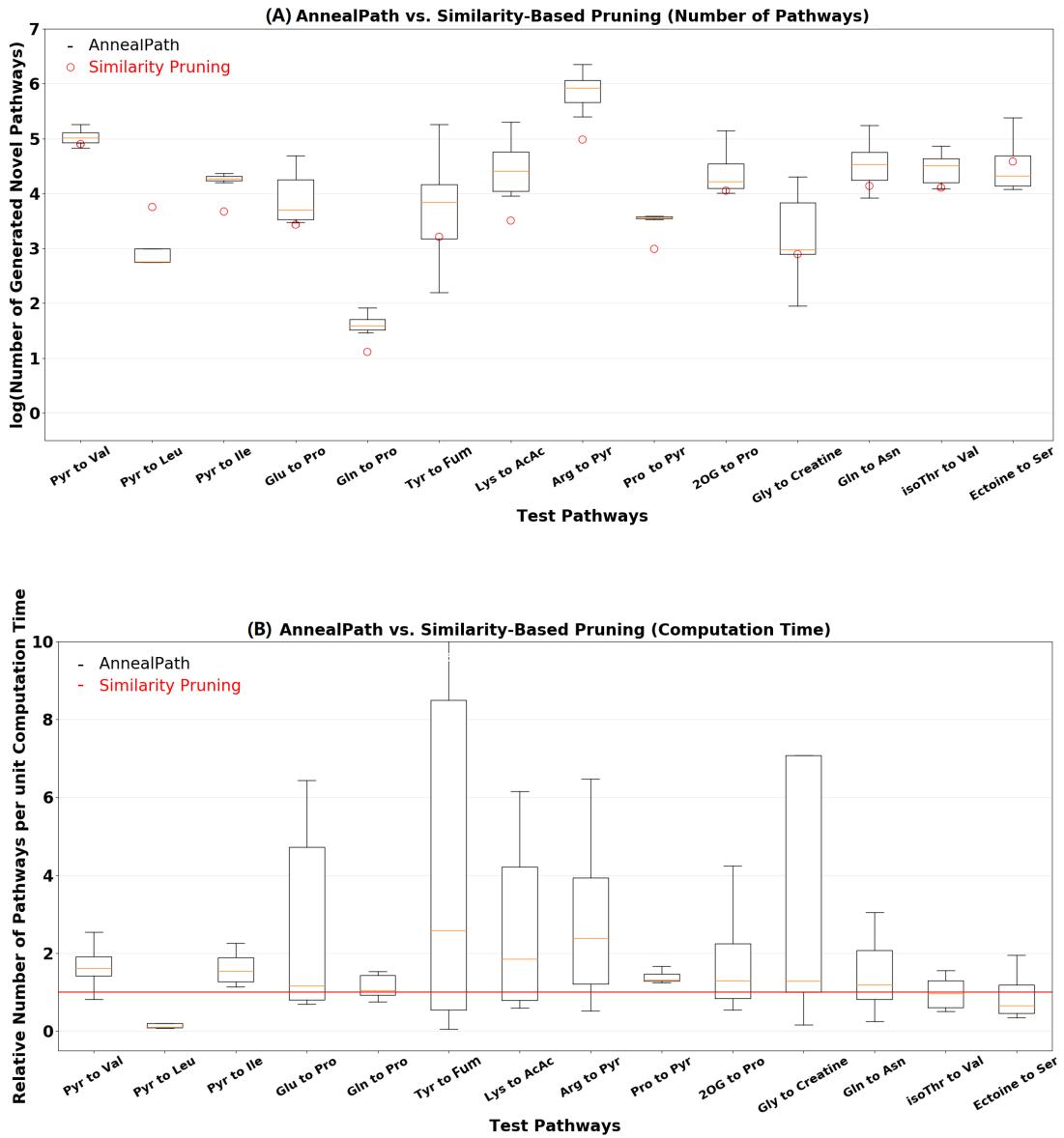


Figure 7: A comparison between the performance of similarity pruning and SA improved pruning (with half of the number of compounds accepted after each step of network expansion). The number of pathways enumerated and computation efficiency are tested by compound pairs found in amino acid biosynthesis/degradation pathways in KEGG/ATLAS (KEGG:map00290, ATLAS:01230, KEGG:map00220, map00330, KEGG:map00350, KEGG:map00310, KEGG:map00330). (A) Number of pathways generated by pruning based on similarity score and the improved pruning algorithm in *Anneal Path*. (B) Relative quantity of number of pathways generated per computation time. The abbreviations for metabolites in the figure were as follows: Pyruvate (Pyr), Valine (Val), Leucine (Leu), Isoleucine (Ile), Glutamate (Glu), Proline (Pro), Glutamine (Gln), Tyrosine (Tyr), Arginine (Arg), Acetoacetate (AcAc), Fumarate (Fum), 2-Oxoglutarate (2OG), Glycine (Gly), Asparagine (Asn), Threonine (Thr), Serine (Ser).

tree classifier) as compared with poor characterization of reaction distances with different similarity metrics. Spearman’s correlation coefficient and different classification evaluation metrics (i.e., Accuracy, weighted average recall, precision and f1-scores) are used to validate the prediction model. Prediction performance of the model trained on synthetic dataset (collected from putative pathways generated by Anneal Path) is also shown in Table 1 together with the comparison to the same list of similarity metrics.

Figure 8 shows a comparison between the stacked distributions of predicted reaction distances (Figure 8A) and similarity scores of compound pairs (Figure 8B) in MetaNetX database. The more separated distributions and reduced overlap areas under different fitted density curves in Figure 8(A) as compared with 8(B) shows a great improvement of the reaction distance characterization made by the new distance metric. Red triangle markers shows the mean value of each distribution which implies a weak (negative) correlation between similarity scores and actual distances for pathway length above six. Heatmaps (confusion matrix) are also plotted for visualizing and comparing the characterizations of reaction distances as shown in Figure 8(C) and 8(D). It can be observed from the heatmap that over 90% of the distance estimations falls in ± 1 range of the real value, which shows its capability of identifying the shortest pathway length between compound pairs.

A same set of four plots as Figure 8 are also made to validate the performance of the distance prediction model trained on the synthetic dataset mentioned earlier (see Figure S27 in Supporting Information). Reaction distance data obtained from the synthetic pathway network allows the model to get trained repeatedly by data generated over combinations of transformations and therefore achieves an even higher prediction accuracy. Training over a pathway network generated by six-step expansion of the reaction rules has achieved an accuracy of 80% and a Spearman’s correlation of 0.863 (also see Table 1).

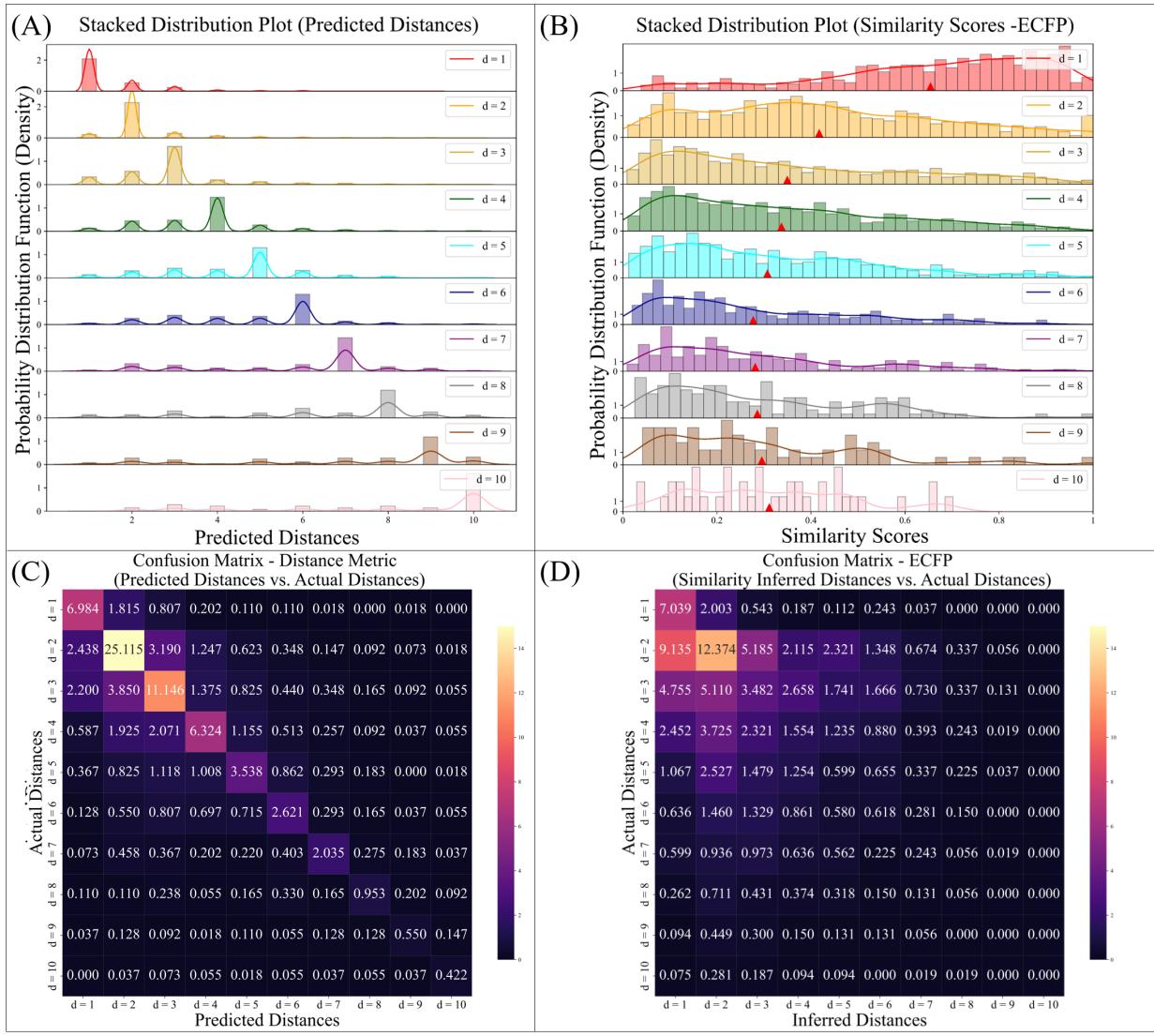


Figure 8: A comparison between the distributions of similarity scores and distributions of the predicted reaction distances for pairs of compounds in MetaNetX data obtained through construct a reaction network using MetaNetX reactions. (A) Stacked distribution of predicted number of reaction steps required for structure transformation between pairs of molecular structures in the test dataset (based on ECFP chemical encodings and decision tree model). (B) Stacked distribution of similarity scores for all compound pairs in the test dataset (similarity scores estimated using ECFP). Red triangle markers are mean values of the similarity scores distributions. The distance prediction model allows a much better estimation (as shown in panel A) of the number of transformation steps between compounds. (C) Heatmap (confusion matrix) showing the trained classification model's performance on the test dataset. Percentages in the diagonal blocks are accurate predictions. It is obvious that most predictions fall into the diagonal blocks. (D) Heatmap that visualize the correlation between actual distances and similarity-inferred distances (calculated by a commonly-used formula that converts a similarity score to a distance value, $similarity(\tau) = \frac{1}{distance(d)+1}$)

Table 1: Full classification results of distance prediction model trained on two dataset (MetaNetX data and synthetic data) and comparison to the characterization of reaction distances by commonly used similarity metrics

Dataset	Evaluation Metrics ^a		<i>R</i>	Accuracy	Recall	Precision	F1
MetaNetX data	New Distance Metric	0.703	0.60	0.60	0.59	0.59	
	ECFP	0.359	0.26	0.26	0.25	0.24	
	Atom-pairs	0.363	0.25	0.25	0.26	0.24	
	Top. Torsions	0.295	0.22	0.22	0.22	0.20	
	MACCS	0.342	0.22	0.21	0.22	0.19	
Synthetic data	New Distance Metric	0.869	0.80	0.80	0.80	0.80	
	ECFP	0.586	0.27	0.27	0.23	0.16	
	Atom-pairs	0.638	0.30	0.30	0.30	0.22	
	Top. Torsions	0.463	0.32	0.32	0.33	0.32	
	MACCS	0.546	0.21	0.21	0.23	0.10	

^aHere, *R* represents Spearman's correlation. Recall, Precision, and F1-score reported are weighted average values.

Validation of the Integrated ML based Distance Predictor and path-way prediction using *Anneal Path*

The ability of finding biologically relevant pathways has been tested on the list of 20 "Golden Set Pathways" tested in *Reinforcement Learning for Bioretrosynthesis*.⁴⁶ All 20 pathways in the golden data set are present in the solution space of Anneal Path which validates its ability of finding biologically correct solutions. Figure S7-S26 in the Supporting Information shows the graphical outputs of the 20 pathways found in the search results.

A comparison of actual pathway length of the 20 pathways and pathway length predicted by the machine learning based reaction distance predictor is shown in Figure 9 below. The predictive performance is a little lower than that of the test set mainly because a number of chemical descriptors (ECFPs) of the compounds in the 20 golden set pathways are not presented in the training set for learning the distance predictor. However, among the 20 two-to-six steps pathways, 16 predictions of pathway length falls in ± 1 range of the real value, which shows a fairly good predictive power. This validates the pathway length predictor's

ability on unknown metabolites that are not in current biochemical databases.

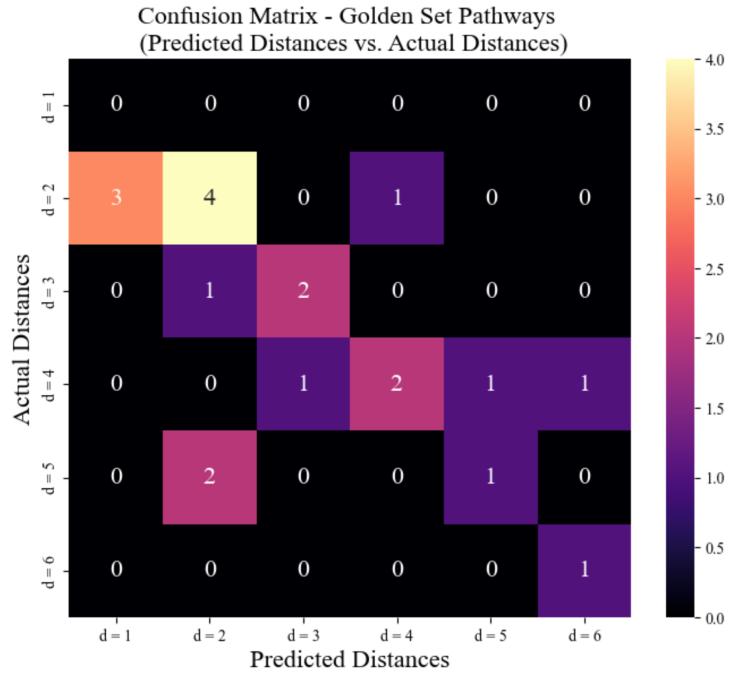


Figure 9: Pathway length of 20 golden set pathways vs. predicted pathway length between each pair of source and target chemicals of 20 golden set pathways

Conclusion

In this study, we described how *Anneal Path* does a more efficient and comprehensive searching for chemically plausible novel pathways through using a hypergraphic putative reaction network and a simulated annealing modified pruning algorithm. We demonstrated how the hypergraph representation of the pathway network can be utilized to explore a broadened solution space and help to identify a greater number of novel pathways than if we had only used simple graph representations. We showed for such a hypergraphic reaction network, screening the nodes using chemical similarity is inefficient. Instead of modifying the chemical similarity metric, we introduced an optimization approach which tolerates the fluctuations in the scoring system which improved the pathway network construction. With a subsequent pathway prioritization process, *Anneal Path* can assist users to efficiently enumerate various potential novel pathways with user-specified sources, target and customized transformation rules.

Acknowledgement

This work was supported by the Natural Sciences and Engineering Research Council (NSERC), the NSERC Industrial Biocatalysis Network (IBN), Biochemicals from Cellulosic Biomass (BioCeB) grant from the Ontario Research Fund (Research Excellence) and a grant from the Genome Canada Genomics Applied Partnership Program (GAPP).

ASSOCIATED CONTENT

Supporting Information available: A set of 393 reaction rules with relatively lower level of specificity from our manually-written in-house reaction database that is used for generating pathways (see https://github.com/LMSE/Anneal-Path/blob/main/rxn_rule/APrules.csv) with details discussed in Supporting Information.

References

- (1) Lee, J.W., Na, D., Park, J.M., Lee, J., Choi, S., Lee, S.Y., (2012) Systems metabolic engineering of microorganisms for natural and non-natural chemicals., *Nat. Chem. Biol.*, **8**, 536–546.
- (2) Prather, K.L.J., Martin, C.H., (2008) De novo biosynthetic pathways: rational design of microbial chemical factories., *Curr. Opin. Biotechnol.*, **19**, 468–474.
- (3) Mederma, M. H., van Raaphorst, R., Takano, E., and Breitling, R., (2012) Computational tools for the synthetic design of biochemical pathways., *Nat. Rev. Microbiol.*, **10**, 191–202.
- (4) Wang, J., Jain, R., Shen, X., Sun, X., Cheng, M., Liao, J.C., Yuan, Q., Yan, Y., (2017) Rational engineering of diol dehydratase enables 1,4-butanediol biosynthesis from xylose., *Metab. Eng.*, **40**, 148–156.
- (5) Hadadi, N., Hatzimanikatis, V., (2015) Design of computational retrobiosynthesis tools for the design of de novo synthetic pathways., *Curr. Opin. Chem. Biol.*, **28**, 99–104.
- (6) Henson, A. B., Gromski, P. S., Cronin, L., (2018) Designing Algorithms To Aid Discovery by Chemical Robots, *ACS Cent. Sci.* , **4**, 7, 793–804.
- (7) Claassen, N. J., Bordnaba-Florit, G., Cotton, A.R., Maria, A., Finger-Bou, Max., Friedeheim, L., Giner-Laguarda, N., Munar-Palmer, M., Newell, W., Scarinci, G., Verbunt, J., de Vries, S. T., Yilmaz, S., Bar-Even, A., (2020) Replacing the Calvin cycle with the reductive glycine pathway in Cupriavidus necator, *Metabolic Engineering* , **62**, 30-41, 1096-7176.
- (8) Carbonell, P., Faulon, J.L.,(2010). Molecular signatures-based prediction of enzyme promiscuity., *Bioinformatics*, **26**, 2012-2019.

- (9) Jeffryes, J., Strutz, J., Henry, C., Tyo, K.E., (2019). Metabolic In Silico Network Expansions to Predict and Exploit Enzyme Promiscuity., *Bioinformatics*, **26**, 2012-2019.
- (10) Hatzimanikatis, V., Li, C., Ionita, J.A., Henry, C.S., Jankowski, M.D., Broadbelt, L.J., (2005) Exploring the diversity of complex metabolic networks., *Bioinformatics*, **21**, 1603–1609.
- (11) Calhoun, S., Korczynska, M., Wichelecki, D.J., Francisco, B.S., Zhao, S., Rodionov, D.A., Vetting, M.W., Al-Obaidi, N.F., Lin, H., O'Meara, M.J., Scott, D.A., Morris, J.H., Russel, D., Almo, S.C., Osterman, A.L., Gerlt, J.A., Jacobson, M.P., Shoichet, B.K., Sali, A., (2018) Prediction of enzymatic pathways by integrative pathway mapping., *eLife*, **7**, e31097.
- (12) Carbonell, P., Planson A.G., Fichera, D., Faulon, J.L., (2011). A retrosynthetic biology approach to metabolic pathway design for therapeutic production., *BMC Syst. Biol.*, **5**, 112.
- (13) Yim, H., Haselbeck, R., Niu, W., Pujol-Baxley, C., Burgard, A., Boldt, J., Khandurina, J., Trawick, J.D., Osterhout, R.E., Stephen, R., Estadilla, J., Teisan, S., Schreyer, H.B., Andrae, S., Yang, T.H., Lee, S.Y., Burk, M.J., Van Dien, S., (2011) Metabolic engineering of Escherichia coli for direct production of 1,4-butanediol., *Nat. Chem. Biol.*, **7**, 445–452.
- (14) Schwaller, P., Petraglia, R., Zullo, V., Nair, V.H., Haeuselmann, R.A., Pisoni, R., Bekas, C., Iuliano, A. and Laino, T., (2020) Predicting retrosynthetic pathways using transformer-based models and a hyper-graph exploration strategy, *Chem. Sci.*, **11**, 3316–3325.
- (15) Hadadi, N., Hafner, J., Shajkofci, A., Zisaki, A., Hatzimanikatis, V., (2016) ATLAS of Biochemistry: A Repository of All Possible Biochemical Reactions for Synthetic Biology and Metabolic Engineering Studies., *ACS Synth. Biol.*, **5**, 1155–1166.

- (16) Campodonico, M.A., Andrews, B.A., Asenjo, J.A., Palsson, B.O., Feist, A.M., (2014) Generation of an atlas for commodity chemical production in *Escherichia coli* and a novel pathway prediction algorithm, GEM-Path. *Metab. Eng.*, **25**, 140-158.
- (17) Sivakumar, T.V., Giri, V., Park, J.H., Kim, T.Y., Bhaduri, A., (2016) ReactPRED: a tool to predict and analyze biochemical reactions., *Bioinformatics*, **32(22)**, 3522-3524.
- (18) Duigou, T., du Lac, M., Carbonell, P., Faulon, J.L., (2018) RetroRules: a database of reaction rules for engineering biology., *Nucleic Acids Res.*, **47**, D1229–D1235.
- (19) Delepine, B., Duigou, T., Carbonell, P., Faulon, J.L., (2017) RetroPath2.0: A retrosynthesis workflow for metabolic engineers., *Metab. Eng.*, **45**, 158–170.
- (20) Carbonell, P., Parutto, P., Baudier, C., Junot, C., Faulon, J.L., (2014) RetroPath: Automated Pipeline for Embedded Metabolic Circuits., (2013) *ACS Synth Biol.*, **3,8**, 565-577.
- (21) Gao, H., Struble, T. J., Coley, C.W., Wang, Y., Green, W. H. and Jensen, K.F.(2018) Using Machine Learning To Predict Suitable Conditions for Organic Reactions, *ACS Cent. Sci.*, **4**, 11, 1465–1476.
- (22) Moriya, Y., Shigemizu, D., Hattori, M., Tokimatsu, T., Kotera, M., Goto, S. and Kanehisa, M.(2010) PathPred: an enzyme-catalyzed metabolic pathway prediction server, *Nucleic Acids Res.*, **38**, pp. W138-W143.
- (23) Chowdhury, A., Maranas, C., (2015) Designing overall stoichiometric conversions and intervening metabolic reactions., *Sci. Rep.*, **5**, 16009.
- (24) Carbonell, P., Fichera, D., Pandit, S.B., Faulon, J.-L. (2012) Enumerating metabolic pathways for the production of heterologous target chemicals in chassis organisms., *BMC. Syst. Biol.*, **6**, 10.

- (25) Heath, A.P., Bennett, G.N., Kavraki, L.E., (2011) An Algorithm for Efficient Identification of Branched Metabolic Pathways., *Comput. Biol.*, **18**, 1575–1597.
- (26) Pertusi, D.A., Stine, A.E., Broadbelt, L.J., Tyo, K.E.J., (2014) Efficient searching and annotation of metabolic networks using chemical similarity., *Bioinformatics*, **31**, 1016–1024.
- (27) O’Boyle, N.M., Banck, M., James, C.A., Morley, C., Vandermeersch, T., Hutchison, G.R., (2011) Open Babel: An open chemical toolbox., *J. Cheminform.*, **3**, 33.
- (28) O’Boyle, N.M., Banck, M., James, C.A., Morley, C., Vandermeersch, T., Hutchison, G.R., (2008) Pybel: a Python wrapper for the OpenBabel cheminformatics toolkit., *Chem. Cent. J.*, **2**, 5.
- (29) Coley, C. W., Rogers, L., Green, W. H. and Jensen, K. F., (2017) Computer-Assisted Retrosynthesis Based on Molecular Similarity, *ACS Cent. Sci.* , **3**, 12, 1237-1245.
- (30) Wang, L., Dash, S., Ng, C.Y. and Maranas, C. D., (2017) A review of computational tools for design and reconstruction of metabolic pathways, *Synthetic and Systems Biotechnology* , **2**, 4, 243-252, ISSN 2405-805X
- (31) Durant, J.L., Leland, B.A., Henry, D.R., Nourse, J.G., (2002) Reoptimization of MDL keys for use in drug discovery., *J Chem Inf Comput Sci.*, **42**, 6, 1273–1280.
- (32) Rogers, D., Hahn, M., (2010) Extended-Connectivity Fingerprints., *J. Chem. Inf. Model.*, **50**, 5, 742-754.
- (33) Carhart, R.E.; Smith, D.H.; Venkataraghavan R., (1985) Atom Pairs as Molecular Features in Structure-Activity Studies: Definition and Applications., *J. Chem. Inf. Comp. Sci.*, **25**, 64-73.

- (34) Raymond, J. W., Willett, P., (2002) Maximum common subgraph isomorphism algorithms for the matching of chemical structures., *J. COMPUT. AID MOL. DES.*, **16**, 7, 521–533.
- (35) Daylight, (2017) Daylight Theory Manual., <http://www.daylight.com/dayhtml/doc/theory/>.
- (36) Landrum, (2016) RDKit: Open-source Cheminformatics., <http://www.rdkit.org/>.
- (37) Koda, P., Hoksza, D., (2015) Exploration of Topological Torsion Fingerprints., *IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pp822-828.
- (38) Gallo, G., Longo, G., (1993) Directed Hypergraphs and Applications., *Discrete Applied Mathematics.*, **42**, 177-201.
- (39) Noor, E., Eden, E., Milo, R., Alon, U., (2010) Central carbon metabolism as a minimal biochemical walk between precursors for biomass and energy., *Mol. Cell.*, **39**, (5) 809-820.
- (40) Moretti, S., Martin, O., Van Du Tran, T., Bridge, A., Morgat, A., Pagni, M., (2015) MetaNetX/MNXref – reconciliation of metabolites and biochemical reactions to bring together genome-scale metabolic networks., *Nucleic Acids Res.*, **44**, D523–D526.
- (41) Kirkpatrick, S., Gelatt, C., Vecchi, M., (1983) Optimization by Simulated Annealing., *Science*, **220**,(4598) 671-680.
- (42) Dolan, W.B., Cummings, P.T., LeVan, M.D., (1989) Process optimization via simulated annealing: Application to network design., *Science*, **35**,(5).
- (43) Rafael, G., Jennifer, W., et al., (2017) Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules, *ACS Cent. Sci.*, **42268**, 276.
- (44) Sutskever, I., Vinyals, O., and Le, Q. V., Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, Eds. Curran (2014) Sequence to sequence learning with neural networks, *Advances in Neural Information Processing Systems*, **27**, pp. 3104–3112.

- (45) Sutskever, I., Vinyals, O., and Le, Q. V., Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, Eds. Curran (2014) Sequence to sequence learning with neural networks, *Advances in Neural Information Processing Systems*, **27**, pp. 3104–3112.
- (46) Mathilde Koch, Thomas Duigou, and Jean-Loup Faulon (2020) Reinforcement Learning for Bioretrosynthesis, *ACS Synthetic Biology*, **9** (1), 157-168.
- (47) Kanehisa, M., Sato, Y., Furumichi, M., Morishima, K., and Tanabe, M., (2019) New approach for understanding genome variations in KEGG., *Nucleic Acids Res.* , **47**,D590-D595.

Table of Contents/Abstract Graphics

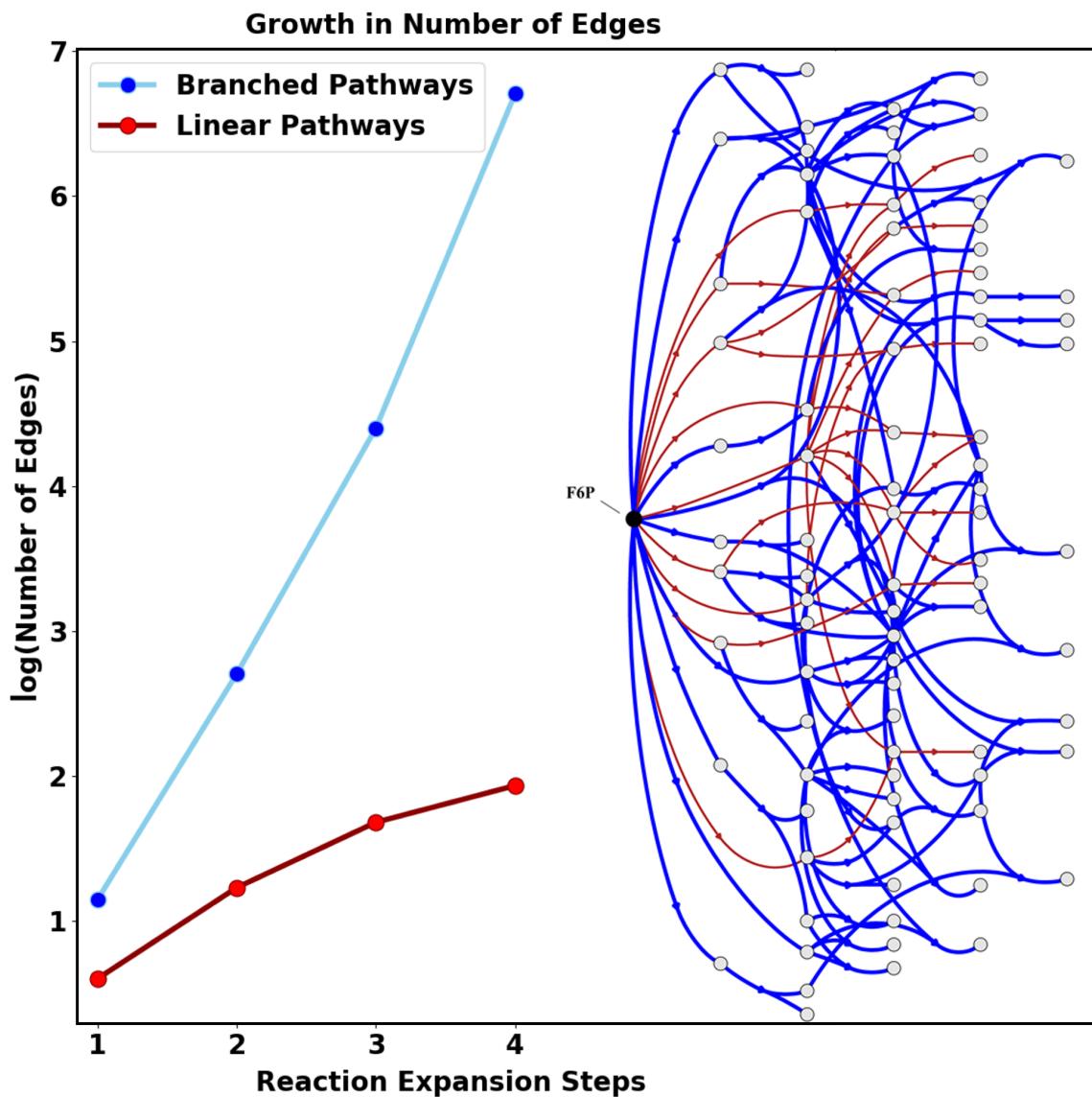


Figure 10: Abstract Graphics