

Discovery and Evaluation of Biosynthetic Pathways for the Production of Five Methyl Ethyl Ketone Precursors

Milenko Tokic¹, Noushin Hadadi¹, Meric Ataman¹, Dário Neves², Birgitta E. Ebert²,

Lars M. Blank², Ljubisa Miskovic^{1,*}, Vassily Hatzimanikatis^{1,**}

¹ Laboratory of Computational Systems Biotechnology (LCSB), Swiss Federal Institute of Technology (EPFL), CH-1015 Lausanne, Switzerland.

² Institute of Applied Microbiology (iAMB), Aachen Biology and Biotechnology (ABBt), RWTH Aachen University, D-52056 Aachen, Germany.

* Corresponding author at Laboratory of Computational Systems Biotechnology (LCSB), École Polytechnique Fédérale de Lausanne (EPFL), CH-1015 Lausanne, Switzerland.

Email: ljubisa.miskovic@epfl.ch,

Phone: + 41 (0)21 693 98 92 Fax: +41 (0)21 693 98 75

** Corresponding author at Laboratory of Computational Systems Biotechnology (LCSB), École Polytechnique Fédérale de Lausanne (EPFL), CH-1015 Lausanne, Switzerland

Email: vassily.hatzimanikatis@epfl.ch

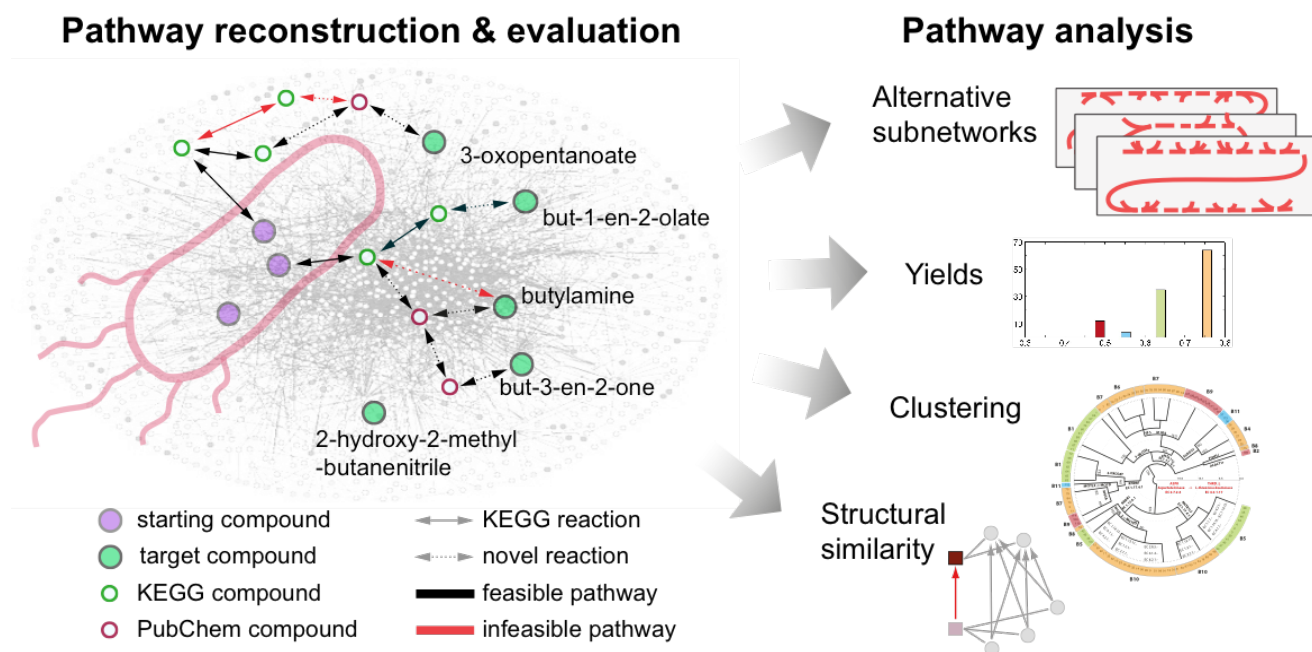
Phone: + 41 (0)21 693 98 70 Fax: +41 (0)21 693 98 75

Abstract

The limited supply of fossil fuels and the establishment of new environmental policies shifted research in industry and academia towards sustainable production of the 2nd generation of biofuels, with Methyl Ethyl Ketone (MEK) being one promising fuel candidate. MEK is a commercially valuable petrochemical with an extensive application as a solvent. However, as of today, a sustainable and economically viable production of MEK has not yet been achieved despite several attempts of introducing biosynthetic pathways in industrial microorganisms. We used BNICE.ch as a retrobiosynthesis tool to discover all novel pathways around MEK. Out of 1'325 identified compounds connecting to MEK with one reaction step, we selected 3-oxopentanoate, but-3-en-2-one, but-1-en-2-olate, butylamine, and 2-hydroxy-2-methyl-butanenitrile for further study. We reconstructed 3'679'610 novel biosynthetic pathways towards these 5 compounds. We then embedded these pathways into the genome-scale model of *E. coli*, and a set of 18'925 were found to be most biologically feasible ones based on thermodynamics and their yields. For each novel reaction in the viable pathways, we proposed the most similar KEGG reactions, with their gene and protein sequences, as candidates for either a direct experimental implementation or as a basis for enzyme engineering. Through pathway similarity analysis we classified the pathways and identified the enzymes and precursors that were indispensable for the production of the target molecules. These retrobiosynthesis studies demonstrate the potential of BNICE.ch for discovery, systematic evaluation, and analysis of novel pathways in synthetic biology and metabolic engineering studies.

Keywords: Novel synthetic pathways, pathway feasibility, pathway similarity,
Methyl Ethyl Ketone, 2-butanone, *E. coli*.

Graphical abstract



Introduction

Limited reserves of oil and natural gas and the environmental issues associated with their exploitation in the production of chemicals sparked off current developments of processes that can produce the same chemicals from renewable feedstocks using microorganisms.¹⁻³ A fair amount of these efforts focuses on a sustainable production of the 2nd generation of biofuels.

Compared to the currently used fossil fuels and bioethanol, these 2nd generation biofuels should provide lower carbon emissions, higher energy density, and should be less corrosive to engines and distribution infrastructures. Recently, a large number of potential candidates has been proposed, such as n-butanol, isobutanol, 2-methyl-1-butanol or 3-methyl-1-butanol⁴, C13 to C17 mixtures of alkanes and alkenes⁵, fatty esters, fatty alcohols¹, and Methyl Ethyl Ketone (MEK) also referred to as 2-butanone⁶.

For many of these chemicals, natural microbial producers are not known and novel biosynthetic pathways for their production are yet to be discovered.^{7, 8} Even when production pathways for target chemicals are known, it is important to find alternatives in order to further reduce cost and greenhouse emissions, and as well to avoid possible patent issues.

Computational tools are needed to assist in the design of novel biosynthetic pathways because they allow exhaustive generation of possible alternatives and evaluation of their properties and prospects for producing target chemicals.⁸ For instance, computational tools can be used to assess the performance of a production pathway operating in one organism in another host organism. They can also be used to predict, prior to experimental pathway implementation, the yields across organisms of a particular pathway in producing a target molecule.

There are different computational tools for pathway prediction available in the literature.⁹⁻¹⁹ An important class of these tools is based on the concept of generalized enzyme reaction rules, which were introduced by Hatzimanikatis and co-workers.^{20, 21} These rules emulate the functions of enzymes, and they can be used to apply *in silico* biotransformations over a wide range of substrates.⁸ Most of the implementations of this concept appear in the context of retrobiosynthesis, where the algorithm generates all possible pathways from a target compound towards desired precursors in an iterative backward manner.^{3, 7-9, 14, 16, 19-25}

In this study, we used the retrobiosynthesis framework of BNICE.ch^{8, 9, 20-25} to explore the biotransformation space around Methyl Ethyl Ketone (MEK). Besides acetone, MEK is the most commercially produced ketone with broad applications as a solvent for paints and adhesives and as a plastic welding agent.²⁶ MEK shows superior characteristics compared to the existing fuels in terms of its thermo-physical properties, increased combustion stability at low engine load, and cold boundary conditions, while decreasing particle emissions.²⁷ There is no known native producer of MEK, but in the recent studies this molecule was produced in *E. coli*^{28, 29} and *S. cerevisiae*⁶ by introducing novel biosynthetic pathways. To convert 2,3-butanediol to MEK, Yoneda *et al.*³⁰ introduced into *E. coli* a B-12 dependent glycerol dehydratase from *Klebsiella pneumoniae*. Srirangan *et al.*²⁹ expressed in *E. coli* a set of promiscuous ketothiolases from *Cupriavidus necator* to form 3-ketovaleryl-CoA, and they further converted this molecule to MEK by expressing acetoacetyl-CoA:acetate/butyrate:CoA transferase and acetoacetate decarboxylase from *Clostridium acetobutylicum*. In *S. cerevisiae*, Ghiaci *et al.*⁶ expressed a B12-dependent diol dehydratase from *Lactobacillus reuteri* to convert 2,3-butanediol to MEK. Alternatively, hybrid biochemical/chemical approaches were proposed where

precursors of MEK were biologically produced through fermentations and then catalytic processes were used to produce MEK.^{30, 31}

We used the BNICE.ch algorithm to generate a network of potential biochemical reactions around MEK, and we identified 159 biochemical and 1'166 chemical compounds one reaction step away from MEK (Table S1 - Supporting Information). We considered as biochemical compounds the ones that we found in the KEGG^{32, 33} database, and as chemical compounds the ones that we found in the PubChem^{34, 35} but not in the KEGG database. A set of 154 compounds appeared in both databases. Out of these 1'325 compounds, 2-hydroxy-2-methyl-butanenitrile (MEKCNH) was the only KEGG compound connected to MEK through a KEGG reaction (KEGG R09358). For further study, we chose MEKCNH along with three KEGG compounds: 3-oxopentanoate (3OXPNT), but-3-en-2-one (MVK) and butylamine (BuNH₂), and one PubChem compound: 1-en-2-olate (1B2OT). The latter four compounds were chosen based on two important properties: (i) their simple chemical conversion to MEK, e.g., 3OXPNT spontaneously decarboxylates to MEK; and (ii) their potential use as precursor metabolites to further produce a range of other valuable chemicals.³⁶⁻³⁸ MVK can be converted to MEK by a 2-enoate reductase from *Pseudomonas putida*, *Kluyveromyces lactis* or *Yersinia bercovieri*,³⁹ however, these reactions are not catalogued in KEGG. Similarly, 3-OXPNT can be decarboxylated to MEK by acetoacetate decarboxylase from *Clostridium acetobutylicum*.²⁹ In contrast, there are no known enzymes that can convert 1B2OT and BuNH₂ to MEK.

We have reconstructed all possible novel biosynthetic pathways (3'679'610 in total) up to a length of 4 reaction steps from the central carbon metabolites of *E. coli* towards the 5 compounds mentioned above. We evaluated the feasibility of these 3'679'610 pathways with respect to the mass and energy balance, and we found

18'925 thermodynamically feasible pathways which we further ranked with respect to their carbon yields. We identified the metabolic subnetworks that were carrying fluxes when the optimal yields were attained, and we determined the minimal sets of precursors and the common routes and enzymes for production of the target compounds.

Results and Discussion

Generated metabolic network around Methyl Ethyl Ketone

We used the retrobiosynthesis algorithm of BNICE.ch to reconstruct the biochemical network around MEK. BNICE.ch^{8, 9, 20-25} is a computational framework that takes advantage of the biochemical knowledge derived from the thousands of known enzymatic reactions to predict all possible biotransformation pathways from known compounds to desired target molecules. We applied BNICE.ch and generated all compounds and reactions that were up to five generations away from MEK (Figure 1).

To start the reconstruction procedure, we provided the initial set of compounds that contained 26 cofactors along with MEK (Table S2 - Supporting Information). In the first BNICE.ch generation, we produced 6 biochemical and 25 chemical compounds connected through 48 reactions to MEK. Interestingly, among these reactions were also the ones proposed by Yoneda *et al.*³⁰, Srirangan *et al.*²⁹ and Ghiaci *et al.*⁶

After five generations, a total of 13'498 compounds were generated (Figure 1.a). Out of these, 749 were biochemical and the remaining 12'749 were chemical compounds. We could also find 665 out of the 749 biochemical compounds in the PubChem database. All generated compounds were involved in 65'644 reactions, out of which 560 existed in the KEGG database and the remaining 65'084 were novel reactions (Figure 1.b). A large majority of the discovered reactions (67%) were

oxidoreductases, 15.4% were lyases, 8.6% were hydrolases, 4.3% transferases, 3.6% isomerases and only 0.72% ligases (Figure 1.c). Out of 361*2 bidirectional generalized enzyme reaction rules of BNICE.ch, 369 were required to generate the metabolic network around MEK with the size of 5 reaction steps. As expected from the statistics on the discovered reactions, most of these rules (38%) described the oxidoreductase biotransformation (Figure 1.d).

Although MEK participated in a total of 1'551 reactions only one reaction was catalogued in the KEGG database (KEGG R09358) which connected MEK to MEKCNH. The generated reactions involved 1'325 compounds (159 biochemical and 1'166 chemical) that could be potentially used as MEK precursors.

Figure 1. Growth of the BNICE.ch generated metabolic network over 5 generations. Compounds: biochemical (Panel a, blue) and chemical (Panel a, red). Reactions: KEGG (Panel b, blue) and novel (Panel b, red). Discovered reactions (Panel c) and utilized generalized enzyme reaction rules (Panel d) organized on the basis of their Enzymatic Commission⁴⁰, EC, class.

Pathway reconstruction towards five target compounds

In the pathway reconstruction process, we used as starting compounds 157 metabolites selected from the generated network, which were identified as native *E. coli* metabolites using the *E. coli* genome-scale model iJO1366⁴¹ (Table S3 - Supporting Information). We performed an exhaustive pathway search on the generated metabolic network, and we reconstructed 3'679'610 pathways towards these five target compounds with pathway lengths ranging from 1 up to 4 reaction steps (Table 1). The reconstructed pathways combined consist of 37'448 reactions,

i.e., 57% of the reactions reproduced from the BNICE.ch generated metabolic network.

More than 58% of the discovered pathways were towards BuNH₂, while only 3.8% of the reconstructed pathways were towards 1B2OT, which was the only PubChem target compound (Table 1). Only 33 reconstructed pathways were of length one, and 28 out of them were towards BuNH₂ and none towards 1B2OT. The majority of reconstructed pathways (> 97%) were of length four. These results suggest that the biochemistry of enzymatic reactions favors smaller changes of a molecule structure over several steps.

Table 1. Reconstructed pathways towards five target compounds.

Target compounds	Reconstructed pathways	Reaction steps				Feasible pathways	
		1	2	3	4	FBA	TFA
3-oxopentanoate (3OXPNT)	641'493	1	198	12'222	629'072	361'187	11'145
but-3-en-2-one (MVK)	438'889	1	136	7'554	431'198	57'173	4'117
Butylamine (BuNH ₂)	2'146'890	28	1'236	53'573	2'092'053	27'211	1'177
but-1-en-2-olate (1B2OT)	140'779	0	53	2'905	137'821	30'689	1'826
2-hydroxy-2-methylbutanenitrile (MEKCNH)	311'559	3	94	6'546	304'916	11'151	660
	3'679'610	33	1'717	82'800	3'595'060	487'411	18'925

Evaluation of reconstructed pathways

We performed a series of studies of the 3'679'610 generated pathways to assess their biological feasibility and performance (Methods). The feasibility of the pathways depends on the metabolic network of the chassis organism. Therefore, we embedded each of the reconstructed pathways in the *E. coli* genome-scale model iJO1366 and performed flux balance analysis (FBA) and thermodynamics-based flux analysis (TFA). The directionality of the reactions is an important factor in FBA and TFA⁴²,

and in our studies, we distinguished the following types the reactions: (*R1*) known and novel reactions for which have no information about their directionality; (*R2*) reactions that have preassigned directionality in iJO1366; and (*R3*) reactions that involve CO₂ as a metabolite. We employed 2 types of constraints (Methods): (*C1*) the preassigned directionalities of the *R2* reactions were removed and the directionality of the *R3* reactions was fixed towards decarboxylation; and (*C2*) the preassigned directionalities of the *R2* reactions were kept and the directionality of the *R3* reactions was fixed towards decarboxylation. In *C1* constraints, we removed the preassigned directionalities to explore the possibilities that might be lost due to assumptions about the catalytic activities of enzymes. The catalytic reversibility or irreversibility of enzymes could be altered through protein and evolutionary engineering and enzyme screening.⁴² Unless stated otherwise, for FBA and TFA we applied *C1* constraints.

Flux balance analysis. We used FBA as a prescreening method and we found the pathways that were incompatible with the host organism as they required co-substrates that were absent in the host organism (based on the iJO1366 reconstruction). Out of all reconstructed pathways, only 13.24% (487'411) were FBA feasible (Table 1). Though the largest number of reconstructed pathways were towards BuNH₂, only 1.27% (27'211) of these were FBA feasible. The number of FBA feasible pathways for MEKCNH was also low (3.59%). In contrast, more than 56% of pathways towards 3OXPNT were FBA feasible.

Thermodynamics-based flux analysis. We used TFA and identified 18'925 thermodynamically feasible pathways (0.51% of all generated pathways, or 3.88% of the FBA feasible pathways). The set of TFA feasible pathways involved 3'269 unique

reactions. These results demonstrate that TFA is important for pathway evaluation and screening.

We found BuNH₂ to have the lowest rate of TFA feasible pathways with 0.05% of reconstructed pathways being TFA feasible (Table 1). The highest rate of TFA feasible pathways was again for 3OXPNT (1.74 %). The shortest TFA feasible pathways consisted of 2 reaction steps (21 pathways), whereas a majority of TFA feasible pathways had length 4 (Table 2). All pathways contained novel reaction steps, and only 19 pathways had one novel reaction step (Table 2). All of these 19 pathways were towards MVK, and they all had as intermediates 2-acetolactate and acetoin. The final reaction step converting acetoin to MVK was novel in all of them.

Table 2. Number of known reaction steps versus all reaction steps. Pathways with one novel reaction step are marked in red. All shown pathways are TFA feasible.

		Reaction steps			Feasibility	
		2	3	4	TFA	
Number of known steps in a pathway	0	14	371	7'059	7'444	3-oxopentanoate (3OXPNT)
	1		118	2'956	3'074	
	2			627	627	
	0	4	72	3'196	3'272	but-3-en-2-one (MVK)
	1	1	13	703	717	
	2		2	110	112	
	3			16	16	
	0	2	35	974	1'011	Butylamine (BuNH ₂)
	1		7	139	146	
	2			20	20	
	0		23	1'576	1'599	but-1-en-

	1	10	196	206	2-hydroxy-2-methyl- butanenitrile (MEKCNH)
	2		21	21	
	0	50	380	430	
	1	2	202	204	
	2		26	26	
		21	703	18'201	18'925

Yield analysis. We used TFA to assess the production yield of the feasible pathways from glucose to the target compounds (Table S4 - Supporting Information). We identified pathways for all target compounds that could operate without a loss of carbon from glucose. More than a half of the pathways towards 3OXPNT (57%) could operate with the maximum theoretical yield of 0.774 g/g, i.e., 1Cmol/1Cmol (Figure 2). In contrast, only 4% (27 out of 660) pathways towards MEKCNH could operate with the maximal theoretical yield of 0.66 g/g (Table S4 - Supporting Information). We found that pathways were distributed into several distinct sets rather than being more spread and continuous, i.e., we obtained eleven sets for 3OXPNT, four sets for MEKCNH, 15 sets for BuNH₂, nine sets for 1B2OT and ten sets for MVK (Table S4 - Supporting Information). Interestingly, a discrete pattern in pathway yields was also observed in a similar retrobiosynthesis study for the production of mono-ethylene glycol in *Moorella thermoacetica* and *Clostridium ljungdahlii*.⁴³

Analysis of alternative assumptions on reaction directionalities. Since we found that the directionality of reactions in the network impacts yields, we investigated how the type of alternative constraints *C2* affected the yield distribution. The *R2* reactions

that could operate in both directions with the *C1* constraints applied were unidirectional with the *C2* constraints applied. As expected, these additional constraints reduced flexibility of the metabolic network and some pathways even became infeasible (Table S5 - Supporting Information). With the *C2* constraints, the yields were in general reduced and their distribution was more spread compared to the one obtained using the *C1* constraints. For example, we found with both sets of constraints three alternative pathways for the production of 3OXPNT from acetate via two intermediate compounds: 2-ethylmalate and (3S)-3-hydroxypentanoate. The three alternative pathways had three different cofactor pairs in the final reaction step that converts (3S)-3-hydroxypentanoate to 3OXPNT (Figure F1 - Supporting Information). With the *C1* constraints, the three pathways had an identical yield of 0.642 g/g. In contrast, with the *C2* constraints, the pathway with NADH/NAD cofactor pair in the final step had a yield of 0.537 g/g, the one with NADPH/NADP had a yield of 0.542 g/g, and the one with H₂O₂/H₂O had a yield of 0.495 g/g. These differences in yields are a consequence of the different costs of cofactor production upon adding supplementary constraints.

These results highlight the importance of the choice of constraints in FBA and TFA as they can influence our conclusions on reaction directionalities. Besides, the reaction directionalities have a critical impact on network properties such as gene essentiality or yields.⁴² This suggests that particular caution should be exercised when using “off-the-shelf” models as some of them have *ad hoc* pre-assigned directionalities.^{42, 44, 45} Additionally, this indicates that there is a need for revisiting assumptions on reaction directionalities in the current genome-scale reconstructions. This task can be performed by integrating thermodynamics in metabolic networks and thus allowing for systematical assigning of reaction directionalities.⁴² However, for an accurate

estimation of the reaction directionalities using thermodynamics, it is crucial to consider the contribution of the activities to the Gibbs free energy of reactions instead of using only the standard values^{44, 45} further emphasizing the importance of integrating metabolomics data.

BridgIT analysis. For each novel reaction from the feasible pathways, we identified the most similar KEGG reaction whose gene and protein sequences were assigned to the novel reaction (Methods). The BridgIT⁴⁶ results can be consulted at <http://lcsb-databases.epfl.ch/pathways/> upon subscription.

Identification and analysis of anabolic subnetworks capable of synthesizing target molecules

In pathway reconstruction, we identified the sequence of the main reactions required to produce the target molecules from precursor metabolites in the core network. However, these reactions require additional co-substrates and cofactors that should become available from the rest of the metabolism. In addition, these reactions produce also side products and cofactors that must be recycled by the genome-scale metabolic network in order to have a biologically feasible and balanced subnetwork for the production of the target molecules. Therefore, we identified the active metabolic subnetworks required to synthesize the corresponding target molecule (Methods). We then divided the active metabolic subnetworks into the *core metabolic network*, which included central carbon metabolism pathways^{47, 48}, and the *active anabolic subnetwork* (Figure 2.a, and Methods). Interestingly, we found that on average there were more than 3 alternative anabolic subnetworks per pathway due to the redundant topology of metabolism (Table 3). For example, for 11'145 feasible pathways towards

3OXPNT we identified 35'013 alternative anabolic subnetworks. Overall, we identified 57'139 active anabolic subnetworks from the 18'925 TFA feasible pathways.

Figure 2. Metabolic network representing the production of 3OXPNT from glucose (Panel a). Black lines: reactions pertaining to the core metabolic network. Red lines: reactions pertaining to the active anabolic metabolic subnetwork. Green nodes: metabolites in the core metabolic network. Orange nodes: metabolites in the active anabolic metabolic subnetwork. Yellow nodes: core precursors, i.e., metabolites that connect the core and active anabolic subnetworks. Alternative pathways connecting ribose-5-phosphate, r5p, with 2-deoxy-D-ribose-1-phosphate, 2dr1p (Panel b). Alternative pathways connecting the core metabolites with propanal, Ppal (Panel c).

Table 3. Alternative anabolic subnetworks for 5 target compounds together with their lumped reactions and precursors.

Target compounds	Feasible pathways	Alternative anabolic subnetworks	Unique lumps	Overlapping sets of precursors	Unique precursors
3-oxopentanoate (3OXPNT)	11'145	35'013	4'517	281	40
but-3-en-2-one (MVK)	4'117	10'162	1'762	126	32
Butylamine (BuNH ₂)	1'177	4'259	1'791	102	30
but-1-en-2-olate (1B2OT)	1'826	5'339	1'536	97	30
2-hydroxy-2-methylbutanenitrile (MEKCNH)	660	2'420	794	37	17
	18'925	57'139	10'400		

Next, we computed a lumped reaction for each of the alternative subnetworks (Methods). Out of the 57'139 computed lumped reactions, only 10'400 were unique (Table 3) similar to previous findings from the analysis of the biomass building

blocks in *E. coli*.⁴⁹ Overall, for the five target compounds, there were, on average, more than 5 alternative subnetworks per lumped reaction. For the compound 3OXPNT, we found the largest diversity in alternative subnetworks per lumped reaction, where, on average, more than 7 alternative subnetworks had the same lumped reactions (35'013 alternative subnetworks were lumped into 4'517 unique reactions). In contrast, we observed the smallest diversity for BuNH₂ with approximately three alternative subnetworks per lumped reaction (2'420 alternative subnetworks lumped into 794 unique reactions) (Table 3). An illustrative example of multiple pathways with the same lumped reaction is provided in Figure F2 in the Supporting Information. This result suggests that the overall chemistry and the cost to produce the corresponding target molecule are the same for many different pathways. Since the cost of producing a target molecule depends of the host organism, this implies that the choice of the host organism is important. On the other hand, the multiple alternative options could also provide useful degrees of freedom for synthetic biology and metabolic engineering design.

The 35'013 active anabolic networks towards the production of 3OXPNT were composed of only 394 unique reactions. Out of these 394 reactions, 132 were common with the pathways leading towards the production of all biosynthetic building blocks (BBB) except chorismate, phenylalanine, and tyrosine. This finding suggests that BBBs could be competing for resources with 3OXPNT and that they could affect the production of this compound.

Origins of diversity of alternative anabolic subnetworks. To better understand the diversity in alternative anabolic subnetworks, we performed an in-depth analysis of the two-step pathway from acetyl-CoA and propanal to 3OXPNT, which presented the largest number of alternative anabolic networks (185) among all reconstructed

pathways (Figure 2.a). The smallest anabolic subnetwork of the 185 alternatives consisted of 14 enzymes, whereas the largest one comprised 22 enzymes (Table S6 - Supporting Information). All 185 subnetworks shared five common enzymes: the two enzymes from the reconstructed pathway converting propanal via (3S)-3-hydroxypentanoate to 3OXPNT (with the BNICE.ch assigned third level Enzymatic Commission⁴⁰, EC, numbers 2.3.3.- and 1.1.1.-), two enzymes involved in acetyl-CoA production (phosphopentomutase deoxyribose (PPM2), and deoxyribose-phosphate aldolase (DRPA)), one enzyme converting propionate to propanal (aldehyde dehydrogenase (ALDD3y)) (Figure 2).

The multiplicity of ways to produce acetyl-CoA and propionate contributed to a large number of alternative subnetworks: there were 102 alternative ways of producing acetyl-CoA from ribose-5-phosphate (r5p) via 2-deoxy-D-ribose-1-phosphate (2dr1p) (Figure 2.b) and 9 different ways of producing propionate (Figure 2.c).

There were two major routes to produce 2dr1p within the 102 alternatives. In the first route with 50 alternatives, r5p is converted either to ribose-1-phosphate (in 31 alternatives) or to D-ribose (in 19 alternatives), which are intermediates in producing nucleosides such as adenosine, guanosine, inosine and uridine. These nucleosides are further converted to deoxyadenosine (dad), deoxyguanosine (dgsn) and deoxyuridine (duri) that are ultimately phosphorylated to 2dr1p. In 26 of the remaining 52 alternatives of the second route, r5p is converted to phosphoribosyl pyrophosphate (prpp), which is followed by a transfer of its phospho-ribose group to nucleotides such as AMP, GMP, IMP and UMP. These nucleotides are then converted to 2dr1p by downstream reaction steps. In the remaining alternatives for the second route, r5p is first converted to AMP in one reaction step, and then to 2dr1p via dad and dgsn.

There were 9 alternative routes to produce propionate. In 4 of these, this compound was produced from pyruvate and succinate (Figure 2.a and 2.c), in 3 routes it was produced from aspartate (Figure 2.c), and in 2 routes it was produced from 3-phosphoglycerate and glutamate.

Core precursors of five target compounds. An abundant availability of precursor metabolites is crucial for an efficient production of target molecules.⁵⁰ Here, we defined as *core precursors* the metabolites that connect the core to the active anabolic metabolic subnetworks (Figure 2.a). We analyzed the different combinations of core precursors that appeared in the alternative subnetworks. Our analysis revealed that the majority of subnetworks were connected to the core network through a limited number of core precursors. We found that all 35'013 alternative subnetworks for the production of 3OXPNT were connected to the core network by 281 sets of different combinations among 40 unique core precursors (Table 3). We ranked these sets based on their number of appearances in the alternative networks. The top ten sets appeared in 24'210 subnetworks, which represented 69% of all identified subnetworks for this compound (Table 4). Moreover, the metabolites from the top set (acetyl-CoA, propionyl-CoA, pyruvate, ribose-5-phosphate, and succinate) were the precursors in 8'510 (24.3%) subnetworks for 3OXPNT (Table 4). Ribose-5-phosphate appeared in 9 out of the top ten sets, and it was a precursor in 32'237 (92%) 3OXPNT producing subnetworks.

Table 4. Top ten core precursor combinations for the production of 3OXPNT. Core precursors: acetate (ac), acetyl-CoA (acCoA), aspartate (asp-L), dihydroxyacetone

phosphate (dhap), propionyl-CoA (ppCoA), pyruvate (pyr), ribose-5-phosphate (r5p), succinate (succ), succinyl-CoA (sucCoA).

ac	acCoA	asp-L	dhap	ppCoA	pyr	r5p	succ	sucCoA	No. of sub-networks	No. of feasible pathways
	✓			✓	✓	✓	✓		8'510	624
		✓		✓		✓			5'409	2790
				✓	✓	✓	✓		3'463	920
✓						✓		✓	1'344	672
					✓	✓	✓		1'049	382
	✓					✓	✓		965	191
	✓					✓		✓	956	478
				✓	✓		✓		915	460
	✓		✓			✓			834	419
			✓			✓			765	387
									24'210	7323

Clustering of feasible pathways

The repeating occurrences of core precursors and lumped reactions in the alternative anabolic subnetworks motivated us to identify common patterns in core precursors, enzymes, and intermediate metabolites required to produce the target molecules. To this end, we used the feasible pathways from acetate to 3OXPNT as the test study, and we performed two types of clustering on these 115 pathways (File M1 – Supporting information).

Clustering based on core precursors and byproducts of lumped reactions. We computed 242 lumped reactions corresponding to 115 pathways and 242 subnetworks from the test study (File M1 – Supporting information). We chose the first lumped reaction returned by the solver for each of the 115 pathways, and we clustered them

based on the structural similarity between their core precursors and byproducts (Methods).

The main clustering condition among the 115 pathways was the presence or absence of thioesters, such as AcCoA, in the set of core precursors (Figure 3). There were 56 pathways with CoA-related precursors and 59 pathways that did not require CoA. The pathways from the former group were further clustered subject to the presence of the precursors acCoA (1 pathway), ppCoA (30 pathways), both acCoA and ppCoA (6 pathways), and sucCoA (19 pathways), or the occurrence of the byproducts malonate (maln) or CO₂. The pathways that did not require CoA were further clustered depending on if they had as precursors formate (for) or dhap (27 pathways) or not (32 pathways). The clustering results for the complete set of 242 lumped reactions are provided in the supplementary material (Figure F3 – Supporting information).

Figure 3. Clustering dendrogram of the 115 reconstructed pathways from acetate to 3OXPNT and their respective yields (inset). Pathways were classified based on core precursors (red) and byproducts (green) of their lumped reactions. (R)-CoA denotes the group of thioesters. Abbreviations: 2-oxoglutarate (akg), acetyl-CoA (acCoA), aspartate (asp), dihydroxyacetone phosphate (dhap), formate (for), glycolate (glyclt), malate (mal), malonate (maln), propionyl-CoA (ppCoA), pyruvate (pyr), succinyl-CoA (sucCoA).

In general, we expect the set of precursors and byproducts to affect the pathway yield. Interestingly, the clustering based on core precursors and byproducts of lumped reactions also separated distinctly the pathways based on their yields (Figure 3, inset). Pathways that have acCoA, ppCoA, dhap, and for as precursors have a maximal theoretical yield of 0.774 g/g. In contrast, pathways with 2-oxoglutarate (akg) or

acCoA as precursors, and maln as the byproduct, had the lowest yield (0.483 g/g) from the set of examined pathways.

The clustering also provided insight into the different chemistries behind the analyzed pathways. For most of the pathways, i.e., the ones classified in groups B1-2 and B4-10, there was a clear link between the core precursors and co-substrates of acetate in the first reaction step of the pathways (Figure 3). For example, the pathways from the group B1 have a common first reaction step (EC 2.8.3.-) that converts acetate and 3-oxoadipyl-CoA to 3-oxoadipate (Figure 3). The clustering grouped these pathways together because sucCoA was the core precursor of 3-oxoadipyl-CoA through 3-oxoadipyl-CoA thiolase (3-OXCOAT). Moreover, 3-oxoadipate, a 6-carbon compound, was converted in downstream reaction steps to 3OXPNT, a 5-carbon compound, and one molecule of CO₂ through 18 alternative routes. Similarly, in the single pathway of group B2 the co-substrate in the first reaction step was (S)-methylmalonyl-CoA, which was produced from sucCoA through methylmalonyl-CoA mutase (MMM). This enzyme, also known as sleeping beauty mutase, is a part of the pathway converting succinate to propionate in *E. coli*.⁵¹ Malonate (maln), a 2-carbon compound, was released in the first reaction step, which resulted in a low yield of this pathway (Figure 3).

Despite sharing the first reaction step in which acetate reacted with 2-oxoglutarate to create 2-hydroxybutane 1-2-4-tricarboxylate, the pathways from group B9 were split in two groups with different yields (Figure 3). These two groups differed in the sequences of reactions involved in the reduction of 2-hydroxybutane 1-2-4-tricarboxylate, a 7-carbon compound, to 3OXPNT. In 11 pathways, the yield was 0.483 g/g due to a release of two CO₂ molecules, whereas in one pathway the yield

was 0.644 g/g due to malate being created as a side-product and recycled back to the system.

Pathways from group B3 utilized different co-substrates, such as ATP and crotonoyl-CoA, along with acetate to produce acetaldehyde in the first reaction step. All these pathways shared a common novel reaction step with acetaldehyde and propionyl-CoA as substrates (EC 2.3.1.-).

Finally, group B11 contained the pathways with the intermediate 2-methylcitrate, which was produced from pyruvate (pyr).

The presented clustering analysis has been shown to be very powerful in identifying the features of the large number of pathways. The classification can further guide us to identify the biochemistry responsible for the properties of pathways. Such deeper understanding can provide further assistance for the design and analysis of novel synthetic pathways.

Clustering based on involved enzymes. Although the clustering based on the core precursors and byproducts provided an insight of the chemistry underlying the production of 3OXPNT from acetate, lumped reactions conceal the identity of the enzymes involved in the active anabolic subnetworks. We analyzed the 115 active subnetworks corresponding to 115 pathways (File M1 – Supporting information), and we found that five enzymes were present in all of them: AMP nucleosidase (AMPN), 5'-nucleotidase (NTD6), purine-nucleoside phosphorylase (PUNP2), PPM2 and DRPA, which participated in the production of acetaldehyde from r5p (Figure 4.b).

To find common enzyme routes in these subnetworks, we performed a clustering based on the structural similarity between their constitutive reactions (Methods). The clustering separated 115 subnetworks in two groups depending on the existence (47

subnetworks) or not (68 subnetworks) of a sequence of six enzymes starting with aspartate kinase (ASPK) and ending with L-threonine deaminase (THRD_L), whose product 2-oxobutanoate was converted downstream to 3OXPNT (Figures 4.a and 4.b).

Figure 4. Panel a: Clustering dendrogram of 115 active subnetworks corresponding to 115 reconstructed pathways from acetate to 3OXPNT. Subnetworks were clustered based on enzymes they involved. Panel b: Network structure of 47 subnetworks containing a sequence of six enzymes starting with aspartate kinase (ASPK) and ending with L-threonine deaminase (THRD_L) (groups B_I and B_{II} in Panel A). Core metabolites are marked in green, while the metabolites from the active anabolic networks are marked in orange.

Both groups were further clustered based on a set of enzymes required to produce deoxyadenosine and the downstream metabolite acetaldehyde (Figures 4.a and 4.b). The first subgroup of enzymes, i.e. ribonucleoside-diphosphate reductase (RNDR1), deoxyadenylate kinase (DADK) and NTD6, converted adp to deoxyadenosine. In the second subgroup, atp was transferred to deoxyadenosine via ribonucleoside-triphosphate reductase (RNTR1c2), nucleoside triphosphate pyrophosphorylase (NTPP5) and NTD6 (Figure 4.b). Then, for both subgroups, deoxyadenosine was converted to 2-deoxy-D-ribose 5-phosphate (2dr5p) that was further transformed to acetaldehyde via PPM2 and DRPA (Figures 2 and 4.b).

The clustering based on enzymes allowed us to identify enzymatic routes corresponding to different yields (Figures 4.a, and Figure 3 inset). For example, all pathways that include ASPK and novel reaction steps that involve oxidoreductases of the third level EC class 1.14.13.- and 1.2.1.-, would provide the maximal theoretical

yield of 0.774 g/g (Figure 4.a). Similarly, pathways that contain ALDD3Y, methylisocitrate lyase (MCITL2), and RNTR1C2, but not 3-OXCOAT and ASPK, would also provide the maximal theoretical yield. In contrast, the clustering also permitted us to identify key enzymes participating in pathways with a reduced yield. For example, pathways that contained 3-OXCOAT had a yield of 0.644 g/g. Furthermore, the clustering based on enzymes allowed us to clarify the link between the precursors and the corresponding sequence of enzymes that needed to be active for producing the target molecule. For example, pathways from group B1, which had sucCoA as a core precursor and CO₂ as a byproduct, had the common reaction step 3-OXCOAT (Figure 4.a). Similarly, all pathways from group B4 with core precursors ppCoA and acCoA contained ALDD3Y.

Ranking of biosynthetic pathways and recommendations

We further ranked the corresponding feasible pathways according to their yield, number of reaction steps and enzymes that could be directly implemented or needed to be engineered (Methods). As we saw earlier (e.g. in Figure 3, inset), there are several distinct maximum yield values that can be achieved with all these alternatives rather than a continuous distribution of yields. The clustering analysis suggests that the reason for the discreet distribution is the loss of the carbon atoms in specific steps along the pathways. We obtained the top candidate pathways for each of the target molecules that were likely to produce these compounds with economically viable yields. The highest ranked candidate pathway among all feasible pathways was from pyruvate to 3OXPNT, and it consisted of two KEGG reactions, R00203 and R02527, and two novel reactions of the third level EC class 2.3.3- and 4.2.1.- (Figure 5). The BridgIT⁴⁶ analysis identified KEGG R00472 as the most similar reaction to 2.3.3.-.

KEGG reports that R00472 can be catalyzed by EC 2.3.3.9. Similarly, KEGG R04441 was identified as the most similar to 4.2.1-, and according to the KEGG database this reaction is catalyzed by 4.2.1.9. Therefore, the BridgIT results suggest that the two novel reactions can be catalyzed by the known enzymes. Finally, the pathway could operate with the maximum theoretical yield of 0.774 g/g.

The top candidates were visualized and can be consulted at <http://lcsb-databases.epfl.ch/pathways/> upon subscription.

Figure 5. The highest ranked candidate pathway for production of 3OXPNT (dashed box) and the corresponding active anabolic subnetwork (red) together with the core network (grey).

Further experimental implementation and pathway optimization

After ranking of the top candidate pathways, the experts can choose the most amenable ones for experimental implementation in the host organism. The implemented pathways typically need to be optimized further for economically viable production titers and rates. The optimization is performed through the Design-Built-Test-(Learn) cycle of metabolic engineering⁵²⁻⁵⁴ where stoichiometric⁵⁵⁻⁵⁷ and kinetic models⁵⁸⁻⁶⁵, genome editing^{66, 67} and phenotypic characterization⁶⁸ are combined to improve recombinant strains for production of biochemicals.

Methods

We employed the BNICE.ch framework^{8, 9, 20-25} to generate biosynthetic pathways towards 5 precursors of Methyl Ethyl Ketone: 3-oxopentanoate (3OXPNT), 2-hydroxy-2-methyl-butanenitrile (MEKCNH), but-3-en-2-one (MVK), 1-en-2-olate

(1B2OT) and butylamine (BuNH₂). We tested the set of reconstructed pathways against thermodynamic feasibility and mass balance constraints, and discarded the pathways that were not satisfying these requirements.⁸ Next, we ranked the pruned pathways based on the several criteria, such as yield, number of known reaction steps and pathway length. The steps of the employed workflow are discussed further (Figure 6).

Figure 6. Computational pipeline for discovery, evaluation and analysis of biosynthetic pathways.

Metabolic network generation

We applied the retrobiosynthesis algorithm of BNICE.ch^{8, 43} to generate a biosynthetic network that contains all theoretically possible compounds and reactions that are up to 5 reaction steps away from MEK. The BNICE.ch network generation algorithm utilizes the expert-curated generalized enzyme reaction rules^{20, 21, 69} for identifying all potential compounds and reactions that lead to the production of the target molecules. The most recent version of BNICE.ch includes 361*2 bidirectional generalized reaction rules capable of reconstructing more than 6'500 KEGG reactions.²³ Starting from MEK and 26 cofactors required for the generalized enzyme reaction rules (Table S2 - Supporting Information), we identified the reactions that lead to MEK along with its potential precursors.⁷⁰

Note that for studies where we need to generate a metabolic network that involves only KEGG compounds, mining the ATLAS of Biochemistry²³ is a more efficient procedure than using BNICE.ch retrobiosynthesis algorithm. The ATLAS of

Biochemistry is a repository that contains all KEGG reactions and over 130'000 novel enzymatic reactions between KEGG compounds.

Pathway reconstruction

We performed a graph-based search to reconstruct all possible pathways that connect the five target molecules with the set of 157 native *E. coli* metabolites (Table S3 - Supporting Information).⁴¹ We reconstructed the exhaustive set of pathways up to the length of 4 reaction steps.

Note: If we were interested in pathways containing only KEGG reactions, we would perform a graph-based search over the network mined from the ATLAS of Biochemistry.

Pathway evaluation

It is crucial to identify and select, out of a vast number of generated pathways, the ones that satisfy physico-chemical constraints, such as mass balance and thermodynamics, or the ones that have an economically viable production yield of the target compounds from a carbon source. Evaluation of pathways is context-dependent, and it is important to perform it in an exact host organism model and under the same physiological conditions as the ones that will be used in the experimental implementation. We performed both Flux Balance Analysis (FBA)⁷¹ and Thermodynamic-based Flux Analysis (TFA)^{42, 44, 45, 72, 73} to evaluate the pathways. We have also used BridgIT⁴⁶ to identify candidate sequences for protein and evolutionary engineering in implementing the pathways. The availability of such sequences for the novel reactions and the ability to engineer them should also serve as a metric in ranking the feasibility of the pathways.

Flux balance and thermodynamic-based flux balance analysis. We embedded the generated pathways one at the time in the genome-scale model of *E. coli*, iJO1366⁴¹ (File M1 – Supporting information) and we performed FBA and TFA on the resulting models. In these analyses, we assumed that the only carbon source was glucose and we applied the following two types of constraints on reaction directionalities:

(C1) We removed the preassigned reaction directionalities⁷⁴ from the iJO1366 model (*R2* reactions) with the exception of ATP maintenance (ATPM), and we assumed that the reactions that involve CO₂ (*R3* reactions) are operating in the decarboxylation direction. The lower bound on ATPM was set to 8.39 mmol/gDCW/hr. The remaining reactions (*R1* reactions) were assumed to be bi-directional for FBA, whereas for TFA the directionality of these reactions was imposed by thermodynamics. The purpose of removing preassigned reaction directionalities was to investigate alternative hypotheses about the catalytic reversibility of the enzymes.

(C2) This type of constraints contains the preassigned directionalities of *R2* reactions together with the constraints from *C1*.

Since FBA is less computationally expensive than TFA, we first performed FBA as a prescreening method to identify and discard the pathways: (i) that are not satisfying the mass balance, e.g., pathways that need co-substrates not present in the model; and (ii) that have a yield from glucose to the target compounds lower than a pre-specified threshold. We then performed TFA on the reduced set of pathways to identify the pathways that are bio-energetically favorable and we computed their yields from glucose to 5 target compounds under thermodynamic constraints.

BridgIT analysis. We used BridgIT to find known reactions with associated genes in databases that were the most structurally similar to novel reactions appearing in the feasible pathways.⁴⁶ BridgIT integrates the information about the structures of substrates and products of a reaction into reaction difference fingerprints.⁷⁵ These reaction fingerprints contain the information about chemical groups in substrates and products that were modified in the course of a reaction. BridgIT compares the reaction fingerprints of novel reactions to the ones of known reactions, and quantifies this comparison with the Tanimoto similarity score. The Tanimoto score of 1 signifies that two compared reactions had a high similarity, whereas the Tanimoto score values close to 0 signify that there was no similarity. We used this score to rank the reactions identified as similar to each of the novel reactions. The gene and protein sequences of the highest ranked reactions were proposed as candidates for either a direct experimental implementation or enzyme engineering.

Subnetwork reconstruction analysis

Once the biologically feasible pathways were identified and ranked, we analyzed the parts of the metabolism that carry fluxes when the target compounds are produced from glucose. We considered that the active parts of metabolism consisted of: (i) the core metabolic network (Figure 2.a), which included the central carbon pathways, such as glycolysis, pentose phosphate pathway, tricarboxylic cycle, electron transport chain; and (ii) the active anabolic metabolic subnetworks (Figure 2.a), which contain the reactions that would carry fluxes when a target molecule is produced, but did not belong to the core metabolic network. We also defined the core precursors as metabolites that are connecting the core and the active anabolic metabolic subnetworks (Figure 2.a).

We derived the core metabolic network from the genome-scale reconstruction iJO1366⁴¹ using the redGEM algorithm⁷⁶, and we then used the lumpGEM⁴⁹ algorithm to identify active anabolic subnetworks, and to compute their lumped reactions. The analysis of lumped reactions allowed us to identify core precursors of the target chemicals. We then performed clustering to uncover core precursors, common enzymes, and intermediate metabolites of the anabolic subnetworks leading to the production of the target chemicals.

Identification and lumping of active anabolic subnetworks. The lumpGEM algorithm was applied to identify the comprehensive set of smallest metabolic subnetworks that were stoichiometrically balanced and capable of synthesizing a target compound from a defined set of core metabolites. The set of core metabolites belongs to the core metabolic network, and it includes also cofactors, small metabolites, and inorganic metabolites (Table S7 - Supporting Information). Then, for each target compound and for each identified subnetwork, we used lumpGEM to generate a corresponding lumped reaction. Within this process, we also identified the stoichiometric cost of core metabolites for the biosynthesis of these target compounds.

Clustering of subnetworks. To better understand the chemistry that leads towards the target compounds, we performed two types of clustering on the identified subnetworks:

- Clustering based on the structural similarity between the core precursors and byproducts of the lumped reactions. For each lumped reaction, we removed all non-carbon compounds, such as H₂, O₂, and phosphate, and the cofactor pairs, such as ATP and ADP, NAD⁺ and NADH, NADP⁺ and NADPH, flavodoxin oxidized and reduced, thioredoxin oxidized and reduced, ubiquinone and

ubiquinol. This way, we created a set of substrates (core precursors) and byproducts of interest for each lumped reaction. We then used the *msim* algorithm from the RxnSim⁷⁷ tool to compare the lumped reactions based on individual similarities of their core precursors and byproducts. We finally used the obtained similarity scores to perform the clustering.

- Clustering based on the structural similarity between reactions that constitute the anabolic subnetworks. We used BridgIT to compute structural fingerprints of reactions that constitute the anabolic subnetworks, and we then performed a pairwise comparison of the anabolic subnetworks as follows.

For a given pair of anabolic subnetworks, we carried out a pairwise comparison of their reactions. As a comparison metric we used the Tanimoto distance of the reaction fingerprints.⁷⁸ Based on this comparison, we found the pair of the most similar reactions in two subnetworks and we stored the corresponding distance score. We then removed this pair of reactions from comparison, and we found the next pair of the most similar reactions, we stored their distance score, and we continued with this procedure until we found all pairs of reactions in two subnetworks. Whenever the number of reactions in two subnetworks was unequal, we ignored the unmatched reactions. The distance score between two compared subnetworks was formed as the sum of the distance scores of compared pairs of reactions. This procedure was repeated for all pairs of subnetworks.

We then used the computed distance scores to perform the subnetworks clustering.

Ranking and visualization of *in silico* pathways

In this step, we identified the pathways that were most likely to produce the target molecules. For scoring and ranking the biologically meaningful pathways we used the following criteria: (i) maximum yield from glucose to the target molecules; (ii) minimal number of novel reactions, i.e., enzymes to be engineered; (iii) minimal number of reaction steps in the production pathway; and (iv) highest similarity scores from BridgIT.

Experimental implementation and pathway optimization

The highest ranked candidate pathways can then be experimentally implemented in the host organism and can further be optimized through the Design-Built-Test-(Learn) cycle of metabolic engineering.⁵²⁻⁵⁴

Conclusions

In this work, we used BNICE.ch to reconstruct, evaluate and analyze more than 3.6 million biosynthetic pathways from the central carbon metabolites of *E. coli* towards five precursors of Methyl Ethyl Ketone (MEK), a 2nd generation biofuel candidate. Our evaluation and analysis showed that more than 18'000 of these pathways are biologically feasible. We ranked these pathways based on process- and physiology-based criteria, and we identified gene and protein sequences of the structurally most similar KEGG reactions to the novel reactions in the feasible pathways, which can be used to accelerate their experimental realization. Implementation of the discovered pathways in *E. coli* will allow the sustainable and efficient production of five precursors of MEK (3OXPNT, MVK, 1B2OT, BuNH₂, and MEKCNH), which can also be used as precursors for the production of other valuable chemicals.³⁶⁻³⁸

The pathway analysis methods developed and used in this work offer a systematic way for classifying and evaluating alternative ways for the production of target molecules. They also provide a better understanding of the underlying chemistry and can be used to guide the design of novel biosynthetic pathways for a wide range of biochemicals and for their implementation into host organisms.

The present study shows the potential of computational retrobiosynthesis tools for discovery and design of novel synthetic pathways for complex molecules, and their relevance for future developments in the area of metabolic engineering and synthetic biology.

Acknowledgments

M.T. was supported by the Ecole Polytechnique Fédérale de Lausanne (EPFL) and the ERASYNBIO1-016 SynPath project funded through ERASynBio Initiative for the robust development of Synthetic Biology. N.H. and M.A were supported through the RTD grant MicroScapesX, no. 2013/158, within SystemX, the Swiss Initiative for System Biology evaluated by the Swiss National Science Foundation. L.M. and V.H. were supported by the Ecole Polytechnique Fédérale de Lausanne (EPFL). D.N, B.E.E. and L.M.B. thank the German Federal Ministry of Education and Research for funding (Grant ID 031A459). We would like to thank Joana Pinto Vieira for her help with editing this manuscript, and Homa Mohammadi-Peyhani for her help with generating reaction fingerprints for the clustering studies.

Supporting information

Tables S1-S7 (XLSX)

S1: List of compounds one step away from MEK.

S2: List of starting compounds used for the retrobiosynthesis of BNICE.ch

S3: List of 157 starting compounds.

S4: Yield histograms for 5 MEK precursors obtained with C1 constraints.

S5: Yield histograms for 5 MEK precursors obtained with C2 constraints.

S6: List of 185 alternative pathways from AcCoA and PpCoA to 3OXPNT.

S7: List of metabolites in the core metabolic network.

Figures F1-F3 (PDF)

F1: Three alternative ways to produce 3OXPNT from acetate through 2 intermediate metabolites: 2-ethylmalate and 3-hydroxypentanoate.

F2: Three different pathways from acetate to 3OXPNT sharing the same lumped reaction.

F3: Clustering of all 242 alternatives for production of 3OXPNT from acetate.

File M1 (Matlab .mat file)

M1: Genome-scale model of *E. coli* iJO1366 and 242 active anabolic subnetworks connecting the core metabolism with 115 pathways from acetate to 3OXPNT together with their lumped reactions and stoichiometry.

Abbreviations

Abbreviation	Reaction	EC number
3OXCOAT	3-oxoadipyl-CoA thiolase	2.3.1.174
ACACT1r	Acetyl-CoA C-acetyltransferase	2.3.1.9
ALDD3y	Aldehyde dehydrogenase (propanal, NADP)	1.2.1.4
AMPN	AMP nucleosidase	3.2.2.4
ASPK	Aspartate kinase	2.7.2.4
DADK	Deoxyadenylate kinase	2.7.4.3

DRPA	Deoxyribose-phosphate aldolase	4.1.2.4
FTHFLi	Formate-tetrahydrofolate ligase	6.3.4.3
LALDO3	L-Lactaldehyde:NADP+ 1-oxidoreductase	1.1.1.283, 1.2.1.49
MCITD	2-methylcitrate dehydratase	4.2.1.79
MCITL2	Methylisocitrate lyase	4.1.3.30
MGSA	Methylglyoxal synthase	4.2.3.3
MMM	Methylmalonyl-CoA mutase	5.4.99.2
NTD6	5'-nucleotidase (dAMP)	3.1.3.89
NTPP5	Nucleoside triphosphate pyrophosphorylase	3.6.1.19
PPM2	Phosphopentomutase 2 (deoxyribose)	5.4.2.7
PUNP2	Purine-nucleoside phosphorylase	2.4.2.1
RNDR1	Ribonucleoside-diphosphate reductase (ADP)	1.17.4.1
RNTR1c2	Ribonucleoside-triphosphate reductase (ATP)	1.17.4.2
THRD_L	L-threonine deaminase	4.1.1.19

Abbreviation	Compound	Alternative name	Database IDs	
			KEGG	PubChem
3OXPNT	3-Oxopentanoate	3-Oxopentanoic acid	C02233	5297
MVK	But-3-en-2-one	Methyl Vinyl Ketone	C20701	172232421
BuNH ₂	Butylamine	Butanamine	C18706	124489380
1B2OT	But-1-en-2-olate	1-Butene-2-olate	x	54444500
MEKCH	2-Hydroxy-2-methylbutanenitrile	Methyl Ethyl Ketone Cyanohydrin	C18796	124489470
2dr1p	2-deoxy-D-ribose-1-phosphate	2-Deoxy-alpha-D-ribose 1-phosphate	C00672	3941
2dr5p	2-deoxy-D-ribose 5-phosphate	2-Deoxy-alpha-D-ribose 5-phosphate	C00673	3942
dad	Deoxyadenosine	2'-Deoxyadenosine	C00559	3839
dgsn	Deoxyguanosine	2'-Deoxyguanosine	C00330	3624
duri	Deoxyuridine	2'-Deoxyuridine	C00526	3809

prpp	Phosphoribosyl pyrophosphate	5-Phospho-alpha-D-ribose 1-diphosphate	C00119	3419
ac	Acetate	Acetic acid	C00033	3335
acCoA	Acetyl-CoA	Acetyl coenzyme A	C00024	3326
akg	2-oxoglutarate	2-Ketoglutaric acid	C00026	3328
asp-L	Aspartate	L-Aspartic acid	C00049	3351
dhap	Dihydroxyacetone phosphate	3-Hydroxy-2-oxopropyl phosphate	C00111	3411
ppCoA	Propionyl-CoA	Propionyl coenzyme A	C00100	3400
pyr	Pyruvate	2-Oxopropanoate	C00022	3324
r5p	Ribose-5-phosphate	alpha-D-Ribose 5-phosphate	C03736	6499
succ	Succinate	Butanedionic acid	C00042	3344
sucCoA	Succinyl-CoA	Succinyl coenzyme A	C00091	3391

References:

- [1] Steen, E. J., Kang, Y., Bokinsky, G., Hu, Z., Schirmer, A., McClure, A., Del Cardayre, S. B., and Keasling, J. D. (2010) Microbial production of fatty-acid-derived fuels and chemicals from plant biomass, *Nature* 463, 559-562.
- [2] Lee, S. K., Chou, H., Ham, T. S., Lee, T. S., and Keasling, J. D. (2008) Metabolic engineering of microorganisms for biofuels production: from bugs to synthetic biology to fuels, *Current Opinion in Biotechnology* 19, 556-563.
- [3] 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., and Van Dien, S. (2011) Metabolic engineering of *Escherichia coli* for direct production of 1,4-butanediol, *Nature Chemical Biology* 7, 445-452.
- [4] Atsumi, S., Hanai, T., and Liao, J. C. (2008) Non-fermentative pathways for synthesis of branched-chain higher alcohols as biofuels, *Nature* 451, 86-89.
- [5] Schirmer, A., Rude, M. a., Li, X., Popova, E., and del Cardayre, S. B. (2010) Microbial biosynthesis of alkanes, *Science (New York, N.Y.)* 329, 559-562.
- [6] Ghiaci, P., Norbeck, J., and Larsson, C. (2014) 2-Butanol and Butanone Production in *Saccharomyces cerevisiae* through Combination of a B12 Dependent Dehydratase and a Secondary Alcohol Dehydrogenase Using a TEV-Based Expression System, *PLOS ONE* 9, e102774.
- [7] Cho, A., Yun, H., Park, J. H., Lee, S. Y., and Park, S. (2010) Prediction of novel synthetic pathways for the production of desired chemicals, *Bmc Systems Biology* 4.
- [8] Hadadi, N., and Hatzimanikatis, V. (2015) Design of computational retrobiosynthesis tools for the design of de novo synthetic pathways, *Current Opinion in Chemical Biology* 28, 99-104.

- [9] Henry, C. S., Broadbelt, L. J., and Hatzimanikatis, V. (2010) Discovery and Analysis of Novel Metabolic Pathways for the Biosynthesis of Industrial Chemicals: 3-Hydroxypropanoate, *Biotechnology and Bioengineering* 106, 462-473.
- [10] 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 Research* 38, W138-W143.
- [11] Cho, A., Yun, H., Park, J. H., Lee, S. Y., and Park, S. (2010) Prediction of novel synthetic pathways for the production of desired chemicals, *BMC Systems Biology* 4, 35.
- [12] Hou, B. K., Ellis, L. B. M., and Wackett, L. P. (2004) Encoding microbial metabolic logic: predicting biodegradation, *Journal of Industrial Microbiology and Biotechnology* 31, 261-272.
- [13] Ellis, L. B. M., Gao, J., Fenner, K., and Wackett, L. P. (2008) The University of Minnesota pathway prediction system: predicting metabolic logic, *Nucleic Acids Research* 36, W427-W432.
- [14] Campodonico, M. A., Andrews, B. A., Asenjo, J. A., Palsson, B. O., and Feist, A. M. (2014) Generation of an atlas for commodity chemical production in *Escherichia coli* and a novel pathway prediction algorithm, GEM-Path, *Metabolic Engineering* 25, 140-158.
- [15] Rodrigo, G., Carrera, J., Prather, K. J., and Jaramillo, A. (2008) DESHARKY: automatic design of metabolic pathways for optimal cell growth, *Bioinformatics* 24, 2554-2556.
- [16] Carbonell, P., Parutto, P., Herisson, J., Pandit, S. B., and Faulon, J. L. (2014) XTMS: pathway design in an eXTended metabolic space, *Nucleic Acids Research* 42, W389-W394.
- [17] Heath, A. P., Bennett, G. N., and Kavraki, L. E. (2010) Finding metabolic pathways using atom tracking, *Bioinformatics* 26, 1548-1555.
- [18] Dale, J. M., Popescu, L., and Karp, P. D. (2010) Machine learning methods for metabolic pathway prediction, *Bmc Bioinformatics* 11.
- [19] Prather, K. L. J., and Martin, C. H. (2008) De novo biosynthetic pathways: rational design of microbial chemical factories, *Current Opinion in Biotechnology* 19, 468-474.
- [20] Hatzimanikatis, V., Li, C., Ionita, J. A., Henry, C. S., Jankowski, M. D., and Broadbelt, L. J. (2005) Exploring the diversity of complex metabolic networks, *Bioinformatics* 21, 1603-1609.
- [21] Hatzimanikatis, V., Li, C. H., Ionita, J. A., and Broadbelt, L. J. (2004) Metabolic networks: enzyme function and metabolite structure, *Curr Opin Struc Biol* 14, 300-306.
- [22] Hadadi, N., Soh, K. C., Seijo, M., Zisaki, A., Guan, X. L., Wenk, M. R., and Hatzimanikatis, V. (2014) A computational framework for integration of lipidomics data into metabolic pathways, *Metabolic Engineering* 23, 1-8.
- [23] Hadadi, N., Hafner, J., Shajkofci, A., Zisaki, A., and Hatzimanikatis, V. (2016) ATLAS of Biochemistry: A Repository of All Possible Biochemical Reactions for Synthetic Biology and Metabolic Engineering Studies, *ACS Synthetic Biology*, 1155-1166.
- [24] Soh, K. C., and Hatzimanikatis, V. (2010) Dreams of Metabolism, *Trends in Biotechnology* 28, 501-508.

- [25] Brunk, E., Neri, M., Tavernelli, I., Hatzimanikatis, V., and Rothlisberger, U. (2012) Integrating computational methods to retrofit enzymes to synthetic pathways, *Biotechnology and Bioengineering* 109, 572-582.
- [26] Hoell, D., Mensing, T., Roggenbuck, R., Sakuth, M., Sperlich, E., Urban, T., Neier, W., and Strehlke, G. (2009) 2-Butanone, In *Ullmann's Encyclopedia of Industrial Chemistry*, Wiley-VCH Verlag GmbH & Co. KGaA.
- [27] Hoppe, F., Burke, U., Thewes, M., Heufer, A., Kremer, F., and Pischinger, S. (2016) Tailor-Made Fuels from Biomass: Potentials of 2-butanone and 2-methylfuran in direct injection spark ignition engines, *Fuel* 167, 106-117.
- [28] Srirangan, K., Liu, X., Akawi, L., Bruder, M., Moo-young, M., and Chou, C. P. (2016) Engineering *Escherichia coli* for Microbial Production of Butanone, 2574-2584.
- [29] Yoneda, H., Tantillo, D. J., and Atsumi, S. (2014) Biological production of 2-butanone in *Escherichia coli*, *ChemSusChem* 7, 92-95.
- [30] Multer, A., McGraw, N., Hohn, K., and Vadlani, P. (2013) Production of Methyl Ethyl Ketone from Biomass Using a Hybrid Biochemical/Catalytic Approach, *Industrial & Engineering Chemistry Research* 52, 56-60.
- [31] Drabo, P., Tiso, T., Heyman, B., Sarikaya, E., Gaspar, P., Förster, J., Büchs, J., Blank, L. M., and Delidovich, I. (2017) Anionic Extraction for Efficient Recovery of Biobased 2,3-Butanediol—A Platform for Bulk and Fine Chemicals, *ChemSusChem* 10, 3252-3259.
- [32] Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M., and Tanabe, M. (2016) KEGG as a reference resource for gene and protein annotation, *Nucleic Acids Research* 44, D457-D462.
- [33] Kanehisa, M., and Goto, S. (2000) KEGG: Kyoto Encyclopedia of Genes and Genomes, *Nucleic Acids Research* 28, 27-30.
- [34] Engel, T. (2007) The structural- and bioassay database PubChem, *Nachr Chem* 55, 521-524.
- [35] Kim, S., Thiessen, P. A., Bolton, E. E., Chen, J., Fu, G., Gindulyte, A., Han, L., He, J., He, S., Shoemaker, B. A., Wang, J., Yu, B., Zhang, J., and Bryant, S. H. (2016) PubChem Substance and Compound databases, *Nucleic Acids Research* 44, D1202-D1213.
- [36] Krumpfer, J. W., Giebel, E., Frank, E., Müller, A., Ackermann, L.-M., Tironi, C. N., Mourgas, G., Unold, J., Klapper, M., Buchmeiser, M. R., and Müllen, K. (2017) Poly(Methyl Vinyl Ketone) as a Potential Carbon Fiber Precursor, *Chemistry of Materials* 29, 780-788.
- [37] Siegel, H., and Eggersdorfer, M. (2000) Ketones, In *Ullmann's Encyclopedia of Industrial Chemistry*, Wiley-VCH Verlag GmbH & Co. KGaA.
- [38] Eller, K., Henkes, E., Rossbacher, R., and Höke, H. (2000) Amines, Aliphatic, In *Ullmann's Encyclopedia of Industrial Chemistry*, Wiley-VCH Verlag GmbH & Co. KGaA.
- [39] Chaparro-Riggers, J. F., Rogers, T. A., Vazquez-Figueroa, E., Polizzi, K. M., and Bommarius, A. S. (2007) Comparison of three enoate reductases and their potential use for biotransformations, *Adv Synth Catal* 349, 1521-1531.
- [40] Lilley, D. M. J., Clegg, R. M., Diekmann, S., Seeman, N. C., Von Kitzing, E., and Hagerman, P. J. (1995) A nomenclature of junctions and branchpoints in nucleic acids, *Nucleic Acids Research* 23, 3363-3364.
- [41] Orth, J. D., Conrad, T. M., Na, J., Lerman, J. A., Nam, H., Feist, A. M., and Palsson, B. O. (2011) A comprehensive genome-scale reconstruction of *Escherichia coli* metabolism-2011, *Molecular Systems Biology* 7.

- [42] Ataman, M., and Hatzimanikatis, V. (2015) Heading in the right direction: thermodynamics-based network analysis and pathway engineering, *Current Opinion in Biotechnology* 36, 176-182.
- [43] Islam, M. A., Hadadi, N., Ataman, M., Hatzimanikatis, V., and Stephanopoulos, G. (2017) Exploring biochemical pathways for mono-ethylene glycol (MEG) synthesis from synthesis gas, *Metabolic Engineering* 41, 173-181.
- [44] Soh, K. C., and Hatzimanikatis, V. (2010) Network thermodynamics in the post-genomic era, *Current Opinion in Microbiology* 13, 350-357.
- [45] Soh, K. S., and Hatzimanikatis, V. (2014) Constraining the flux space using thermodynamics and integration of metabolomics data, *Methods in Molecular Biology* 1191, 49-63.
- [46] Hadadi, N., MohamadiPeyhani, H., Miskovic, L., Seijo, M., and Hatzimanikatis, V. (2017) Knowledge of the Neighborhood of the Reactive Site up to Three Atoms Can Predict Biochemistry and Protein Sequences, *bioRxiv*, <https://doi.org/10.1101/210039>.
- [47] Neidhardt, F. C., Ingraham, J. L., and Schaechter, M. (1990) *Physiology of the bacterial cell : a molecular approach*, Sinauer Associates, Sunderland, Mass.
- [48] Sudarsan, S., Dethlefsen, S., Blank, L. M., Siemann-Herzberg, M., and Schmid, A. (2014) The Functional Structure of Central Carbon Metabolism in *Pseudomonas putida* KT2440, *Applied and Environmental Microbiology* 80, 5292-5303.
- [49] Ataman, M., and Hatzimanikatis, V. (2017) lumpGEM: Systematic generation of subnetworks and elementally balanced lumped reactions for the biosynthesis of target metabolites, *PLOS Computational Biology* 13, e1005513.
- [50] Chen, Y., Daviet, L., Schalk, M., Siewers, V., and Nielsen, J. (2013) Establishing a platform cell factory through engineering of yeast acetyl-CoA metabolism, *Metab Eng* 15, 48-54.
- [51] Haller, T., Buckel, T., Rétey, J., and Gerlt, J. A. (2000) Discovering New Enzymes and Metabolic Pathways: Conversion of Succinate to Propionate by *Escherichia coli*, *Biochemistry* 39, 4622-4629.
- [52] Petzold, C. J., Chan, L. J. G., Nhan, M., and Adams, P. D. (2015) Analytics for Metabolic Engineering, *Frontiers in Bioengineering and Biotechnology* 3, 135.
- [53] Campbell, K., Xia, J., and Nielsen, J. The Impact of Systems Biology on Bioprocessing, *Trends Biotechnol* 35, 1156-1168.
- [54] Miskovic, L., Alff-Tuomala, S., Soh, K. C., Barth, D., Salusjarvi, L., Pitkanen, J. P., Ruohonen, L., Penttila, M., and Hatzimanikatis, V. (2017) A design-build-test cycle using modeling and experiments reveals interdependencies between upper glycolysis and xylose uptake in recombinant *S. cerevisiae* and improves predictive capabilities of large-scale kinetic models, *Biotechnology for Biofuels* 10.
- [55] Bordbar, A., Monk, J. M., King, Z. A., and Palsson, B. O. (2014) Constraint-based models predict metabolic and associated cellular functions, *Nat Rev Genet* 15, 107-120.
- [56] Maarleveld, T. R., Khandelwal, R. A., Olivier, B. G., Teusink, B., and Bruggeman, F. J. (2013) Basic concepts and principles of stoichiometric modeling of metabolic networks, *Biotechnol J* 8, 997-U952.
- [57] Garcia-Albornoz, M. A., and Nielsen, J. (2013) Application of Genome-Scale Metabolic Models in Metabolic Engineering, *Industrial Biotechnology* 9, 203-214.

- [58] Miskovic, L., Tokic, M., Fengos, G., and Hatzimanikatis, V. (2015) Rites of passage: requirements and standards for building kinetic models of metabolic phenotypes, *Current Opinion in Biotechnology* 36, 146-153.
- [59] Miskovic, L., and Hatzimanikatis, V. (2010) Production of biofuels and biochemicals: in need of an ORACLE, *Trends Biotechnol* 28, 391-397.
- [60] Andreozzi, S., Chakrabarti, A., Soh, K. C., Burgard, A., Yang, T. H., Van Dien, S., Miskovic, L., and Hatzimanikatis, V. (2016) Identification of metabolic engineering targets for the enhancement of 1,4-butanediol production in recombinant E. coli using large-scale kinetic models, *Metabolic Engineering* 35, 148-159.
- [61] Chakrabarti, A., Miskovic, L., Soh, K. C., and Hatzimanikatis, V. (2013) Towards kinetic modeling of genome-scale metabolic networks without sacrificing stoichiometric, thermodynamic and physiological constraints, *Biotechnology journal* 8, 1043-1057.
- [62] Savoglidis, G., dos Santos, A. X. D., Riezman, I., Angelino, P., Riezman, H., and Hatzimanikatis, V. (2016) A method for analysis and design of metabolism using metabolomics data and kinetic models: Application on lipidomics using a novel kinetic model of sphingolipid metabolism, *Metabolic Engineering* 37, 46-62.
- [63] Stanford, N. J., Lubitz, T., Smallbone, K., Klipp, E., Mendes, P., and Liebermeister, W. (2013) Systematic Construction of Kinetic Models from Genome-Scale Metabolic Networks, *PLOS One* 8.
- [64] Dash, S., Khodayari, A., Zhou, J., Holwerda, E. K., Olson, D. G., Lynd, L. R., and Maranas, C. D. (2017) Development of a core Clostridium thermocellum kinetic metabolic model consistent with multiple genetic perturbations, *Biotechnology for Biofuels* 10.
- [65] Lee, Y., Lafontaine Rivera, J. G., and Liao, J. C. (2014) Ensemble Modeling for Robustness Analysis in engineering non-native metabolic pathways, *Metabolic Engineering* 25, 63-71.
- [66] Esvelt, K. M., and Wang, H. H. (2013) Genome-scale engineering for systems and synthetic biology, *Molecular Systems Biology* 9.
- [67] Barrangou, R., and Doudna, J. A. (2016) Applications of CRISPR technologies in research and beyond, *Nature Biotechnology* 34, 933-941.
- [68] Bochner, B. R. (2009) Global phenotypic characterization of bacteria, *Fems Microbiol Rev* 33, 191-205.
- [69] Li, C. H., Henry, C. S., Jankowski, M. D., Ionita, J. A., Hatzimanikatis, V., and Broadbelt, L. J. (2004) Computational discovery of biochemical routes to specialty chemicals, *Chemical Engineering Science* 59, 5051-5060.
- [70] Corey, E. J. (1991) The Logic of Chemical Synthesis - Multistep Synthesis of Complex Carbogenic Molecules, *Angew Chem Int Edit* 30, 455-465.
- [71] Orth, J. D., Thiele, I., and Palsson, B. Ø. (2010) What is flux balance analysis?, *Nature biotechnology* 28, 245-248.
- [72] Henry, C. S., Broadbelt, L. J., and Hatzimanikatis, V. (2007) Thermodynamics-based metabolic flux analysis, *Biophys J* 92, 1792-1805.
- [73] Henry, C. S., Jankowski, M. D., Broadbelt, L. J., and Hatzimanikatis, V. (2006) Genome-scale thermodynamic analysis of Escherichia coli metabolism, *Biophys J* 90, 1453-1461.
- [74] Frainay, C., and Jourdan, F. (2017) Computational methods to identify metabolic sub-networks based on metabolomic profiles, *Briefings in Bioinformatics* 18, 43-56.

- [75] James, C. A., and Weininger, D. Daylight Theory Manual, *Daylight Chemical Information Systems, Inc.: Irvine, CA*.
- [76] Ataman, M., Gardiol, D. H. F., Fengos, G., and Hatzimanikatis, V. (2017) redGEM: Systematic Reduction and Analysis of Genome-scale Metabolic Reconstructions for Development of Consistent Core Metabolic Models, *PLOS Computational Biology*, e1005444.
- [77] Giri, V., Sivakumar, T. V., Cho, K. M., Kim, T. Y., and Bhaduri, A. (2015) RxnSim: a tool to compare biochemical reactions, *Bioinformatics* 31, 3712-3714.
- [78] Bajusz, D., Rácz, A., and Héberger, K. (2015) Why is Tanimoto index an appropriate choice for fingerprint-based similarity calculations?, *Journal of Cheminformatics* 7, 20.

Figures

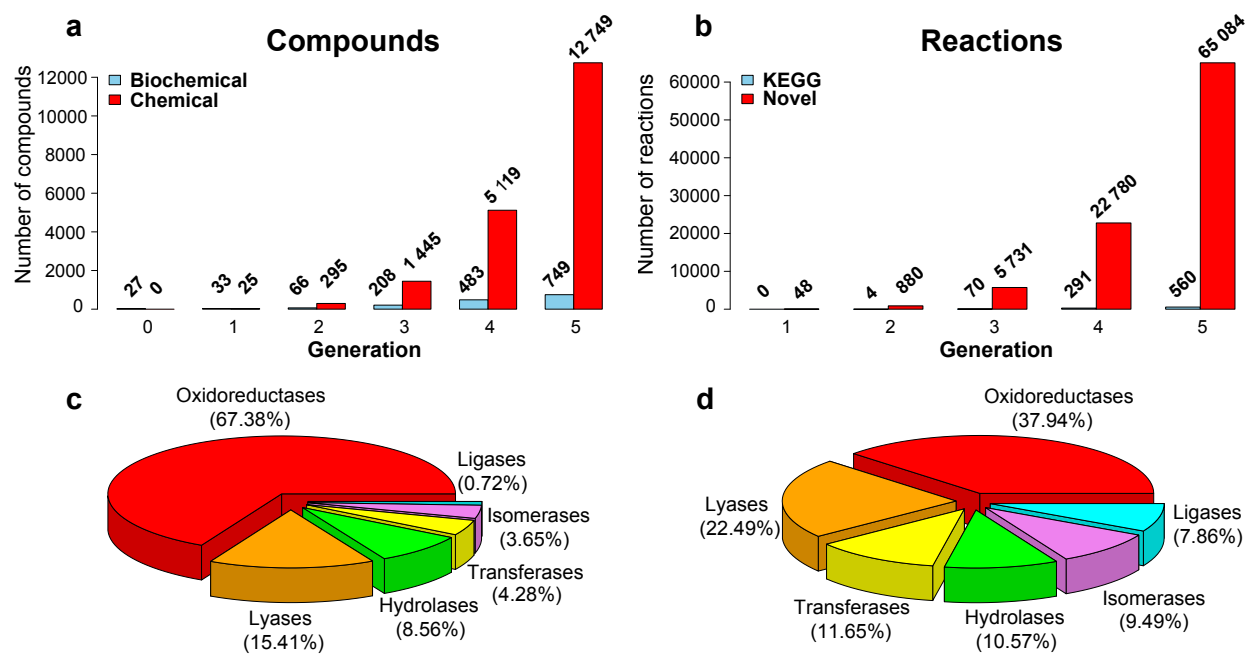
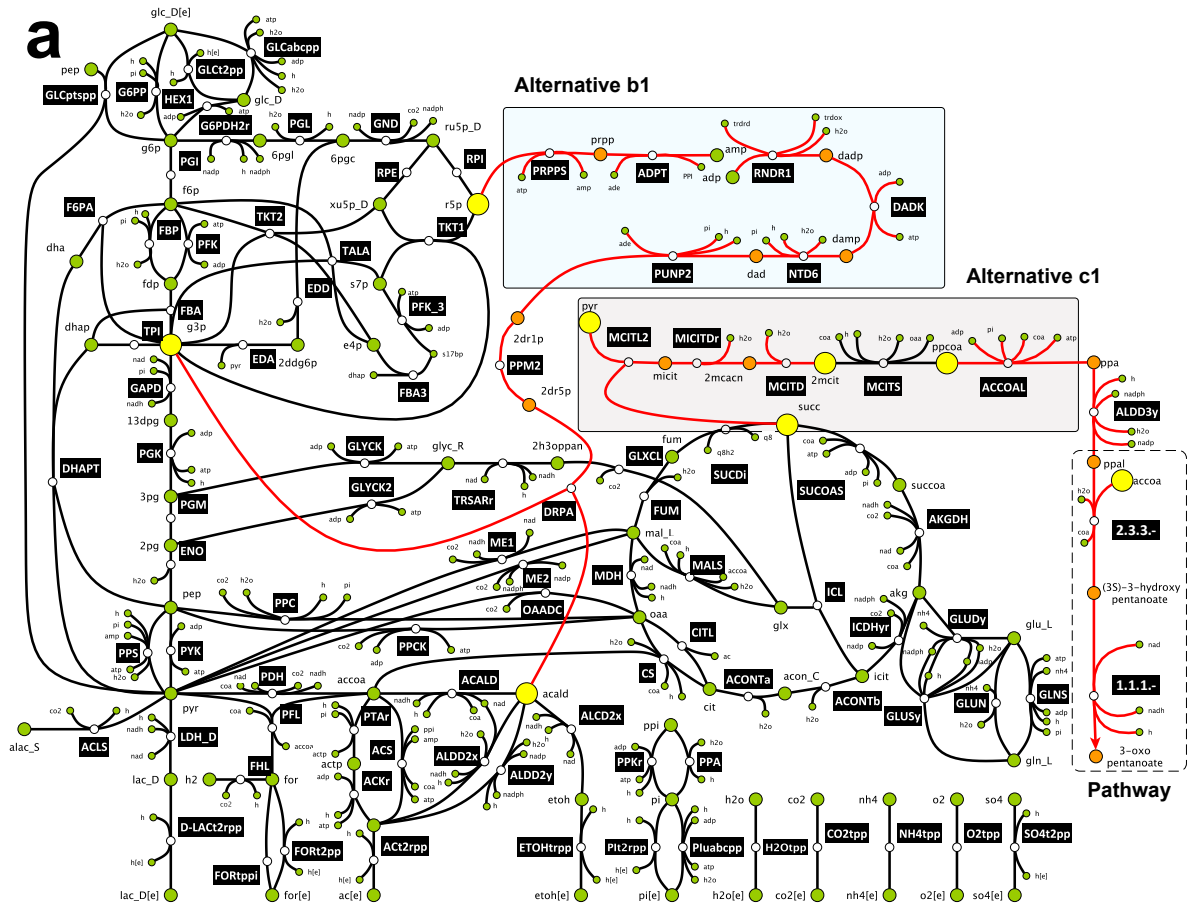
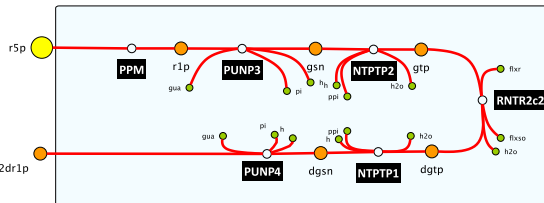


Figure 1

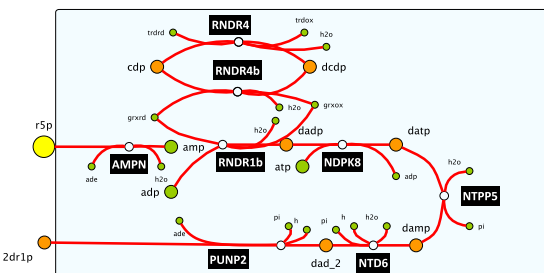


b

Alternative b2

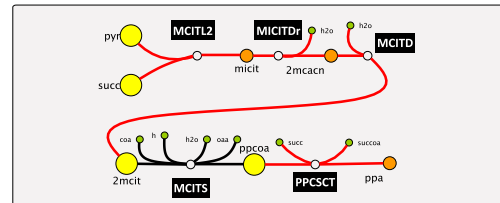


Alternative b102



c

Alternative c2



Alternative c9

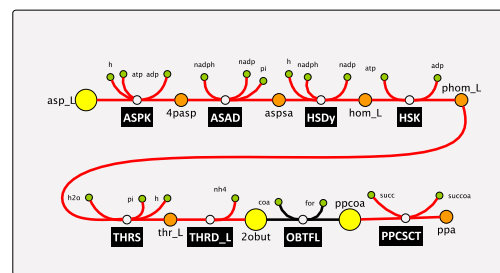


Figure 2

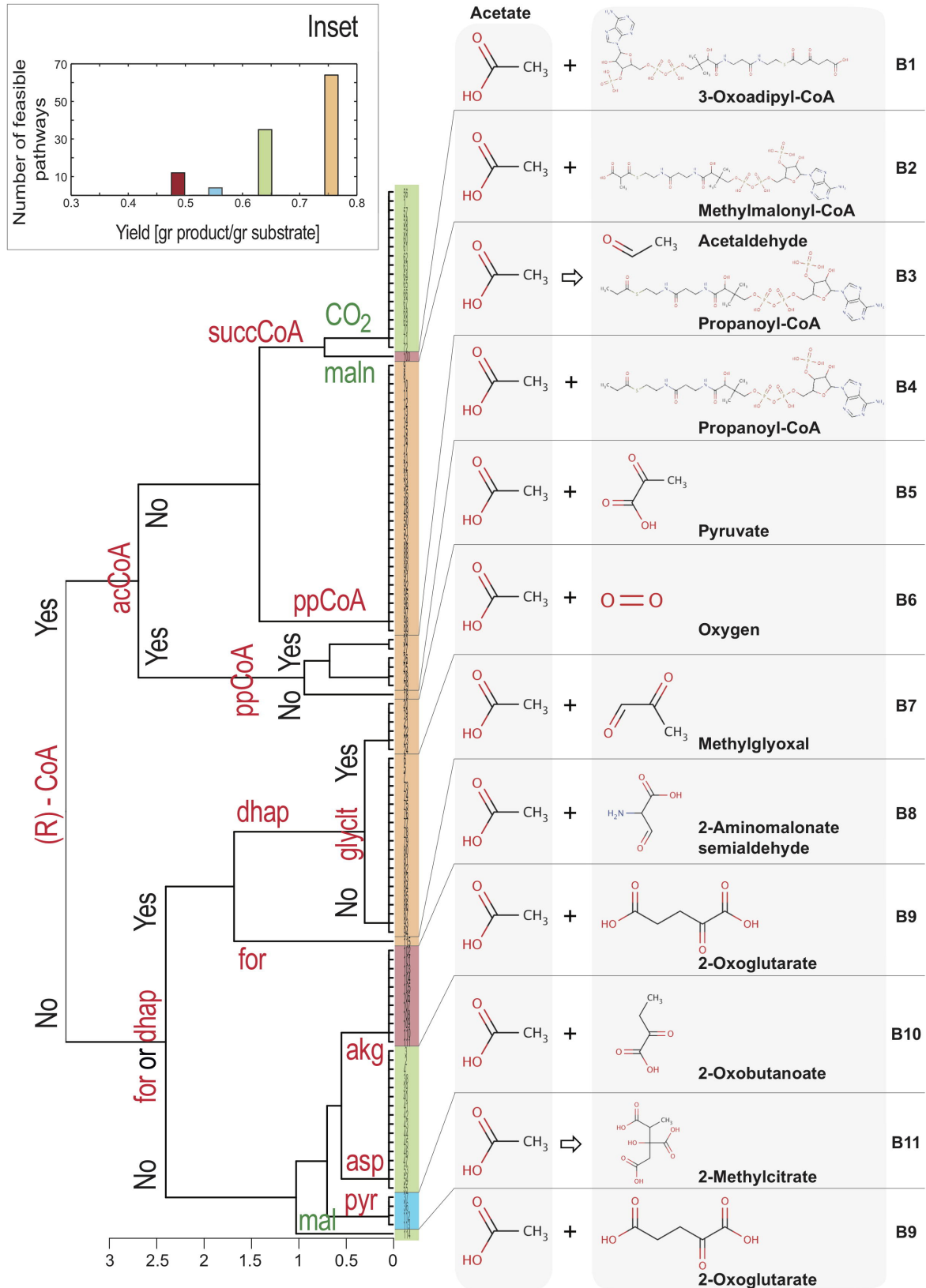
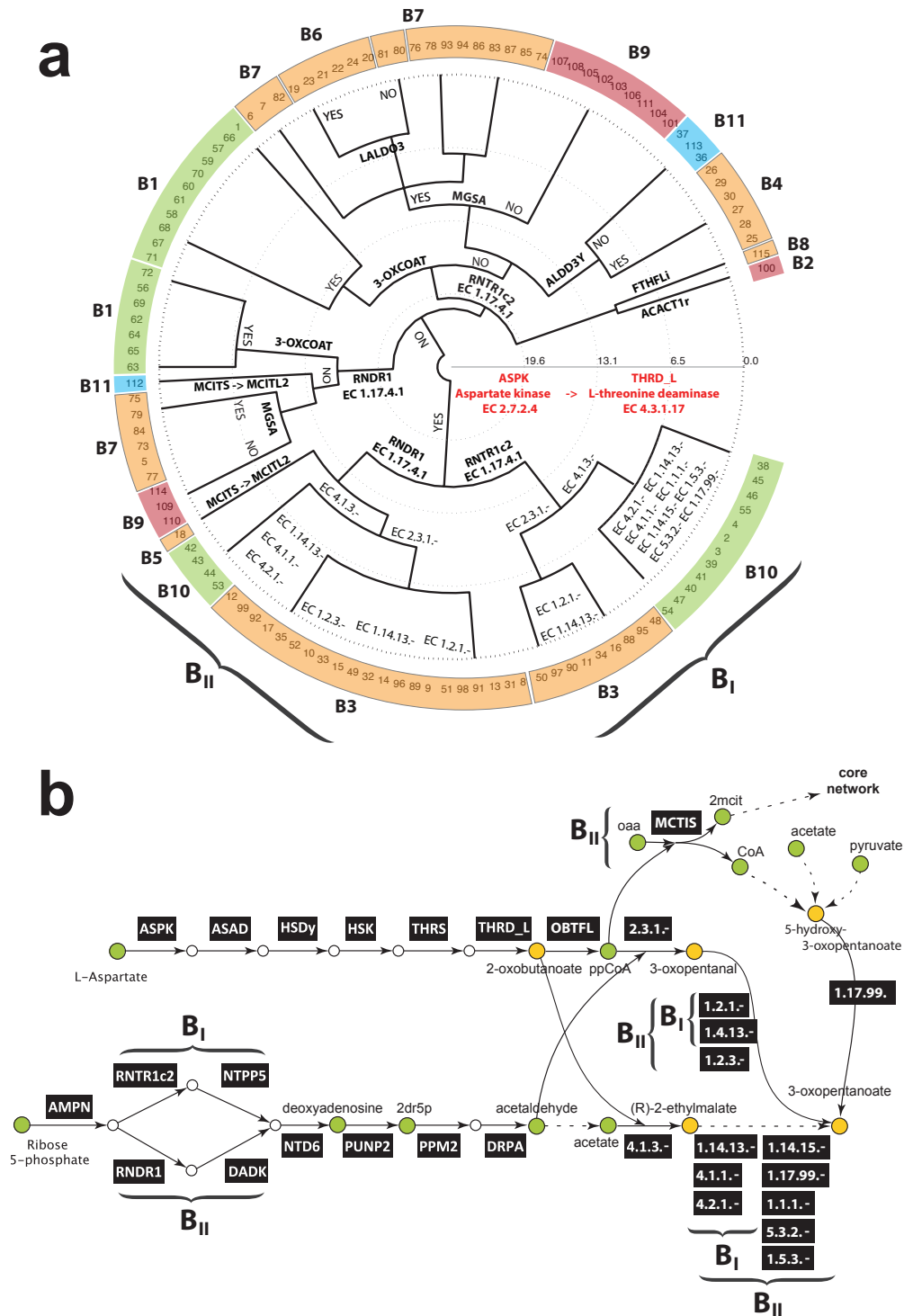


Figure 3



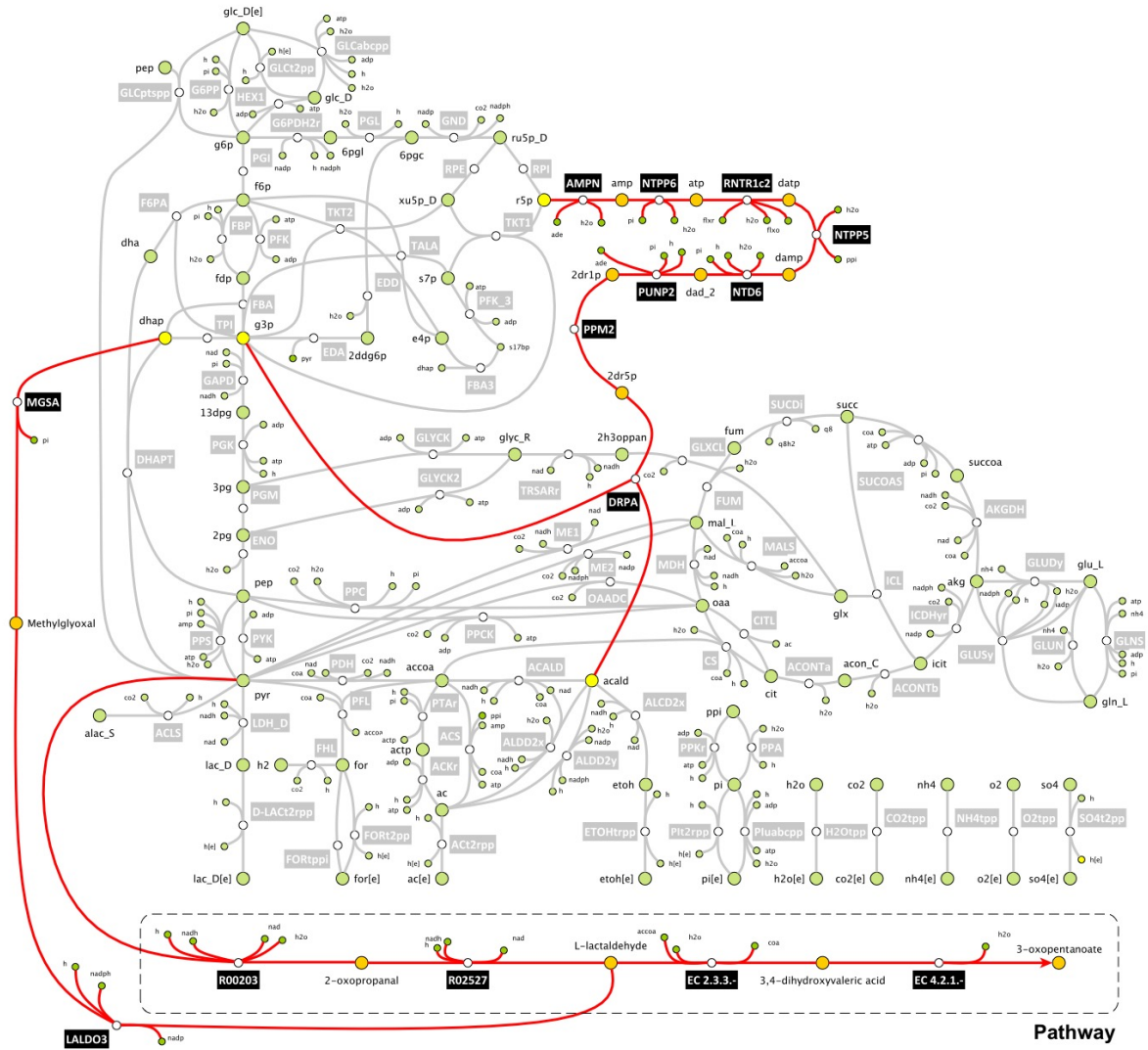


Figure 5

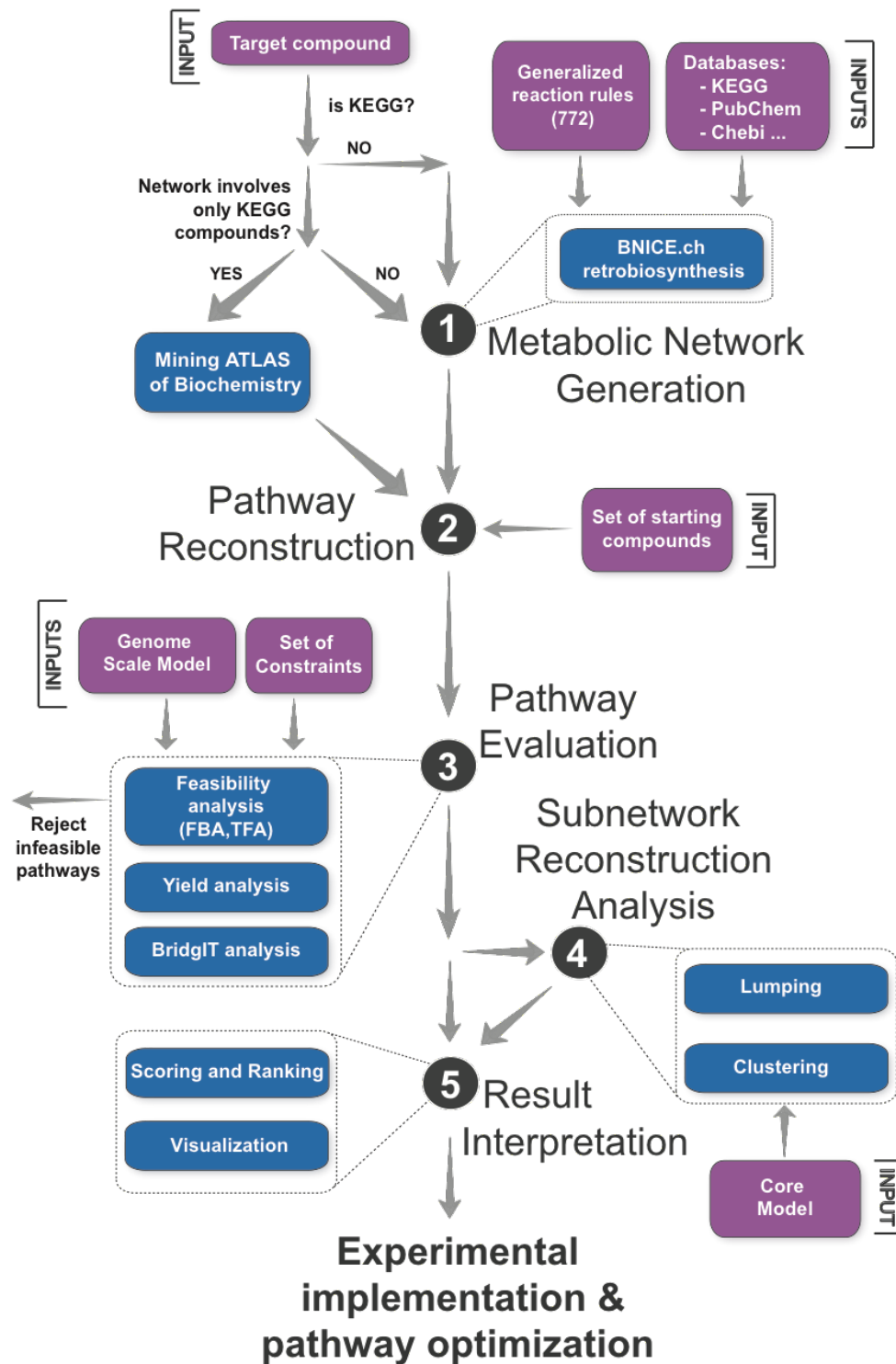


Figure 6