

## ARTICLE

# Establishing synthesis pathway-host compatibility via enzyme solubility

Sara A. Amin<sup>1</sup> | Venkatesh Endalur Gopinarayanan<sup>2</sup> | Nikhil U. Nair<sup>2</sup>  |  
Soha Hassoun<sup>1,2</sup> 

<sup>1</sup>Department of Computer Science, Tufts University, Medford, Massachusetts

<sup>2</sup>Department of Chemical and Biological Engineering, Tufts University, Medford, Massachusetts

**Correspondence**

Nikhil U. Nair, Department of Chemical and Biological Engineering, Science and Technology Center #276, 4 Colby St, Medford, MA 02155.

Email: nikhil.nair@tufts.edu

Soha Hassoun, Department of Computer Science, Tufts University, Halligan Hall #232, 161 College Avenue, Medford, MA 02155.

Email: soha@cs.tufts.edu

**Funding information**

National Science Foundation, Grant/Award Number: CCF-1421972; Tufts University, Grant/Award Number: Tufts Collaborates; National Institutes of Health, Grant/Award Number: DP2HD091798

**Abstract**

Current pathway synthesis tools identify possible pathways that can be added to a host to produce the desired target molecule through the exploration of abstract metabolic and reaction network space. However, not many of these tools explore gene-level information required to physically realize the identified synthesis pathways, and none explore enzyme-host compatibility. Developing tools that address this disconnect between abstract reactions/metabolic design space and physical genetic sequence design space will enable expedited experimental efforts that avoid exploring unprofitable synthesis pathways. This work describes a workflow, termed Probabilistic Pathway Assembly with Solubility Confidence Scores (*ProPASS*), which links synthesis pathway construction with the exploration of the physical design space as imposed by the availability of enzymes with predicted characterized activities within the host. Predicted protein solubility propensity scores are used as a confidence level to quantify the compatibility of each pathway enzyme with the host *Escherichia coli* (*E. coli*). This study also presents a database, termed Protein Solubility Database (*ProSol DB*), which provides solubility confidence scores in *E. coli* for 240,016 characterized enzymes obtained from *UniProtKB/Swiss-Prot*. The utility of *ProPASS* is demonstrated by generating genetic implementations of heterologous synthesis pathways in *E. coli* that target several commercially useful biomolecules.

**KEYWORDS**

metabolic engineering, pathway design, pathway implementation, solubility, synthesis pathway, synthetic biology

## 1 | INTRODUCTION

Synthesis pathways have been engineered within microbial hosts to produce commercially useful biomolecules including polyesters (Fidler & Dennis, 1992), building blocks for industrial polymers (Tong, Liao, & Cameron, 1991), biofuels (Nawabi, Bauer, Kyrpides, & Lykidis, 2011; Radakovits, Jinkerson, Darzins, & Posewitz, 2010; Steen et al., 2010), and therapeutic natural products derived from isoprenoids (Martin, Pitera, Withers, Newman, & Keasling, 2003; Pitera, Paddon, Newman, & Keasling, 2007; Watts, Mijts, & Schmidt-Dannert, 2005), polyketides (Peirú, Menzella, Rodríguez, Carney, &

Gramajo, 2005; Pfeifer, Admiraal, Gramajo, Cane, & Khosla, 2001), and nonribosomal peptides (Takahashi et al., 2007). Synthesis pathways consist of a series of nonnative enzyme-controlled reactions from metabolites within the host to a target molecule. Often, many possible biochemical routes to a target molecule exist such as in the production of 3-hydroxypropionic acid through glycerol,  $\beta$ -alanine, or acrolein pathways (Cheng, Jiang, Wu, Li, & Ye, 2016; Luo et al., 2016; Song, Kim, Cho, & Lee, 2016). Further, the genetic implementation of each pathway is not unique because genes can be sourced from various organisms. For example, production of thebaine and hydrocodone in *Saccharomyces cerevisiae* required

construction of a heterologous pathway with enzymes from a bacterium (*Pseudomonas putida*), a mammal (*Rattus norvegicus*), and several plants (*Papaver somniferum*, *Papaver bracteatum*, *Coptis japonica*, *Eschscholzia californica*) (Galanie, Thodey, Trenchard, Interante, & Smolke, 2015).

Needless to state, experimental efforts to explore all possible pathways and all possible organism-specific enzyme selections are prohibitive. Unfortunately, current pathway synthesis computational tools provide solutions based on exploring the abstract metabolic space without regards to possible genetic implementation. Identifying best genetic implementations are typically ad hoc, and there is a lack of systematic design tools and workflows that support the coexploration of the metabolic and genetic implementation design spaces. Pathway synthesis tools that consider genetic compatibility between biosynthetic pathways and heterologous hosts can significantly expedite the metabolic engineering cycle by focusing on design spaces with a greater fraction of profitable options.

In this light, pathway synthesis tools must conceptually perform two distinct tasks. The first task involves identifying possible synthesis pathways from the host to a target molecule and evaluating potential yield. This is known as the pathway construction or pathway identification problem, where sequences of reactions that synthesize the target molecule from a host metabolite are selected based on data from multi-organism databases such as Kyoto Encyclopedia of Genes and Genomes (KEGG; Kanehisa & Goto, 2000), MetaCyc (Caspi et al., 2008), and SEED (Overbeek et al., 2005). Pathway construction utilizes either rule-based techniques (e.g. BNICE; Wu, Wang, Assary, Broadbelt, & Krilov, 2011), PathPred (Moriya et al., 2010), Meta (Klopman, Dimayuga, & Talafous, 1994; Klopman, Tu, & Talafous, 1997), Meteor (Greene, Judson, Langowski, & Marchant, 1999; Marchant, Briggs, & Long, 2008), and UM-PPS (Ellis & Wackett, 2012; Hou, Wackett, & Ellis, 2003) or graph-based techniques such as *ProPath* (Yousofshahi, Lee, & Hassoun, 2011), DESHARKY (Rodrigo, Carrera, Prather, & Jaramillo, 2008), Metaroute (Blum & Kohlbacher, 2008). Thus, there are currently multiple tools available for pathway construction.

The second synthesis task concerns the parts selection problem or, more broadly, the pathway implementation problem. This task involves identifying gene sequences from specific organisms that best realize reactions along the synthesis pathway. There are several factors to consider. One set of factors is associated with transcription and translation, including promoter strength and efficiency of ribosome binding. Another set of factors is dependent on the interactions of the enzyme's protein coding sequence with the cellular metabolic machinery to ensure correct conformation and function while avoiding misfolding or aggregation, known as protein solubility. Overcoming protein solubility issues can be achieved by codon optimization, coexpressing molecular chaperones (Tréaugues et al., 2004; Xia et al., 2016), lowering culture temperature (Makrides, 1996) or modifying growth media (Makrides, 1996). These strategies can improve solubility by finding optimal refolding conditions, but may also lead to activity loss of a protein (Singh & Panda, 2005). Further, when multiple enzymes in a pathway are insoluble, a single strategy may not be applicable. It is

therefore desirable to select high-solubility enzymes to implement synthesis pathways.

To identify profitable synthesis pathway designs for experimental implementation, we couple in this paper the two synthesis tasks: pathway construction and pathway implementation. Using sequences from *UniProtKB* database (UniProt Consortium., 2017), we assemble a Protein Solubility Data base (*ProSol DB*), that comprises confidence scores for the likelihood of an enzyme being soluble in *Escherichia coli* (*E. coli*) and also allows for quick lookup using enzyme commission (EC) numbers. We used *ccSQL omics* (Agostini, Cirillo, Livi, Ponti, & Tartaglia, 2014) to compute solubility propensity scores for various proteins from *UniProtKB* in *E. coli*. Our database stores solubility confidence scores for 240,016 sequences, of which 34,046 sequences are associated with commonly used organisms. We develop a new workflow, termed Probabilistic Pathway Assembly with Solubility confidence Scores (*ProPASS*), to couple *ProPath* (Yousofshahi et al., 2011), a method for constructing synthesis pathways, and *ProSol DB*. The workflow consists of first identifying synthesis pathways from a host to a target metabolite. Pathways are then ordered increasingly based on their length and also ordered decreasingly based on yield. This ordering step allows speedy identification of short-length and high-yielding pathways, which are desirable for experimental validation. *ProPASS* then recommends sequences based on their solubility confidence scores in *ProSol DB* to implement reactions along each pathway. We apply our workflow to identify implementations of synthesis pathways for seven target molecules. We analyze three test cases in detail. In all three cases, we show that *ProPASS* identifies one or more implementation that is predicted soluble. We further show that implementations published in the literature are recommended by *ProPASS* if cataloged in *UniProtKB*. As far as we know, this is the first description of a systematic workflow or method that explores using gene-specific sequence information to systematically identify host-compatible enzymes that can be used to implement synthesis pathways.

## 2 | METHODS

### 2.1 | Assembling the protein solubility database (*ProSol DB*)

To create a database of solubility confidence scores for various enzymes, we utilize annotated sequences in the *UniProtKB* database (UniProt Consortium., 2017). Protein sequences that have been either experimentally validated or manually analyzed are considered "reviewed" within the *UniProtKB/Swiss-Prot* database. While *UniProtKB/Swiss-Prot* currently contains over 550,000 sequences, only 240,016 sequences are associated with 4,652 EC numbers. *ProSol DB* is designed such that the EC numbers serve as a key for database lookups, and the returned value contains the *UniProtKB/Swiss-Prot* sequence IDs, solubility confidence scores, and associated organisms. Solubility confidence scores are computed only once, using *ccSQL omics*, and stored locally in the database, eliminating the need to predict scores repeatedly, thus speeding up enzyme selection. *ccSQL*

*omics* calculates solubility propensity of fragments within the sequences and utilizes neural networks to estimate confidence in solubility based on Fourier transform coefficients of the sequences (Agostini et al., 2014; Tartaglia, Cavalli, & Vendruscolo, 2007; Tartaglia, Pechmann, Dobson, & Vendruscolo, 2009).

While several solubility prediction tools are available, we calculate solubility confidence scores using *ccSOL omics*. The prediction accuracy of *ccSOL omics* reported is high (74%) for three independent datasets (Agostini et al., 2014). Further, the availability of the *ccSOL omics* code without a web interface facilitated the direct construction of our database. Several recent papers (Chang, Song, Tey, & Ramanan, 2014; Habibi, Hashim, Norouzi, & Samian, 2014; Khurana et al., 2018), provide an assessment of the accuracy of various tools; however, each tool is trained on a different data set. As our goal is to predict the solubility of enzyme sequences from a broad range of organisms reported in *UniProtKB/Swiss-Prot*, we identified a subset of sequences in *UniProtKB/Swiss-Prot* for which experimental solubility labels are readily available. The labels were collected by comparing sequences from *UniProtKB/Swiss-Prot* to those in Target Track DB (Chen, Oughtred, Berman, & Westbrook, 2004). The total number of sequences in our independent data set, referred to as “the enzyme solubility test set”, consisted of 716 sequences, where 462 are soluble and 254 are insoluble. For the enzyme solubility test set, the accuracy of *ccSOL omics* was 73.46% using a 30% confidence threshold, compared with 46% for DeepSol S1, a recent tool for predicting solubility (Khurana et al., 2018). Supporting Information File 6 provides detailed accuracy results at various confidence thresholds for *ccSOL omics* and for DeepSol S1. While relatively small, the enzyme solubility test set is the most representative sample of

*UniProtKB/Swiss-Prot*. The high prediction accuracy using *ccSOL omics* justifies its use when computing solubility for *ProSol DB*. Further, based on the analysis of the enzyme solubility test, we suggest using a solubility confidence score equal to or greater than 30% as “high confidence of solubility”.

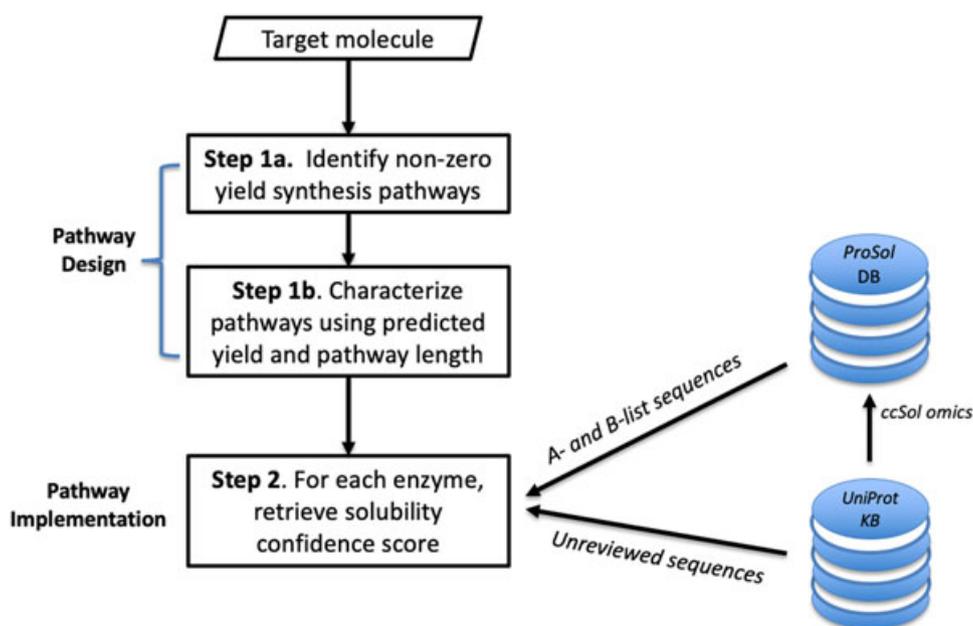
## 2.2 | ProPASS: Coupling probabilistic pathway construction with protein solubility prediction

With the ability to lookup protein solubility confidence scores in *ProSol DB* based on EC numbers, we developed *ProPASS*, a workflow that couples pathway construction with the identification of organism-specific enzymes for each step. The workflow (Figure 1) consists of a design step followed by an implementation step. Additional details can be found in Supporting Information File 1.

Step 1 – Synthesis pathway design.

Step1a. Synthesis pathway identification. Given a target molecule, *ProPath* synthesizes pathways. *ProPath*, a graph-based probabilistic search algorithm, selects synthesis pathways, from a target molecule terminating at a molecule within the host, using all reactions extracted from the KEGG database. We selected *ProPath* for the pathway construction as it is computationally efficient and was shown effective in generating synthesis pathways with yield distributions similar to those obtained using limited-in-depth exhaustive search. In addition, *ProPath* reproduces validated pathways published in the literature.

Step 1b. Synthesis pathway ordering. Identified pathways are ordered by two metrics: pathway yield and pathway length. These metrics are reported to users to allow the examination of different



**FIGURE 1** An overview of the steps used by *ProPASS* to identify synthesis pathways and their enzymatic implementations for a given target molecule. Solubility confidence scores from *ProSol DB* are retrieved for all sequences associated with a specific enzyme from the A-list or B-list organisms or computed on-the-fly using *ccSOL omics* for non-reviewed sequences in *UniProtKB*. *ProPASS*: Probabilistic Pathway Assembly with Solubility confidence Scores; *ProSol DB*: Protein Solubility Database [Color figure can be viewed at wileyonlinelibrary.com]

designs and their tradeoffs when identifying suitable pathway implementations.

**Step 2 – Synthesis pathway implementation.** Given a synthesis pathway, the EC number for each of its enzymatic reaction is utilized to look up the associated *UniProtKB/Swiss-Prot* protein IDs in *ProSol DB*. The retrieved data also contains solubility confidence scores and organisms associated with each protein ID. The scores provide a confidence level that determines if a protein is soluble in *E. coli*. The scores should not be taken as any measure of relative solubility of a given enzyme in *E. coli*. In the case of nonmatching protein IDs in *ProSol DB*, we allow the option to retrieve nonreviewed sequences from *UniProtKB*, and solubility confidence score for each sequence is computed using *ccSOL omics*. If available, each reaction will have one or more recommended implementation sequence, along with its solubility confidence score and source organism.

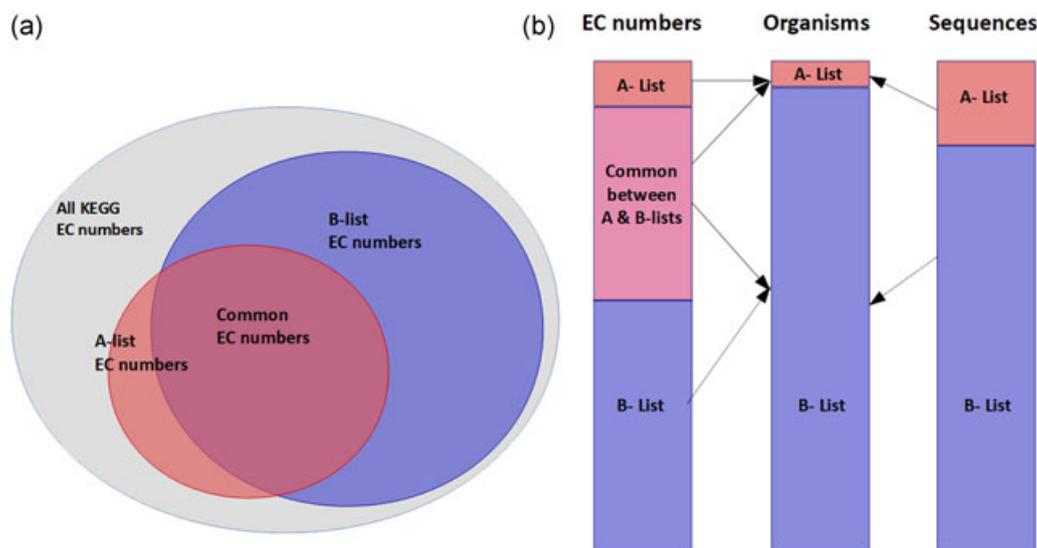
### 3 | RESULTS

#### 3.1 | Overview of Protein Solubility Database (*ProSol DB*)

*ProSol DB* is considered the primary source of solubility confidence scores for *ProPASS* and consists of 240,016 sequences. It is built using all reviewed protein sequences in *UniProtKB/Swiss-Prot*. The protein sequences in *ProSol DB* are associated with 4,652 EC numbers out of 6,896 EC numbers reported in KEGG (Figure 2a). We divide protein sequences organisms into two sets. The first set consists of 20 commonly used and well-studied organisms (e.g., *E. coli*, *S. cerevisiae*,

*B. subtilis*, etc.) from which protein sequences can be easily acquired for experimental validation (A-list; Table S1). The second set, the B-list, contains all other organisms associated with the remaining protein sequences associated with EC numbers. Protein sequences associated with the A-list organisms cover 2,427, or 52%, of all unique EC numbers available in *ProSol DB*, and 13% of all sequences in the database even though it consists of only 0.33% of the organismal diversity (Figure 2b). Example of solubility confidence scores present in our database can be found in Table 1. The first column lists the EC number. The second column shows the associated *UniProtKB/Swiss-Prot* IDs for protein sequences associated with the EC number in the first column. The sequences are identified by their *UniProtKB/Swiss-Prot* IDs to make it easier to obtain any information related to these sequences from *UniProtKB/Swiss-Prot*. The third column contains a solubility confidence score for each protein sequence. The fourth and fifth columns show organisms associated with each protein sequence. When an EC number is queried, a list of protein IDs associated with the A-list organisms is returned. If none are associated with the A-list, a list of protein IDs associated with the B-list is returned. Further examples of solubility confidence scores within *ProSol DB* can be found in Supporting Information File 3.

We analyzed correlations of confidence scores across EC classes and sourcing organisms. If such correlations existed, they could guide pathway synthesis and implementation algorithms by favoring ECs or organisms with an increased propensity for soluble proteins. Since *ProSol DB* is the largest database of solubility confidence scores in *E. coli*, we used it as a resource to investigate any potential correlation between the propensity for soluble expression and EC



**FIGURE 2** Statistics related to *ProSol DB*. (a) Venn diagram showing sets of all EC numbers in KEGG (6,896), those covered by the A-list (2,427), and those covered by the B-list (4,350). The intersection between the A-list and the B-list consists of 2,125 EC numbers covering 46% of EC numbers available in *ProSol DB*. (b) Each stack shows information related to the A-list and B-list organisms. The first bar shows a tally of EC numbers in A-list, B-list, and their intersection. The second bar shows that the A-list consists of 20 organisms, which is only 0.33% of the organism diversity yet covers 52% of all the EC numbers in the database. The third bar shows that 13% of all protein sequences in *ProSol DB* are covered by the A-list organisms. These statistics demonstrate that the number of verified enzyme sequences and EC numbers are overrepresented in the A-list, which comprises many model organisms. EC: enzyme commission; KEGG: Kyoto Encyclopedia of Genes and Genomes; *ProSol DB*: protein solubility database [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

**TABLE 1** Sample from the *ProSol DB* for sequences associated with the acetolactate decarboxylase enzyme

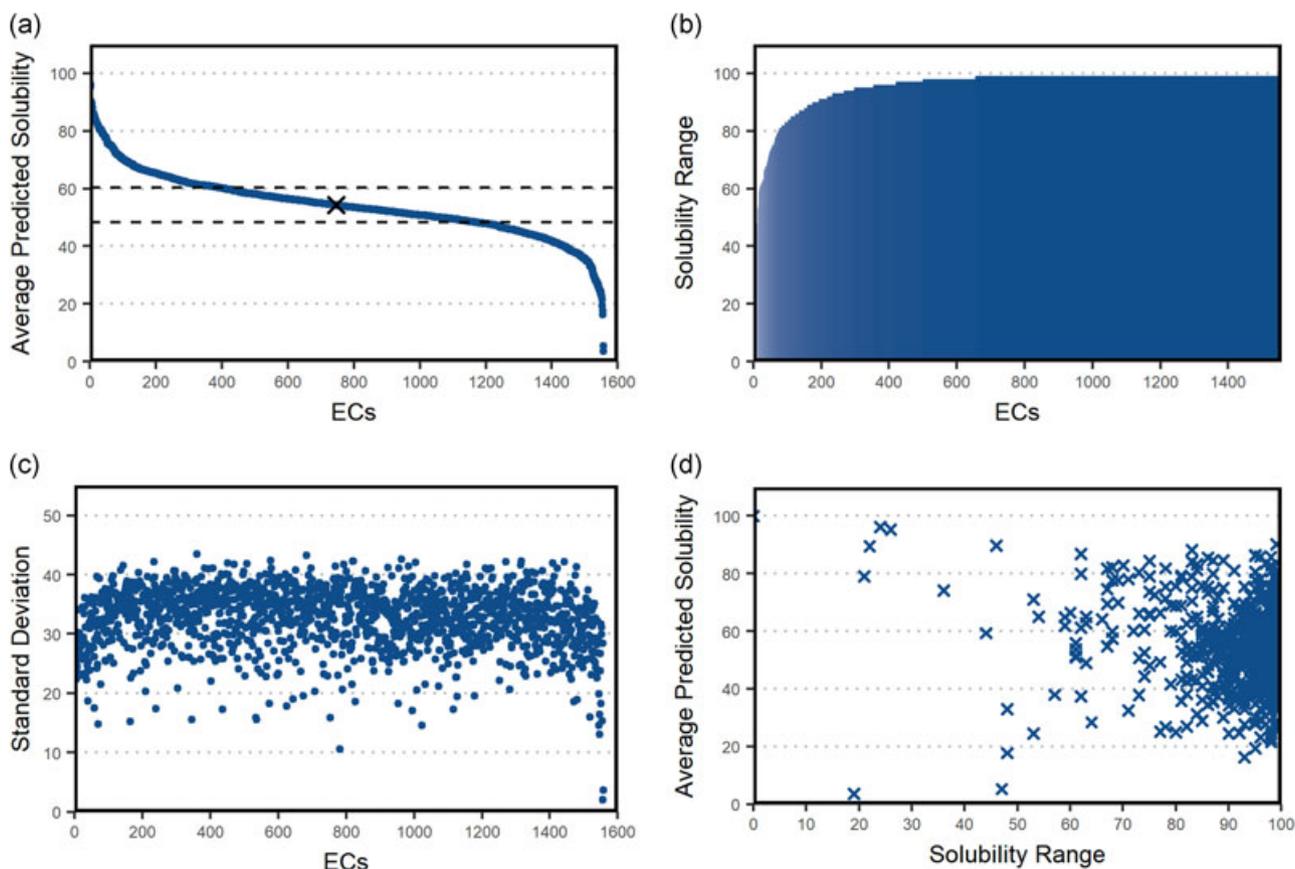
Key		Value		
EC #	UniProt ID	Solubility confidence score	A-List organisms	B-List organisms
4.1.1.5 (acetolactate decarboxylase)	P77880	89	<i>Lactococcus lactis</i>	-
	Q04518	6	N/A	<i>Raoultella terrigena</i>
	Q8PZ55	32	N/A	<i>Methanosarcina mazei</i>
	Q65E52	96	N/A	<i>Bacillus licheniformis</i>
	P95676	83	<i>Lactococcus lactis</i>	-
	Q04777	96	<i>Bacillus subtilis</i>	-
	P23616	34	N/A	<i>Brevibacillus brevis</i>
	Q8L208	68	N/A	<i>Streptococcus thermophilus</i>

Note. N/A, no enzyme with corresponding EC# found in A-list.

-, EC# not queried in B-list since an A-list sequence is found. *ProSol DB*: protein solubility database.

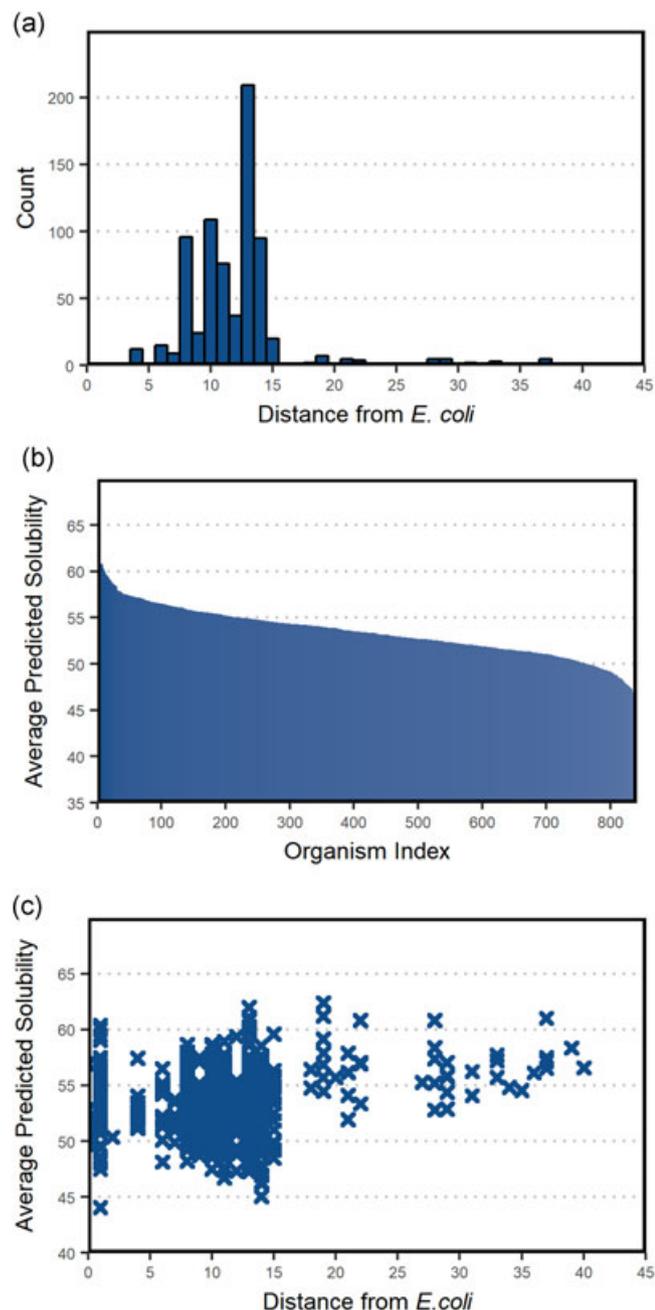
numbers (Figure 3). For this analysis, we only considered EC numbers that are associated with 10 or more protein sequences in *ProSol DB*. The average predicted solubility confidence score for the 1,600 enzymes that meet this criterion is 54.3 (Figure 3a). For each EC number, the range of solubility confidence scores (solubility range) varied tremendously (Figure 3b), with over 96% of ECs having scores

that varied by a difference of 75 or more. Analyzing the standard deviation of the average solubility confidence score per EC number shows that they are also spread over similar wide ranges (Figure 3c). There was a low correlation between the score ranges and average solubility confidence scores (Figure 3d) as analyzed using Spearman's rank test (Spearman's rank coefficient = 0.00897).



**FIGURE 3** Correlation analysis between solubility confidence scores and EC numbers for enzymes with 10 or more sequences in *ProSol DB*. (a) Average solubility confidence scores for enzymes sorted from highest to lowest average score. (b) Predicted solubility confidence score ranges for each enzyme, sorted from smallest to largest. (c) The standard deviation of solubility confidence scores for each enzyme sorted as in panel (a). (d) Scatter plot of solubility confidence score ranges and average scores. There is a low correlation between the score ranges and average solubility confidence scores as determined using Spearman's rank test (Spearman's rank coefficient = 0.00897). EC: enzyme commission; *ProSol DB*: Protein Solubility Database [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

Further, we investigated the correlation between solubility confidence scores and the similarity of the source organisms to *E. coli* (Figure 4). The phylogenetic distance from *E. coli* was computed using phyloT, a tool for the annotation and visualization of phylogenetic trees (Letunic & Bork, 2016). We only considered organisms that had at least 100 solubility confidence scores in



**FIGURE 4** Correlation analysis between predicted solubility scores and sourcing organisms for those with 100 or more sequences in *ProSol* DB. (a) Histogram of phylogenetic distances from *E. coli*. (b) Average predicted solubility per organism, ranked from largest to smallest predicted solubility values. (c) Scatter plot of average predicted solubility per organism versus phylogenetic distance from *E. coli* indicates low positive correlation as calculated using Spearman's rank test (Spearman's rank coefficient = 0.176919). *ProSol* DB: protein solubility database [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

*ProSol* DB. The distances varied between 0 and 40 (Figure 4a), where a score of 0 indicated a specific strain of *E. coli*, and a score of 40 indicated a very distant species such as *Drosophila melanogaster* (fruit fly). The average solubility confidence score per organism was 53.43 (Figure 4b). The correlation between the phylogenetic distance and average solubility confidence score per organism (Figure 4c) showed a low positive correlation as calculated using the Spearman's rank test (Spearman's rank coefficient = 0.176919). We conclude from Figures 3 and 4 that there is low to no correlation between the propensity for soluble expression and EC numbers and between solubility confidence scores and the similarity of the source organisms to *E. coli*. Therefore, when implementing a pathway, protein sequences are selected with respect to their individual solubility confidence scores.

### 3.2 | Validating *ProPASS* workflow

We used the *ProPASS* workflow to identify synthesis pathways and their potential implementation for several commercially useful biomolecules. The targets include industrial precursors used in the production of plastics, printing ink, pharmaceuticals (2,3-butanediol [Cho et al., 2015], *cis*-muconic acid [Weber et al., 2012]), biofuels (triacylglycerols [Radakovits et al., 2010] and fatty acid methyl esters [Nawabi et al., 2011]), a commodity chemical that is widely used in industry (3-hydroxypropanoic acid [Della Pina, Falletta, & Rossi, 2011]) and biological precursors (isopentenyl diphosphate [Kim & Keasling, 2001] and *myo*-inositol [Shiue & Prather, 2014]). The test cases show that 20%–100% of the identified pathways have solubility confidence scores for enzymes catalyzing all their reactions from the reviewed sequences present in *ProSol* DB (Table 2). For example, for target molecule *myo*-inositol, *ProPASS* identified 384 pathways, of which 235 pathways (61% of the identified pathways) had solubility confidence scores for all reactions. The coverage is expected to further increase if the user wishes to utilize unreviewed sequences. Three of the seven test cases mentioned above were analyzed in detail of which two are presented below (3-hydroxypropanoic acid and 2,3-butanediol), and one in Supporting Information Files 4 and 5 (*cis*-muconic acid). In all three test cases, we constructed the pathways using *ProPath*. The identified pathways were ordered twice, once based on pathway yield, and again based on pathway length. Then, for the five highest yielding and for the five shortest pathways, we used *ProPASS* to identify potential implementations based on solubility confidence scores from *ProSol* DB. We compared the recommended implementations with the published implementation of pathways reported in the literature.

### 3.3 | 3-Hydroxypropanoic acid

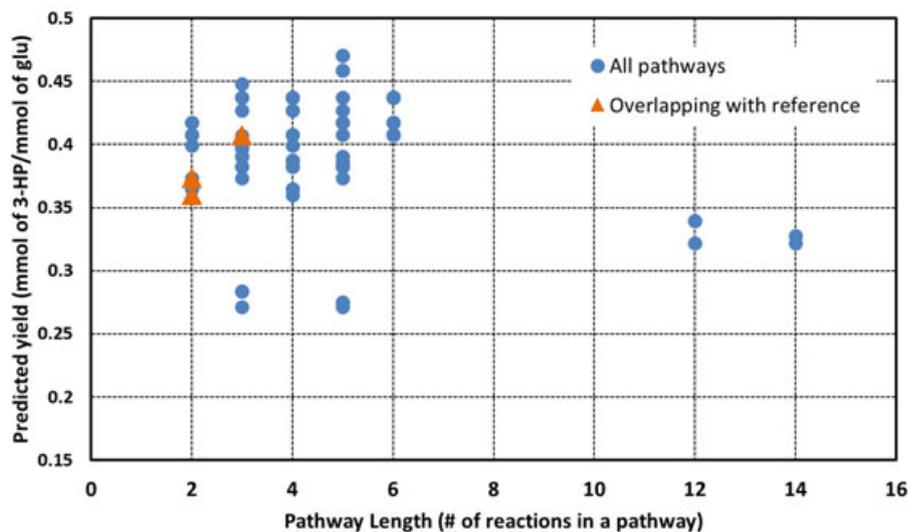
We applied *ProPASS* to synthesize 3-hydroxypropanoic acid (3-HP). Since 3-HP production through the *Pdu* pathway is the only known pathway from glycerol, we only explored pathways originating from

**TABLE 2** A list of commercially useful molecules used to evaluate the ProPASS workflow

Target molecule	# Pathways identified by ProPASS	# Pathways with solubility confidence scores for all reactions	Pathways with solubility confidence scores for all their reactions (%)
2,3-Butanediol	1	1	100
cis-Muconic acid	147	29	20
Triacylglycerols	83	38	46
Fatty acid methyl ester	48	28	58
3-Hydroxypropanoic acid	112	40	36
Isopentenyl diphosphate	33	16	48
myo-Inositol	384	235	61

Note. ProPASS: Probabilistic Pathway Assembly with Solubility Confidence Scores.

**FIGURE 5** Predicted yields and corresponding lengths for pathways producing 3-hydroxypropanoic acid identified by ProPASS. Each point on the plot represents one pathway. The three orange triangles coincide with pathways reported in the literature. ProPASS: Probabilistic Pathway Assembly with Solubility Confidence Scores [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]



glucose and identified 112 synthesis pathways (Figure 5). Pathways lengths varied from 2 to 14 and yield values varied from 0.271 to 0.41 mmol/mmol of glucose.

When pathways were ordered based on yield (Table 3a), the top pathway had five steps and a yield of 0.408 mmol/mmol of glucose. However, two of the five reactions had no solubility data in *ProSol DB*. Other pathways in the top five ordered pathways had three or five steps. Two of the three-step pathways, yielding a value of 0.407 mmol/mmol of glucose, had complete solubility confidence scores.

When pathways were ordered based on length (Table 3b), the top five had only two reaction steps—from malonate semialdehyde to 3-HP. All five pathways had similar yields (0.365–0.373 mmol/mmol of glucose). Only two of the five pathways had complete solubility data.

We compared the predicted solubility confidence scores of the enzymes used in the published literature to the highest solubility confidence score alternatives identified by ProPASS (Figure 6). Of the identified pathways, three have been experimentally used previously for 3-HP production either through the propionyl-CoA pathway (Figure 6a; Luo et al., 2016), the malonyl-CoA pathway (Figure 6b; Cheng et al., 2016), or  $\beta$ -alanine pathway (Figure 6c; Song et al., 2016). In the case of the enzyme malonyl-CoA reductase in the malonyl-CoA pathway (Figure 6b), ProPASS identified a solubility

confidence score of 74, whereas the corresponding enzyme from Cheng et al. (2016) has a score of 97. *UniProtKB* did not catalog the sequences used by Cheng et al. (2016), which precluded their enzyme from the ProPASS workflow results. For 3-HP synthesis through the propionyl-CoA and the  $\beta$ -alanine pathways, our workflow identified sequences with high or similar solubility confidence scores when compared with sequences used in the literature.

### 3.4 | 2,3-Butanediol

ProPASS identified only one two-step synthesis pathway that produces 2,3-butanediol using the (*R,R*)-butanediol dehydrogenase and acetolactate decarboxylase enzymes (Table 4). The highest solubility confidence scores for sequences identified by ProPASS from the A-list organisms were 87 (*S. cerevisiae*) and 83 (*Lactococcus lactis*), respectively. Expanding the search to include sequences from the B-list did not yield higher scores. Upon expanding the search to include reviewed and non-reviewed sequences from *UniProtKB*, ProPASS identified sequences with solubility confidence scores of 100 (*L. lactis*) for (*R,R*)-butanediol dehydrogenase and 99 (*Enterobacter aerogenes*) for acetolactate decarboxylase enzymes. We compared our findings with those in the literature and found that this two-step pathway is the only known synthesis pathway for 2,3-butanediol. Further, the selected enzymes

**TABLE 3** Implementation options for (a) five highest yielding pathways and (b) five shortest pathways. The columns list the pathway yield, length, names of reactions) in the pathway, the number of sequences with a solubility confidence score that was identified using *ProSol DB* for each reaction step, and the maximum solubility confidence score value per reaction step. Entries in parentheses indicate that multiple enzymes catalyze the reaction step, and in this case, the number of sequences and the maximum solubility is reported per enzyme for that reaction step. A “-” indicates that no sequence was identified in both *ProSol DB* and nonreviewed sequences in *UniProtKB*

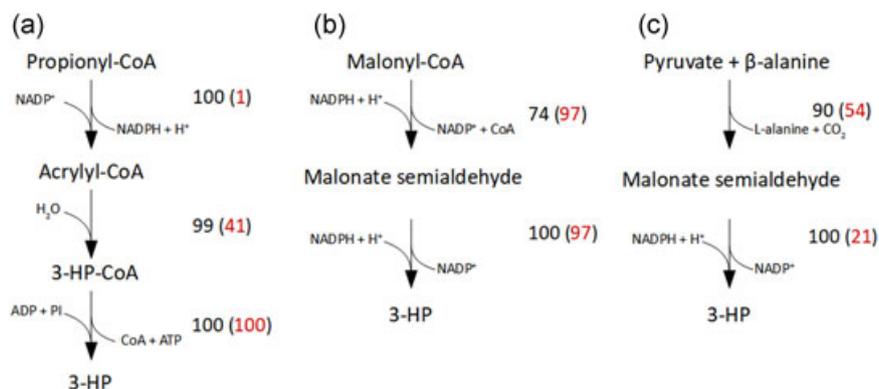
(a)				
Pathway yield	Pathway length	Names of reactions in pathway	# Sequences per reaction	Max sol. conf. score per step
0.408	5	3-Hydroxypropionate:CoA ligase (AMP-forming),	2,	34,
		3-Hydroxypropionyl-CoA:NADP+ oxidoreductase,	350,	100,
		3-Oxopropionyl-CoA hydrolase,	1,	42,
		3-Oxopropanoate:NADP+ oxidoreductase,	-,	-,
		Acetyl-CoA:malonate CoA transferase	-	-
0.407	3	3-Hydroxypropionate:CoA ligase (AMP-forming),	2,	34,
		3-Hydroxypropionyl-CoA hydro-lyase,	(255,1),	(100,37),
		Propanoyl-CoA:NADP+ 2-oxidoreductase	(5,163)	(95,100)
0.407	3	3-Hydroxypropionyl-CoA synthetase (ADP-forming),	159,	100,
		3-Hydroxypropionyl-CoA hydro-lyase	(255,1),	(100,37),
		Propanoyl-CoA:NADP+ 2-oxidoreductase	(5,163)	(95,100)
0.397	3	3-Hydroxypropionyl-CoA synthetase (ADP-forming),	159,	100,
		3-Hydroxypropionyl-CoA hydro-lyase,	(255,1),	(100, 37),
		Propanoyl-CoA:NAD+ oxidoreductase	1	99
0.391	5	KEGG R03158(no name),	14,	100,
		3-Hydroxypropionyl-CoA:NADP+ oxidoreductase,	350,	100,
		3-Oxopropionyl-CoA hydrolase,	1,	42,
		3-Oxopropanoate:NAD+ oxidoreductase,	-,	-,
		Acetyl-CoA:malonate CoA-transferase	-	-
(b)				
Pathway yield	Pathway length	Names of reactions in pathway	# sequences per reaction	Max sol. conf. score per step
0.373	2	3-Hydroxypropanoate:NAD+ oxidoreductase,	-,	-,
		L-Alanine:3-oxopropanoate aminotransferase	2	90
0.373	2	3-Hydroxypropanoate:NAD+ oxidoreductase,	-,	-,
		$\beta$ -Alanine:2-oxoglutarate aminotransferase	120	99
0.367	2	3-Hydroxypropanoate:NAD+ oxidoreductase,	-,	-,
		3-Oxopropanoate:NADP+ oxidoreductase (CoA-malonylating)	2	74
0.365	2	3-Hydroxypropionate:NADP+ oxidoreductase,	386,	100,
		L-Alanine:3-oxopropanoate aminotransferase	2	90
0.365	2	3-Hydroxypropionate:NADP+ oxidoreductase,	386,	100,
		$\beta$ -Alanine:2-oxoglutarate aminotransferase	20	99

Note. KEGG: Kyoto Encyclopedia of Genes and Genomes; *ProSol DB*: Protein Solubility Database

match an experimentally validated pathway by Yan, Lee, and Liao (2009), where enzyme sequences were selected from two A-list organisms, *Bacillus subtilis* (*B. subtilis*) and *Klebsiella pneumoniae*. Both sequences have high predicted solubility confidence scores, comparable to those identified in the nonreviewed sequences. As these sequences were not cataloged in *UniProtKB/Swiss-Prot*, *ProPASS* did not identify them as possible sequence options.

## 4 | DISCUSSION

Productive integration of heterologous pathways into model host organisms is vital for metabolic engineering. Key design steps where computational tools can significantly expedite the engineering cycle are pathway construction to identify reactions required to connect the host to a target metabolite, and parts selection to choose specific



**FIGURE 6** Three pathways for producing 3-hydroxypropanoic acid (3-HP) in *E. coli* were identified by *ProPASS* and proposed by (a) Luo et al. (2016), (b) Cheng et al. (2016), and (c) Song et al. (2016). Two predicted solubility scores are provided for each enzymatic reaction. The black score is for a sequence from organism recommended by *ProPASS* and the red score is for the corresponding sequence that was used experimentally. *ProPASS*: Probabilistic Pathway Assembly with Solubility confidence Scores [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

biochemical parts (such as enzymes, promoters, etc.) to express in the host. While current pathway construction tools can identify high-yielding pathways, they cannot be directly translated into experiments without the parts selection step. Next-generation synthesis tools must provide detailed biochemical and biophysical information about the biological parts that can be directly translated and integrated into the engineering design-build-test-learn cycle. Such information includes, but is not limited to, the sequence of heterologous enzymes, their activity, solubility, stability, and so forth.

Aiming to create a next-generation synthesis tool, we expanded in this work the scope of conventional pathway construction to identify potential sources from which each of the pathway enzymes can be isolated using protein solubility as a criterion of parts selection and host-compatibility. A specific amino acid sequence of enzymes in the pathway determines many outcomes such as kinetics, stability, and actual yield. Solubility, however, is a barrier in achieving these outcomes. When a recombinant enzyme is expressed as an insoluble aggregate (known as inclusion bodies), they retain none or a fraction of the catalytic activity, except in a few cases (Choi et al., 2011; Lin, Zhou, Wu, Xing, & Zhao, 2013; Nahalka & Nidetzky, 2007; Zhou, Xing, Wu, Zhang, & Lin, 2012). An expressed recombinant enzyme, therefore, must be soluble in the heterologous host to ensure product yield and titer.

Current metabolic engineering practices utilize an ad hoc approach to enzyme sourcing based on domain knowledge or practical considerations such as accessibility to organisms from which enzymes can be sourced. For example, the *n*-butanol synthesis pathway is often sourced entirely from natural producers, *Clostridia* (Inui et al., 2008; Nielsen et al., 2009). However, when solubility and expression of heterologous enzymes (i.e., host compatibility) are taken into consideration, several groups have found that a “mix-and-match” approach is often more profitable than sourcing all enzymes from a single organism, as was demonstrated when implementing the *n*-butanol pathway in yeast (Steen et al., 2008) and *E. coli* (Bond-Watts, Bellerose, & Chang, 2011). More drastic examples of such sourcing paradigm were demonstrated by both Keasling and Smolke groups when synthesizing complex natural products (Galanie et al., 2015; Ro et al., 2006). Assembly of such “mix-and-match” pathways has been significantly aided by decreasing costs of gene synthesis. However, such empowerment also leads to choice overload. We posit that host-compatibility in the form of solubility can provide a powerful guide to prune down the options to most profitable choices. Enzyme engineering, for improved catalytic efficiency, substrate and product preferences, and solubility, is traditionally used as a strategy to overcome issues with product synthesis. It can, however, be time-consuming, even when aided by computational protein engineering tools (e.g., Hassanpour, Ullah, Yousofshahi, Nair, & Hassoun, 2017; Sachsenhauser & Bardwell, 2018),

**TABLE 4** The pathway identified to produce 2,3-butanediol consists of two reaction steps, each catalyzed by a single enzyme. The table shows solubility confidence scores for enzymes from organisms recommended by *ProPASS* and for enzymes from organisms used by Yan et al. (2009)

Enzyme	Pathway identified by <i>ProPASS</i>				From Yan et al. (2009)	
	Reviewed A-list organisms	Predicted sol. conf. score	All sequences (reviewed & nonreviewed)	Predicted sol. conf. score	Organisms	Predicted sol. conf. score
( <i>R,R</i> )-butanediol dehydrogenase	<i>S. cerevisiae</i>	87	<i>L. lactis</i>	100	<i>B. subtilis</i>	96
acetolactate decarboxylase	<i>L. lactis</i>	83	<i>E. aerogenes</i>	99	<i>K. pneumoniae</i>	98

Note. *B. subtilis*: *Bacillus subtilis*; *E. aerogenes*: *Enterobacter aerogenes*; *K. pneumoniae*: *Klebsiella pneumoniae*; *L. lactis*: *Lactococcus lactis*; *ProPASS*: Probabilistic Pathway Assembly with Solubility confidence Scores; *S. cerevisiae*: *Saccharomyces cerevisiae*

especially if several enzymes within a pathway are insoluble. Thus, implementing pathways with the knowledge that the enzymes selected have a high propensity for solubility in the host can reduce the traditional trial-and-error approach to enzyme selection. The product of *ProPASS* is a listing of synthesis pathways and recommended implementations. *ProPASS*, therefore, is a tool to aid in exploring the design space of synthesis pathways and their implementation.

Protein solubility is used as a rating and/or selection criterion before the implementation phase. In most cases, multiple implementation options based on solubility confidence scores are available. Since solubility confidence scores do not affect maximum theoretical yield, they simply provide a guideline for the implementation phase. Thus, *ProPASS* advances the state of pathway synthesis by exploring the underlying biophysical design space, leading to more profitable design implementation options when compared to using pathway synthesis tools that only explore the abstract metabolic space. Such incorporation of sequence information is distinct from previous descriptions where protein sequence and structure information has been used to inform putative promiscuous enzymatic activity to assemble nonnatural synthetic pathways (Brunk, Neri, Tavernelli, Hatzimanikatis, & Rothlisberger, 2012; Erb, Jones, & Bar-Even, 2017).

To allow for enzyme selection, we created a database of predicted solubility confidence scores, *ProSol DB*, for reviewed protein sequences from the *UniProtKB/Swiss-Prot* database. *ProSol DB* contains 240,016 sequences and links EC numbers with protein sequences, their source organism, and their solubility confidence scores in *E. coli*. The solubility confidence scores were calculated using *ccSOL omics*. Future improvements in solubility prediction algorithms can be used to update *ProSol DB*, thus providing more confidence regarding enzyme selection. Analysis of solubility confidence scores and enzymes revealed that sequences associated with a particular enzyme have a wide range of associated scores. Further analysis showed no correlation between phylogenetic distance from *E. coli* and solubility confidence scores. Our findings emphasize that solubility is a property of the encoding sequence, and not a function of the host organism or associated EC number. As far as we know, *ProSol DB* is the largest database of enzyme solubility confidence scores for heterologous enzymes in *E. coli* and serves as a resource for pathway synthesis tools as well as for meta-analysis of sequence-independent factors that affect the solubility of enzymes in *E. coli*.

We validated the *ProPASS* workflow by comparing predicted pathway implementations with those published in the literature for three target molecules, 3-hydroxypropanoic acid, 2,3-butanediol, and *cis*-muconic acid. In all cases, we identified a significant number of pathway implementations for which solubility confidence scores were available through *ProSol DB*. The availability of data was dependent on the studied target molecule. In the case of *myo*-inositol, 61% of identified pathways had confidence scores for each enzyme, whereas in the case of *cis*-muconic acid only 20% of identified pathways had such solubility data. Several observations can be drawn based on examining the three detailed test cases. When considering

solubility, longer pathways have a higher chance of having enzymes with low solubility confidence scores. For example, the *cis*-muconic acid synthesis pathway with the highest yield has a pathway length of 10, with one of the steps having a low solubility confidence score of 22. In contrast, the pathway with the highest confidence of having a soluble enzyme for each step is a three-step pathway but has only half the yield of the highest-yield pathway. *ProPASS* can provide several implementation alternatives, from which one or more pathways can be selected for experimental testing.

When compared to solubility confidence scores of published implementations for the three test cases, *ProPASS* recommended the published implementations provided they were in *UniProtKB*. In some cases, *ccSOL omics* predicted the published successful enzyme implementation as insoluble, with solubility confidence scores less than 30. This was the case for three (e.g., the first step for producing 3-HP from propionyl-CoA, Figure 6a, the second step of producing 3-HP from pyruvate and  $\beta$ -alanine, Figure 6c, and the first step of producing *cis*-muconic acid from *para*-hydroxybenzoic acid; Figure S2b) of 12 comparisons with known soluble enzymes reported in the literature (Figure 6, Table 4, and Figure S2). Such mispredictions are expected considering that the accuracy of *ccSOL omics* for the enzyme solubility test set was evaluated to be 73.46%. When the solubility confidence score for a given enzyme was low, it was possible to utilize nonreviewed sequences in *UniProtKB* to identify higher score alternatives. This was the case when synthesizing 2,3-butanediol. As more reviewed sequences are added to *UniProtKB/Swiss-Prot*, the reliance on nonreviewed sequences will diminish. Importantly, *ProPASS* can provide alternative implementations that would not have been easily discoverable otherwise. All test cases discussed here are limited to using *E. coli* as a host. However, the development of solubility prediction tools for other organisms would allow extending *ProPASS* to other hosts.

Strain optimization, the process by which opportunities for metabolic redirection are identified to ensure maximal target formation, is a crucial step that ensures the profitable integration of heterologous synthesis pathways into model host organisms. While several computational tools (Burgard, Pharkya, & Maranas, 2003; Patil, Rocha, Förster, & Nielsen, 2005; Ranganathan, Suthers, & Maranas, 2010) were proposed to identify the most profitable modifications, almost none take into account implementation concerns during optimization. One exception is Chance Constrained Optimization (*CCOpt*; Yousofshahi, Orshansky, Lee, & Hassoun, 2013), which addresses the uncertainty in precise tuning of enzyme levels and imposes probabilistic flux capacity constraints that capture the uncertainty in tuning enzyme levels during implementation. In this work, we treated the two design steps (pathway synthesis and strain optimization) as independent. However, a competitive (in terms of yield) synthesis pathway during the pathway synthesis step may provide low or no yield once implementations are investigated. When developed judiciously to limit the design space, co-optimization tools that are aware of design choices associated with the underlying implementation can potentially yield more profitable experimental guidance.

## ACKNOWLEDGEMENTS

We thank Emily Chicklis (REU student) and Trevor B. Nicks for their help with some figure illustrations. This material is based upon work supported by the National Science Foundation under grant no. CCF-1421972, National Institutes of Health under grant no. DP2HD091798, and a Tufts Collaborates grant.

## ORCID

Nikhil U. Nair  <http://orcid.org/0000-0001-7737-1385>

Soha Hassoun  <http://orcid.org/0000-0001-9477-2199>

## REFERENCES

- Agostini, F., Cirillo, D., Livi, C. M., Delli ponti, R., & Tartaglia, G. G. (2014). ccSOL omics: A webserver for solubility prediction of endogenous and heterologous expression in *Escherichia coli*. *Bioinformatics*, 30(20), 2975–2977.
- Blum, T., & Kohlbacher, O. (2008). MetaRoute: Fast search for relevant metabolic routes for interactive network navigation and visualization. *Bioinformatics*, 24(18), 2108–2109.
- Bond-Watts, B. B., Bellerose, R. J., & Chang, M. C. (2011). Enzyme mechanism as a kinetic control element for designing synthetic biofuel pathways. *Nature Chemical Biology*, 7(4), 222–227.
- Brunk, E., Neri, M., Tavernelli, I., Hatzimanikatis, V., & Rothlisberger, U. (2012). Integrating computational methods to retrofit enzymes to synthetic pathways. *Biotechnology and Bioengineering*, 109(2), 572–582. <https://doi.org/10.1002/bit.23334>
- Burgard, A. P., Pharkya, P., & Maranas, C. D. (2003). Optknock: A bilevel programming framework for identifying gene knockout strategies for microbial strain optimization. *Biotechnology and Bioengineering*, 84(6), 647–657.
- Caspi, R., Foerster, H., Fulcher, C. A., Kaipa, P., Krummenacker, M., Latendresse, M., & Tissier, C. (2008). The metacyc database of metabolic pathways and enzymes and the biocyc collection of pathway/genome databases. *Nucleic Acids Research*, 36(suppl 1), D623–D631.
- Chang, C. C. H., Song, J., Tey, B. T., & Ramanan, R. N. (2014). Bioinformatics approaches for improved recombinant protein production in *Escherichia coli*: Protein solubility prediction. *Briefings in Bioinformatics*, 15(6), 953–962.
- Chen, L., Oughtred, R., Berman, H. M., & Westbrook, J. (2004). TargetDB: A target registration database for structural genomics projects. *Bioinformatics*, 20(16), 2860–2862.
- Cheng, Z., Jiang, J., Wu, H., Li, Z., & Ye, Q. (2016). Enhanced production of 3-hydroxypropionic acid from glucose via malonyl-CoA pathway by engineered *Escherichia coli*. *Bioresource Technology*, 200, 897–904.
- Cho, S., Kim, T., Woo, H. M., Kim, Y., Lee, J., & Um, Y. (2015). High production of 2, 3-butanediol from biodiesel-derived crude glycerol by metabolically engineered *Klebsiella oxytoca* M1. *Biotechnology for Biofuels*, 8(1), 1.
- Choi, S.-L., Lee, S. J., Ha, J.-S., Song, J. J., Rhee, Y. H., & Lee, S.-G. (2011). Generation of catalytic protein particles in *Escherichia coli* cells using the cellulose-binding domain from *Cellulomonas fimi* as a fusion partner. *Biotechnology and Bioengineering*, 16(6), 1173–1179.
- Della Pina, C., Falletta, E., & Rossi, M. (2011). A green approach to chemical building blocks. The case of 3-hydroxypropanoic acid. *Green chemistry*, 13(7), 1624–1632.
- Ellis, L. B., & Wackett, L. P. (2012). Use of the University of Minnesota Biocatalysis/Biodegradation Database for study of microbial degradation. *Microbial Informatics and Experimentation*, 2(1), 1.
- Erb, T. J., Jones, P. R., & Bar-Even, A. (2017). Synthetic metabolism: Metabolic engineering meets enzyme design. *Current Opinion in Chemical Biology*, 37, 56–62. <https://doi.org/10.1016/j.cbpa.2016.12.023>
- Fidler, S., & Dennis, D. (1992). Polyhydroxyalkanoate production in recombinant *Escherichia coli*. *FEMS Microbiology Reviews*, 9(2-4), 231–235.
- Galanie, S., Thodey, K., Trenchard, I. J., Interrante, M. F., & Smolke, C. D. (2015). Complete biosynthesis of opioids in yeast. *Science*, 349(6252), 1095–1100.
- Greene, N., Judson, P. N., Langowski, J. J., & Marchant, C. A. (1999). Knowledge-based expert systems for toxicity and metabolism prediction: DEREK, StAR and METEOR. *SAR and QSAR in Environmental Research*, 10(2-3), 299–314. <https://doi.org/10.1080/10629369908039182>
- Habibi, N., Hashim, S. Z. M., Norouzi, A., & Samian, M. R. (2014). A review of machine learning methods to predict the solubility of over-expressed recombinant proteins in *Escherichia coli*. *BMC Bioinformatics*, 15(1), 134.
- Hassanpour, N., Ullah, E., Yousofshahi, M., Nair, N. U., & Hassoun, S. (2017). Selection Finder (SelFi): A computational metabolic engineering tool to enable directed evolution of enzymes. *Metabolic Engineering Communications*, 4, 37–47. <https://doi.org/10.1016/j.meten.2017.02.003>
- Hou, B. K., Wackett, L. P., & Ellis, L. B. M. (2003). Microbial pathway prediction: A functional group approach. *Journal of Chemical Information and Computer Sciences*, 43(3), 1051–1057. <https://doi.org/10.1021/ci034018f>
- Inui, M., Suda, M., Kimura, S., Yasuda, K., Suzuki, H., Toda, H., ... Yukawa, H. (2008). Expression of *Clostridium acetobutylicum* butanol synthetic genes in *Escherichia coli*. *Applied Microbiology and Biotechnology*, 77(6), 1305–1316.
- Kanehisa, M., & Goto, S. (2000). KEGG: Kyoto encyclopedia of genes and genomes. *Nucleic Acids Research*, 28(1), 27–30.
- Khurana, S., Rawi, R., Kunji, K., Chuang, G.-Y., Bensmail, H., & Mall, R. (2018). DeepSol: A deep learning framework for sequence-based protein solubility prediction. *Bioinformatics*, 34, 2605–2613.
- Kim, S. W., & Keasling, J. (2001). Metabolic engineering of the nonmevalonate isopentenyl diphosphate synthesis pathway in *Escherichia coli* enhances lycopene production. *Biotechnology and Bioengineering*, 72(4), 408–415.
- Klopman, G., Dimayuga, M., & Talafous, J. (1994). META. 1. A program for the evaluation of metabolic transformation of chemicals. *Journal of Chemical Information and Computer Sciences*, 34(6), 1320–1325. <https://doi.org/10.1021/ci00022a014>
- Klopman, G., Tu, M., & Talafous, J. (1997). META. 3. A genetic algorithm for metabolic transform priorities optimization. *Journal of Chemical Information and Computer Sciences*, 37(2), 329–334. <https://doi.org/10.1021/ci9601123>
- Letunic, I., & Bork, P. (2016). Interactive tree of life (iTOL) v3: An online tool for the display and annotation of phylogenetic and other trees. *Nucleic Acids Research*, 44(W1), W242–W245. <https://doi.org/10.1093/nar/gkw290>
- Lin, Z., Zhou, B., Wu, W., Xing, L., & Zhao, Q. (2013). Self-assembling amphipathic alpha-helical peptides induce the formation of active protein aggregates in vivo. *Faraday Discussions*, 166, 243–256.
- Luo, H., Zhou, D., Liu, X., Nie, Z., Quiroga-Sánchez, D. L., & Chang, Y. (2016). Production of 3-Hydroxypropionic acid via the propionyl-CoA pathway using recombinant *Escherichia coli* strains. *PLOS One*, 11(5), e0156286.
- Makrides, S. C. (1996). Strategies for achieving high-level expression of genes in *Escherichia coli*. *Microbiological Reviews*, 60(3), 512–538.
- Marchant, C. A., Briggs, K. A., & Long, A. (2008). *In silico* tools for sharing data and knowledge on toxicity and metabolism: Derek for windows,

- meteor, and vitic. *Toxicology Mechanisms and Methods*, 18(2-3), 177-187. <https://doi.org/10.1080/15376510701857320>
- Martin, V. J. J., Pitera, D. J., Withers, S. T., Newman, J. D., & Keasling, J. D. (2003). Engineering a mevalonate pathway in *Escherichia coli* for the production of terpenoids. *Nature Biotechnology*, 21(7), 796-802.
- Moriya, Y., Shigemizu, D., Hattori, M., Tokimatsu, T., Kotera, M., Goto, S., & Kanehisa, M. (2010). PathPred: An enzyme-catalyzed metabolic pathway prediction server. *Nucleic Acids Research*, 38, gkq318-W143.
- Nahalka, J., & Nidetzky, B. (2007). Fusion to a pull-down domain: A novel approach of producing *Trigonopsis variabilis* D-amino acid oxidase as insoluble enzyme aggregates. *Biotechnology and Bioengineering*, 97(3), 454-461.
- Nawabi, P., Bauer, S., Kyrpides, N., & Lykidis, A. (2011). Engineering *E. coli* for biodiesel production utilizing a bacterial fatty acid methyltransferase. *Applied and environmental microbiology*, AEM, 77, 05046-05011.
- Nielsen, D. R., Leonard, E., Yoon, S.-H., Tseng, H.-C., Yuan, C., & Prather, K. L. J. (2009). Engineering alternative butanol production platforms in heterologous bacteria. *Metabolic Engineering*, 11(4-5), 262-273.
- Overbeek, R., Begley, T., Butler, R. M., Choudhuri, J. V., Chuang, H.-Y., Cohoon, M., ... Edwards, R. (2005). The subsystems approach to genome annotation and its use in the project to annotate 1000 genomes. *Nucleic Acids Research*, 33(17), 5691-5702.
- Patil, K. R., Rocha, I., Förster, J., & Nielsen, J. (2005). Evolutionary programming as a platform for *in silico* metabolic engineering. *BMC Bioinformatics*, 6(1), 308.
- Peirú, S., Menzella, H. G., Rodríguez, E., Carney, J., & Gramajo, H. (2005). Production of the potent antibacterial polyketide erythromycin C in *Escherichia coli*. *Applied and Environmental Microbiology*, 71(5), 2539-2547.
- Pfeifer, B. A., Admiraal, S. J., Gramajo, H., Cane, D. E., & Khosla, C. (2001). Biosynthesis of complex polyketides in a metabolically engineered strain of *E. coli*. *Science*, 291(5509), 1790-1792.
- Pitera, D. J., Paddon, C. J., Newman, J. D., & Keasling, J. D. (2007). Balancing a heterologous mevalonate pathway for improved isoprenoid production in *Escherichia coli*. *Metabolic Engineering*, 9(2), 193-207.
- Radakovits, R., Jinkerson, R. E., Darzins, A., & Posewitz, M. C. (2010). Genetic engineering of algae for enhanced biofuel production. *Eukaryotic Cell*, 9(4), 486-501.
- Ranganathan, S., Suthers, P. F., & Maranas, C. D. (2010). OptForce: An optimization procedure for identifying all genetic manipulations leading to targeted overproductions. *PLOS Computational Biology*, 6(4), e1000744.
- Ro, D.-K., Paradise, E. M., Ouellet, M., Fisher, K. J., Newman, K. L., Ndungu, J. M., ... Kirby, J. (2006). Production of the antimalarial drug precursor artemisinic acid in engineered yeast. *Nature*, 440(7086), 940-943.
- Rodrigo, G., Carrera, J., Prather, K. J., & Jaramillo, A. (2008). DESHARKY: Automatic design of metabolic pathways for optimal cell growth. *Bioinformatics*, 24(21), 2554-2556.
- Sachsenhauser, V., & Bardwell, J. C. (2018). Directed evolution to improve protein folding in vivo. *Current Opinion in Structural Biology*, 48, 117-123.
- Shiue, E., & Prather, K. L. (2014). Improving D-glucaric acid production from myo-inositol in *E. coli* by increasing MIOX stability and myo-inositol transport. *Metabolic Engineering*, 22, 22-31.
- Singh, S. M., & Panda, A. K. (2005). Solubilization and refolding of bacterial inclusion body proteins. *Journal of Bioscience and Bioengineering*, 99(4), 303-310.
- Song, C. W., Kim, J. W., Cho, I. J., & Lee, S. Y. (2016). Metabolic engineering of *Escherichia coli* for the production of 3-hydroxypropionic acid and malonic acid through  $\beta$ -alanine route. *ACS Synthetic Biology*, 5, 1256-1263.
- Steen, E. J., Chan, R., Prasad, N., Myers, S., Petzold, C. J., Redding, A., ... Keasling, J. D. (2008). Metabolic engineering of *Saccharomyces cerevisiae* for the production of n-butanol. *Microbial Cell Factories*, 7(1), 36.
- Steen, E. J., Kang, Y., Bokinsky, G., Hu, Z., Schirmer, A., McClure, A., ... Keasling, J. D. (2010). Microbial production of fatty-acid-derived fuels and chemicals from plant biomass. *Nature*, 463, 559-562.
- Takahashi, H., Kumagai, T., Kitani, K., Mori, M., Matoba, Y., & Sugiyama, M. (2007). Cloning and characterization of a *Streptomyces* single module type non-ribosomal peptide synthetase catalyzing a blue pigment synthesis. *Journal of Biological Chemistry*, 282(12), 9073-9081.
- Tartaglia, G. G., Cavalli, A., & Vendruscolo, M. (2007). Prediction of local structural stabilities of proteins from their amino acid sequences. *Structure*, 15(2), 139-143.
- Tartaglia, G. G., Pechmann, S., Dobson, C. M., & Vendruscolo, M. (2009). A relationship between mRNA expression levels and protein solubility in *E. coli*. *Journal of Molecular Biology*, 388(2), 381-389.
- Tong, I.-T., Liao, H. H., & Cameron, D. C. (1991). 1, 3-Propanediol production by *Escherichia coli* expressing genes from the *Klebsiella pneumoniae* dha regulon. *Applied and Environmental Microbiology*, 57(12), 3541-3546.
- Trésaugues, L., Collinet, B., Minard, P., Henckes, G., Aufrère, R., Blondeau, K., ... van Tilbeurgh, H. (2004). Refolding strategies from inclusion bodies in a structural genomics project. *Journal of Structural and Functional Genomics*, 5(3), 195-204.
- UniProt Consortium. (2017). UniProt: The universal protein knowledge-base. *Nucleic Acids Research*, 45(D1), D158-D169.
- Watts, K. T., Mijts, B. N., & Schmidt-Dannert, C. (2005). Current and emerging approaches for natural product biosynthesis in microbial cells. *Advanced Synthesis & Catalysis*, 347(7-8), 927-940. <https://doi.org/10.1002/adsc.200505062>
- Weber, C., Brückner, C., Weinreb, S., Lehr, C., Essl, C., & Boles, E. (2012). Biosynthesis of *cis*, *cis*-Muconic acid and its aromatic precursors, catechol and protocatechuic acid, from renewable feedstocks by *Saccharomyces cerevisiae*. *Applied and Environmental Microbiology*, 78(23), 8421-8430.
- Wu, D., Wang, Q., Assary, R. S., Broadbelt, L. J., & Krilov, G. (2011). A computational approach to design and evaluate enzymatic reaction pathways: Application to 1-butanol production from pyruvate. *Journal of Chemical Information and Modeling*, 51(7), 1634-1647.
- Xia, P. F., Zhang, G. C., Liu, J. J., Kwak, S., Tsai, C. S., Kong, I. I., ... Jin, Y. S. (2016). GroE chaperonins assisted functional expression of bacterial enzymes in *Saccharomyces cerevisiae*. *Biotechnology and Bioengineering*, 113, 2149-2155.
- Yan, Y., Lee, C.-C., & Liao, J. C. (2009). Enantioselective synthesis of pure (*R*, *R*)-2, 3-butanediol in *Escherichia coli* with stereospecific secondary alcohol dehydrogenases. *Organic & Biomolecular Chemistry*, 7(19), 3914-3917.
- Yousoufshahi, M., Lee, K., & Hassoun, S. (2011). Probabilistic pathway construction. *Metabolic Engineering*, 13(4), 435-444.
- Yousoufshahi, M., Orshansky, M., Lee, K., & Hassoun, S. (2013). Probabilistic strain optimization under constraint uncertainty. *BMC Systems Biology*, 7, 29. <https://doi.org/10.1186/1752-0509-7-29>
- Zhou, B., Xing, L., Wu, W., Zhang, X.-E., & Lin, Z. (2012). Small surfactant-like peptides can drive soluble proteins into active aggregates. *Microbial Cell Factories*, 11(1), 10.

## SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section at the end of the article.

**How to cite this article:** Amin SA, Endalur Gopinarayanan V, Nair NU, Hassoun S. Establishing synthesis pathway-host compatibility via enzyme solubility. *Biotechnology and Bioengineering*. 2019;116:1405-1416. <https://doi.org/10.1002/bit.26959>