

1 **Towards creating an extended metabolic model (EMM) for *E. coli* using enzyme**  
2 **promiscuity prediction and metabolomics data**

3

4 Sara A. Amin\*, Department of Computer Science, Tufts University, Medford, MA,  
5 [sara.amin@tufts.edu](mailto:sara.amin@tufts.edu)

6 Elizabeth Chavez\*, Department of Biology, University of North Carolina, Chapel Hill, NC  
7 [celiz@live.unc.edu](mailto:celiz@live.unc.edu)

8 Vladimir Porokhin, Department of Computer Science, Tufts University, Medford, MA,  
9 [vladimir.porokhin@tufts.edu](mailto:vladimir.porokhin@tufts.edu)

10 Nikhil U. Nair<sup>†</sup>, Department of Chemical and Biological Engineering, Tufts University,  
11 Medford, MA, [nikhil.nair@tufts.edu](mailto:nikhil.nair@tufts.edu), and

12 Soha Hassoun<sup>†</sup>, Department of Computer Science and Department of Chemical & Biological  
13 Engineering, Tufts University, Medford, MA, [soha.hassoun@tufts.edu](mailto:soha.hassoun@tufts.edu)

14 \*Equal contributions

15 <sup>†</sup>Co-corresponding authors

16

## 1 **Abstract**

## 2 **Background**

3 Metabolic models are indispensable in guiding cellular engineering and in advancing our  
4 understanding of systems biology. As not all enzymatic activities are fully known and/or  
5 annotated, metabolic models remain incomplete, resulting in suboptimal computational analysis  
6 and leading to unexpected experimental results. We posit that one major source of unaccounted  
7 metabolism is promiscuous enzymatic activity. It is now well-accepted that most, if not all,  
8 enzymes are promiscuous – i.e., they transform substrates other than their primary substrate.  
9 However, there have been no systematic analyses of genome-scale metabolic models to predict  
10 putative reactions and/or metabolites that arise from enzyme promiscuity.

## 11 **Results**

12 Our workflow utilizes PROXIMAL – a tool that uses reactant-product transformation patterns  
13 from the KEGG database – to predict putative structural modifications due to promiscuous  
14 enzymes. Using iML1515 as a model system, we first utilized a computational workflow,  
15 referred to as Extended Metabolite Model Annotation (EMMA), to predict promiscuous  
16 reactions catalyzed, and metabolites produced, by natively encoded enzymes in *E. coli*. We  
17 predict hundreds of new metabolites that can be used to augment iML1515. We then validated  
18 our method by comparing predicted metabolites with the *Escherichia coli* Metabolome Database  
19 (ECMDB).

## 20 **Conclusions**

21 We utilized EMMA to augment the iML1515 metabolic model to more fully reflect cellular  
22 metabolic activity. This workflow uses enzyme promiscuity as basis to predict hundreds of  
23 reactions and metabolites that may exist in *E. coli* but may have not been documented in

1 iML1515 or other databases. We provide detailed analysis of 23 predicted reactions and 16  
2 associated metabolites. Interestingly, nine of these metabolites, which are in ECMDB, have not  
3 previously been documented in any other *E. coli* databases. Four of the predicted reactions  
4 provide putative transformations parallel to those already in iML1515. We suggest adding  
5 predicted metabolites and reactions to iML1515 to create an Extended Metabolic Model (EMM)  
6 for *E. coli*.

7

## 8 **Keywords**

9 Metabolic engineering, enzyme promiscuity, extended metabolic model, systems biology,  
10 enzyme activity prediction

11

## 12 **Background**

13 The engineering of metabolic networks has enabled the production of high-volume commodity  
14 chemicals such as biopolymers and fuels, therapeutics, and specialty products [1-5]. Producing  
15 such compounds requires transforming microorganisms into efficient cellular factories [6-9].  
16 Biological engineering has been aided via computational tools for constructing synthesis  
17 pathways, strain optimization, elementary flux mode analysis, discovery of hierarchical  
18 networked modules that elucidate function and cellular organization, and many others (e.g., [10-  
19 14]). These design tools rely on organism-specific metabolic models that represent cellular  
20 reactions and their substrates and products. Model reconstruction tools [15, 16] use homology  
21 search to assign function to Open Reading Frames obtained through sequencing and annotation.  
22 Once the function is identified, the corresponding biochemical transformation is assigned to the  
23 gene. Additional biological information such as gene-protein-reaction associations is utilized to

1 refine the models. Exponential growth in sequencing has resulted in an “astronomical”, or better  
2 yet, “genomical”, number of sequenced organisms [17]. There are now databases (e.g., KEGG  
3 [18], BioCyc [19], and BiGG [20]) that catalogue organism-specific metabolic models. Despite  
4 progress in sequencing and model reconstruction, the complete characterizing of cellular activity  
5 remains elusive, and metabolic models remain incomplete. One major source of uncatalogued  
6 cellular activity is attributed to orphan genes. Because of limitations of homology-based  
7 prediction of protein function, there are millions of protein sequences that are not assigned  
8 reliable functions [21]. Integrated strategies that utilize structural biology, computational  
9 biology, and molecular enzymology continue to address assigning function to orphan genes [22].  
10  
11 We focus in this paper on another major source of uncatalogued cellular activity – promiscuous  
12 enzymatic activity, which has recently been referred to as ‘underground metabolism’ [23-25].  
13 While enzymes have widely been held as highly-specific catalysts that only transform their  
14 annotated substrate to product, recent studies show that enzymatic promiscuity – enzymes  
15 catalyzing reactions other than their main reactions – is not an exception but can be a secondary  
16 task for enzymes [26-31]. More than two-fifths (44%) of KEGG enzymes are associated with  
17 more than one reaction [32]. Promiscuous activities however are not easily detectable *in vivo*  
18 since, i) metabolites produced due to enzyme promiscuity may be unknown, ii) product  
19 concentration due to promiscuous activity may be low, iii) there is no high-throughput way to  
20 relate formed products to specific enzymes, and iv) it is difficult to identify potentially unknown  
21 metabolites in complex biological samples. Outside of *in vitro* biochemical characterization  
22 studies to predict promiscuous activities, there are few resources that record details about  
23 promiscuous enzymes such as MINEs Database [33], and ATLAS [34]. Despite the current

1 wide-spread acceptance of enzyme promiscuity, and its prominent utilization to engineer  
2 catalyzing enzymes in metabolic engineering practice [35-38], promiscuous enzymatic activity is  
3 not currently fully documented in metabolic models. Advances in computing and the ability to  
4 collect large sets of metabolomics data through untargeted metabolomics provide an exciting  
5 opportunity to develop methods to identify promiscuous reactions, their catalyzing enzymes, and  
6 their products that are specific to the sample under study. The identified reactions can then be  
7 used to complete existing metabolic models.

8  
9 We describe in this paper a computational workflow that aims to extend preexisting models with  
10 reactions catalyzed by promiscuous native enzymes and validate the outcomes using published  
11 metabolomics datasets. We refer to the augmented models as Extended Metabolic Models  
12 (EMMs), and to the workflow to create them as EMMA (EMM Annotation). Each metabolic  
13 model is assumed to have a set of reactions and their compounds and KEGG reaction IDs. Each  
14 reaction, and thus transformation, is assumed to be reversible unless indicated otherwise. EMMA  
15 utilizes PROXIMAL [39], a method for creating biotransformation operators from KEGG  
16 reactions IDs using RDM (Reaction Center, Difference Region, and Matched Region) patterns  
17 [40], and then applying the operators to given molecules. While initially developed to investigate  
18 products of Cytochrome P450 (CYP) enzymes, highly promiscuous enzymes utilized for  
19 detoxification, the PROXIMAL method is generic. To create an EMM for a known metabolic  
20 model, PROXIMAL generates biotransformation operators for each reaction in the model and  
21 then applies the operators to known metabolites within the model. The outcome of our workflow  
22 is a list of putative metabolites due to promiscuous enzymatic activity and their catalyzing  
23 enzymes and reactions. In this work, we apply EMMA to iML1515, a genome-scale model of

1 *Escherichia coli* MG1655 [41]. EMMA predicts hundreds of putative reactions and their  
2 products due to promiscuous activities in *E. coli*. The putative products are then compared to  
3 measured metabolites as reported in *Escherichia coli* Metabolome Database, ECMDB [42, 43].  
4 We identify 23 new reactions and 16 new metabolites that we recommend adding to the *E. coli*  
5 model iML1515. Four of these reactions have not been catalogued prior for *E. coli* or other  
6 organisms, suggesting novel undocumented promiscuous transformations, while five other  
7 reactions are catalogued for species other than *E. coli*. Further, there were ten reactions that were  
8 catalogued in other *E. coli* databases (e.g. EcoCyc [44], and KEGG), but not in iML1515. These  
9 19 reactions led to the addition of the 16 metabolites that are new to iML1515. Additionally,  
10 there were four new reactions that present putative transformation routes that are in parallel to  
11 existing reactions in *E. coli*. No new metabolites are added due to these four reactions.

12

## 13 **Results**

14 The application of PROXIMAL to iML1515 yielded a lookup table with 1,875 biotransformation  
15 operator entries. The operators were applied on two sets of metabolites. One set consisted of 106  
16 iML1515 metabolites with predicted or measured concentration values above 1  $\mu\text{M}$  [45]. We  
17 focused on these metabolites as the assumption is that high concentration metabolites are more  
18 likely to undergo transformation by promiscuous enzymatic activity and form detectible  
19 derivatives. When applied to this set, the operators predicted the formation of 1,423 known (with  
20 PubChem IDs) metabolites of which 57 were identified to exist in *E. coli* per ECMDB. After  
21 manual curation (per Step 1 in the Methods section), our workflow recommended 16 new  
22 metabolites and 23 reactions that can be used to augment the iML1515 model. The second set of  
23 metabolites consisted of the non-high concentration metabolites in iML1515. Our workflow

1 predicted the formation of 3,694 known (with PubChem IDs) metabolites. Out of the predicted  
2 metabolites of the second set 210 derivatives are found in ECMDB. We provide a listing of all  
3 derivatives in **Supplementary File 1**. For the remainder of the Results section, we focus on  
4 detailed analysis of derivative products due to high-concentration metabolites. Results of Flux  
5 Balance Analysis and Flux Variability Analysis for the added EMMA reactions are reported in  
6 **Supplementary File 2**.

7  
8 Identified reactions were divided into four categories, C1–C4. The rationale for the various  
9 categories is explained using a decision tree (**Fig. 1**), a machine learning model that classifies  
10 data into groupings that share similar attributes [46]. With the exception of leaf nodes, each node  
11 in the tree tests the presence or absence of a particular attribute. Left branches represent the  
12 presence of the attribute, while the right branch represents the attribute's absence. Each leaf node  
13 represents a classification category and is associated with a subset of the 23 reactions. At the root  
14 node of the decision tree, we tested if a PROXIMAL predicted metabolite is in the iML1515  
15 model. If it is, and if the enzyme catalyzing the reaction within iML1515 model producing this  
16 metabolite is different than the enzyme PROXIMAL used to predict the relevant  
17 biotransformation, then it is classified in Category 1 (C1). Reactions belonging to C1 are parallel  
18 transformation to the ones in the model. They represent novel biotransformation routes between  
19 existing metabolites since they are generated using a different gene/enzyme than what is reported  
20 in iML1515. If previous conditions do not apply to the predicted product, then it is discarded as  
21 the reaction is already in iML1515.

22

1 If a predicted metabolite is not one of the known metabolites in iML1515, the decision tree  
2 determines whether the predicted metabolite and reaction set is associated with *E. coli* in other  
3 databases (KEGG and EcoCyc). If the biotransformation is present in KEGG or EcoCyc, then the  
4 predicted metabolite is classified into Category 2 (C2), reflecting a curation issue where some  
5 reactions were not included in the iML1515 model. If the predicted metabolite is not in iML1515  
6 and not associated with *E. coli* in KEGG nor listed in EcoCyc, then the decision tree determines  
7 if the same chemical transformation (same substrate and same product) is documented to occur  
8 in other organisms. Predicted biotransformations documented in KEGG for organisms other than  
9 *E. coli* are classified in Category 3 (C3). While biotransformations not found in KEGG are  
10 classified as Category 4 (C4).

11  
12 Each Category consists of a set of reactions. C1 consists of four reactions that are predicted to be  
13 catalyzed by enzymes that are different than those in iML1515. The details of the predicted  
14 reactions are shown in **Fig 2**, and **Table 1** details a comparison between those predicted reactions  
15 and their parallel reactions in iML1515. The phosphoribosyltransferase reaction between  
16 cytosine and cytidine-5'-monophosphate (CMP) is predicted to occur in *E. coli* due to EC  
17 2.4.2.10 (orotate phosphoribosyltransferase) (**Fig. 2A**) and that between 2-oxoglutarate and 2-  
18 hydroxyglutarate by EC 1.1.1.79 (glyoxylate reductase) (**Fig. 2B**). We also predict the  
19 transformation between bicarbonate and carboxyphosphate catalyzed by EC 3.6.1.7  
20 (acylphosphatase) (**Fig. 2C**). While carboxyphosphate is not in iML1515, the transformation is  
21 considered parallel to a reaction catalyzed by EC 6.3.5.5 that is documented to occur for *E. coli*  
22 in KEGG (see **Fig. 3J**). The last prediction is the coenzyme A transferase reaction between  
23 acetoacetyl-CoA and acetoacetate due to EC 2.8.3.10 (citrate CoA-transferase) (**Fig. 2D**).

1  
 2 **Table 1:** List of C1 reactions predicted by EMMA and their parallel reactions in *E. coli*  
 3 iML1515. Each of the Predicted/iML1515 reaction pair occurs between the same substrate and  
 4 product but utilize different co-substrate or cofactors.

	<b>EC number (gene)</b>	<b>Reaction</b>
<b>Predicted</b>	2.4.2.10 (b3642)	cytosine + 5-phospho- $\alpha$ -D-ribose-1-diphosphate $\rightleftharpoons$ CMP + diphosphate
<b>iML1515</b>	3.2.2.10 (b2795)	cytosine + D-ribose-5-phosphate $\rightleftharpoons$ CMP + H <sub>2</sub> O
<b>Predicted</b>	1.1.1.79 (b1033)	2-oxoglutarate + NADPH + H <sup>+</sup> $\rightleftharpoons$ 2-hydroxyglutarate + NADP <sup>+</sup>
<b>iML1515</b>	1.1.1.95 (b2913)	2-oxoglutarate + NADH + H <sup>+</sup> $\rightleftharpoons$ 2-hydroxyglutarate + NAD <sup>+</sup>
<b>Predicted</b>	3.6.1.7 (b0968)	bicarbonate + orthophosphate $\rightleftharpoons$ carboxyphosphate + H <sub>2</sub> O
<b>KEGG</b>	6.3.5.5 (b0032 or b0033)	bicarbonate + ATP $\rightleftharpoons$ carboxyphosphate + ADP
<b>Predicted</b>	2.8.3.10 (b0615)	acetoacetyl-CoA + citrate $\rightleftharpoons$ acetoacetate + (3S)-citryl-CoA
<b>iML1515</b>	2.8.3.8 (b2221+b2222 or b1694) or 2.8.3.9 (b2221+b2222)	acetoacetyl-CoA + acetate $\rightleftharpoons$ acetoacetate + acetyl-CoA

5  
 6  
 7 C2 consists of 10 reactions known to be in *E. coli* but missing from the iML1515 model. The  
 8 first predicted reaction is the aminoacyltransferase reaction between L-glutamate and  $\gamma$ -glutamyl-  
 9  $\beta$ -cyanoalanine due to EC 2.3.2.2 ( $\gamma$ -glutamyltransferase) (**Fig. 3A**). The second is a predicted  
 10 ligase reaction between L-glutamic acid and THF to form/consume THF-L-glutamic acid by EC

1 6.3.2.17 (tetrahydrofolate synthase) (**Fig. 3B**). The third is an acyltransferase transformation  
2 between propanoyl-CoA and 2-methylacetoacetyl-CoA catalyzed by EC 2.3.1.9 (acetoacetyl-  
3 CoA thiolase) (**Fig. 3C**). Fourth, PROXIMAL predicted the phosphotransferase reaction between  
4 of D-ribulose-5-phosphate and D-ribulose-1,5-bisphosphate by EC 2.7.1.19  
5 (phosphoribulokinase) (**Fig. 3D**). The fifth predicted reaction known to be in *E. coli* is the redox  
6 transformation of D-gluconic acid to 2-keto-D-gluconic acid by EC 1.1.1.215 (gluconate 2-  
7 dehydrogenase) (**Fig. 3E**). The workflow also predicted glycosyltransferase transformation of 5-  
8 amino-4-imidazolecarboxamide to/from 1-(5'-phosphoribosyl)-5-amino-4-  
9 imidazolecarboxamide by EC 2.4.2.7 (AMP pyrophosphorylase) (**Fig. 3F**). The seventh  
10 predicted reaction is the transformation between pyruvate and 4-hydroxy-2-oxoglutarate by EC  
11 4.1.3.24 (**Fig. 3G**). The eighth reaction is catalyzed by EC 2.4.2.10 to transform guanine to/from  
12 GMP (**Fig. 3H**). Also, PROXIMAL predicted the transformation between glycerate and tartrate  
13 by EC 4.1.1.73 (**Fig. 3I**). Lastly, bicarbonate is transformed to/from carboxyphosphate by EC  
14 3.6.1.7 (**Fig. 3J**).

15  
16 C3 consists of five predicted reactions that are not documented in *E. coli* but are known in other  
17 organisms. The first of these, the transformation between pyruvate and 4-carboxy-4-hydroxy-2-  
18 oxoadipate (**Fig. 4A**) catalyzed by EC 4.1.3.17 (HMG aldolase), is present in many organisms,  
19 including bacteria, as part of the benzoate degradation pathway (KEGG R00350). The  
20 transformation is predicted to occur in *E. coli* due to EC 4.1.3.34 (citryl-CoA lyase). Both EC  
21 4.1.3.17 and EC 4.1.3.34 are lyases enzymes that form carbon-carbon bonds. 4-Carboxy-4-  
22 hydroxy-2-oxoadipate is known to be formed/consumed by EC 4.2.1.80 (2-keto-4-pentenoate  
23 hydratase) in *E. coli* (KEGG R04781). Another predicted reaction is the (de)aminating redox

1 transformation between L-histidine and imidazol-5-yl-pyruvate, catalyzed by EC 1.4.1.4  
2 (glutamate dehydrogenase) (**Fig. 4B**). Imidazol-5-yl-pyruvate is not known to be produced in any  
3 other way in *E. coli*, according to ECMDB and KEGG databases. The transformation of L-  
4 histidine to/from imidazol-5-yl-pyruvate is known to occur in the bacterium *Delftia acidovorans*  
5 by EC 2.6.1.38 (histidine transaminase) [47]. C3 also includes the predicted aryltransferase  
6 reaction between geranyl diphosphate and geranyl hydroxybenzoate by EC 2.5.1.39 (4-  
7 hydroxybenzoate transferase) (**Fig. 4C**). While the general reaction of all-*trans*-polyprenyl  
8 diphosphate to 4-hydroxy-3-polyprenylbenzoate is known to occur in *E. coli*, the specific  
9 transformation between geranyl diphosphate to geranyl hydroxybenzoate is known to occur in  
10 plants as part of shikonin biosynthesis, by EC 2.5.1.93 (4-hydroxybenzoate geranyltransferase)  
11 [48]. The fourth predicted reaction is the redox transformation between D-alanine and 2-  
12 aminoacrylic acid (**Fig. 4D**). This reaction is predicted to be catalyzed by EC 1.3.1.98 (UDP-*N*-  
13 acetylmuramate dehydrogenase). While 2-aminoacrylic acid is not known to be produced in *E.*  
14 *coli* in any other way, the transformation between D-alanine and 2-aminoacrylic acid occurs in  
15 other organisms such as *Staphylococcus aureus* [49]. Lastly, our workflow predicts the  
16 transformation between phenylpyruvate and phenyllactate by EC 1.1.1.100 (**Fig. 4E**). This  
17 transformation is known to occur in plants by EC 1.1.1.237 [50].

18  
19 C4 consists of four predicted reactions that are not currently catalogued in KEGG for any  
20 organism (**Fig. 5**). The first reaction (**Fig. 5A**) is the oxidoreductive interconversion between  
21 aminomalonate and L-serine by EC 1.1.1.23 (histidinol dehydrogenase). There is one reaction  
22 (KEGG R02970) catalyzed by EC 2.6.1.47 (L-alanine:oxomalonate aminotransferase) that  
23 produces aminomalonate; but it is not a redox reaction and is associated with rat and silkworm,

1 not *E. coli* [51]. The second is a hydrolytic decarboxylation reaction between *N*-acetylputrescine  
2 and *N*-acetylornithine (**Fig. 5B**) predicted to be catalyzed by EC 4.1.1.36 (PPC decarboxylase).  
3 The product, *N*-acetylputrescine, is involved in a number of enzymatic reactions – ECs 1.4.3.4  
4 (monoamine oxidase), 2.3.1.57 (spermidine acetyltransferase), and 3.5.1.62 (acetylputrescine  
5 deacetylase) – in many organisms that include both eukaryotes and bacteria [16]. The third  
6 reaction in this category is the hydrolytic decarboxylation reaction between 3-ureidopropionate  
7 and *N*-carbamoyl-L-aspartate also catalyzed by EC 4.1.1.36 (PPC decarboxylase). 3-  
8 Ureidopropionate is present in eukaryotes and bacteria (but not *E. coli*) and is involved in  
9 reactions catalyzed by ECs 3.5.1.6 ( $\beta$ -ureidopropionase) and 3.5.2.2 (dihydropyrimidinase). The  
10 last reaction is the transformation between D-gluconic acid and D-galactarate by EC 1.1.1.23. D-  
11 Galactarate is involved in reactions catalyzed by 4.2.1.158 that is present in *Oceanobacillus*  
12 *ihayensis* [52].

13

## 14 **Discussion**

15 Current practices for reconstructing genome-scale metabolic models, which are derived using  
16 sequencing and functional annotation, can be improved by utilizing metabolomics data.  
17 However, directly utilizing metabolomics measurements to augment existing models is  
18 challenging. Not every metabolite is measurable due to limited resolution and fidelity of mass  
19 spectrometry instruments. Further, assigning chemical identities to measured metabolites  
20 remains a challenge. Even if new metabolites are identified, their formation cannot be easily  
21 assigned to enzymes without significant experimental effort involving either genetic or  
22 biochemical screens. Additionally, metabolomics data alone cannot differentiate reactions  
23 catalyzed by different enzymes yet between the same substrates-product pairs without additional

1 experimental efforts. Computational tools and workflows, as presented in this paper, can  
2 significantly guide such studies and aid in metabolic model construction and augmentation based  
3 on metabolomics data.

4  
5 The workflow we developed here is designed to identify metabolites that can form due to  
6 promiscuous enzymatic activity within a specific model organism. Further, the workflow  
7 provides balanced reactions to document such enzymatic activities. We utilized PROXIMAL  
8 [39], which first identifies patterns of structural transformations associated with enzymes in the  
9 biological sample and then applies these transformations to known sample metabolites to predict  
10 putative metabolic products. Using PROXIMAL in this way allows attributing putative  
11 metabolic products to specific enzymatic activity and deriving balanced biochemical reactions  
12 that capture the promiscuous activity. Using PROXIMAL offers an additional advantage – the  
13 derived promiscuous transformations are specific to the sample under study and are not limited  
14 to hand-curated biotransformation operators as in prior works [33, 34]. PROXIMAL therefore  
15 allows exploration of a variety of biotransformations that are commensurate with the  
16 biochemical diversity of the biological sample. The EMMA workflow, which utilized  
17 PROXIMAL, was previously developed to engineer a candidate set from a metabolic model for  
18 metabolite identification [53]. EMMA did not aim to augment existing metabolic models or  
19 derive balanced reactions as utilized in this study.

20  
21 Future experimental and computational efforts can further advance this work. Experimentally,  
22 the list of putative products generated by PROXIMAL but not documented in any metabolomics  
23 databases can be used as a resource to identify as yet unidentified metabolites. Experimental

1 validation of reactions in the C1, C3 and C4 categories would provide further evidence of the  
2 suggested reactions, and would provide a means for expanding existing databases such as KEGG  
3 and EcoCyc. Computationally, PROXIMAL can be upgraded to consider enzymes that act on  
4 more than one Reaction Center (R) within a metabolite (e.g. transketolase). This would produce  
5 multiple operators per reaction and generate a more comprehensive list of putative reactions and  
6 products. When applying PROXIMAL, we did not consider whether products of promiscuous  
7 reactions can themselves act as new substrates for promiscuous reactions. This is due to the large  
8 number of putative products. We are currently developing machine learning techniques to  
9 improve the prediction accuracy of PROXIMAL.

10

## 11 **Conclusion**

12 This study investigates creating Extended Metabolic Models (EMMs) through the augmentation  
13 of existing metabolic models with reactions due to promiscuous enzymatic activity. Our  
14 workflow, EMMA, first utilizes PROXIMAL to predict putative metabolic products, and then  
15 compares these products against metabolomics data. EMMA was applied to iML1515, the  
16 genome-scale model of *E. coli* MG1655. PROXIMAL generated 1,875 biochemical operators  
17 based on reactions in iML1515 and predicted 1,368 derivatives of 106 high-concentration  
18 metabolites. To validate these products, EMMA compared the set of putative derivatives with the  
19 set of metabolites documented in ECMDB as part of *E. coli* metabolism. For the overlapping set,  
20 we generated corresponding atom-balanced reactions by adding suitable cofactors and/or co-  
21 substrates to the substrate-derivative pair suggested by PROXIMAL. The balanced reactions  
22 were compared with data recorded in EcoCyc and KEGG. Our workflow generated a list of 23  
23 new reactions that should be utilized to extend the iML1515 model, including parallel reactions

1 between existing metabolites, novel routes to existing metabolites, and new paths to new  
2 metabolites. Importantly, this study is foundational in providing a systemic way of coupling  
3 computational predictions with metabolomics data to explore the complete metabolic repertoire  
4 of organisms. The described workflow can be applied to any organism utilizing its metabolic  
5 model to predict sample-specific promiscuous enzymatic byproducts. Applying this workflow to  
6 other biological samples and their metabolomics data promise to enhance our understanding of  
7 natural, synthetic, and xenobiotic metabolism.

8

## 9 **Methods**

10 The EMMA workflow was customized to augment the *E. coli* iML1515 model based on the  
11 availability of the metabolic measurements in ECMDB, and the availability of cataloged  
12 reactions and metabolites for *E. coli* in other databases (EcoCyc and KEGG) (**Fig. 6**). The  
13 iML1515 model consists of 1,877 metabolites, 2,712 reactions and 1,516 genes. Our workflow  
14 consists of the following three steps.

15

### 16 **Step 1 – Predict promiscuous products using PROXIMAL**

17 EMMA used PROXIMAL to predict putative products that can be added to the model.  
18 PROXIMAL utilizes RDM patterns [40] specific to the model's reactions to create lookup tables  
19 that map reaction centers to structural transformation patterns. An RDM pattern specifies local  
20 regions of structural similarities/differences for reactant-product pairs based on a given  
21 biochemical reaction. An RDM pattern consists of three parts: i) A Reaction Center (R) atom  
22 exists in both the substrate and reactant molecule and is the center of the molecular  
23 transformation. ii) Difference Region (D) atoms are adjacent to the R atom and are distinct

1 between substrate and product. iii) Matched Region (M) atoms are adjacent to the R atom but  
2 remain unmodified by the transformation. All atoms are labelled using KEGG atom types [54].  
3 PROXIMAL constructs a lookup table of all possible biotransformations that can occur due to  
4 promiscuous activity of enzymes based on the RDM patterns of reactions catalyzed by enzymes  
5 associated with genes in the iML1515 gene list. The “key” in the lookup table consisted of the R  
6 and M atom(s) in the reactant, while the “value” is the R and D atom(s) in the product. RDM  
7 patterns were initially available through the (RPAIR) database, but they are now catalogued in  
8 KEGG’s RClass database. The biotransformation operators in the lookup table were then applied  
9 to model metabolites. The outcome of this step is a list of predicted products due to putative  
10 enzymatic activity.

11

## 12 **Step 2 – Compare predicted products with metabolomics dataset**

13 Metabolites predicted by PROXIMAL were compared with measured metabolic data in  
14 ECMDB. ECMDB contains 3,760 metabolites detected in *E. coli* strain K-12 and related  
15 information such as reactions, enzymes, pathways, and other properties. This information was  
16 either collected from resources and databases such as EcoCyc, KEGG, EchoBase [55], UniProt  
17 [56, 57], YMDB [58], and CCDB [59], or from literature, or validated experimentally by the  
18 creators of ECMDB. Partial information about metabolites such as KEGG compound IDs,  
19 metabolites cell location, and chemical formulas is provided in ECMDB.

20

21 For each putative product, a mol file was generated and then converted to a SMILES string using  
22 Pybel [60], a python wrapper for the chemical toolbox Open Babel [61]. Based on the SMILES  
23 string, we initially retrieved the corresponding PubChem ID and InchiKey from PubChem using

1 Pybel. To ensure consistency, we confirmed that retrieved PubChem IDs and InchiKeys of  
2 PROXIMAL predicted metabolites matched the corresponding entries in ECMDB. During this  
3 process, we noted some discrepancies. In some cases, the information retrieved from PubChem,  
4 such as InchiKeys did not match those in ECMDB. In cases of a mismatch, we sought additional  
5 information to confirm metabolite identities of ECMDB products. We utilized the values of the  
6 CAS ID, BioCyc ID, Chebi ID and KEGG ID fields to retrieve PubChem IDs using Pybel. The  
7 retrieved PubChem IDs are used to determine the ID through a majority vote. For example, if the  
8 PubChem ID associated with InchiKey, KEGG ID and CAS ID matched, but did not match the  
9 PubChem ID provided in ECMDB, then we considered the one retrieved by Pybel as the correct  
10 PubChem ID. Out of 3,760 metabolites in ECMDB, we identified 3,397 metabolites with  
11 consistent information with data retrieved from PubChem. Once PubChem IDs were identified  
12 for ECMDB metabolites, we compared our predicted metabolites against ECMDB metabolites  
13 using PubChem IDs.

14

### 15 **Step 3 – Curation of stoichiometric reactions**

16 If a metabolite predicted by PROXIMAL was in ECMDB, then steps 1 and 2 resulted in the  
17 identification of a *verifiable* predicted promiscuous transformation of an *E. coli* metabolite.  
18 Each predicted transformation was manually examined and compared against the RDM pattern  
19 causing the transformation. Transformations were discarded if they seemed infeasible, if the  
20 substrate was a cofactor, or if the RPAIR entry associated with the PROXIMAL operator  
21 required the presence of more than one Reaction Center (R). For each valid verifiable predicted  
22 transformation by PROXIMAL, we developed a new reaction by examining the reaction(s)  
23 template associated with the enzymatic transformation and adding suitable cofactors to the

1 reactant and product of the biotransformation identified. The set of developed balanced reactions,  
2 where the added cofactors to a reaction caused the number of atoms of reactants and products to  
3 match on both sides of the reaction, are then compared to reactions recorded in EcoCyc, KEGG,  
4 or the literature.

5  
6 The outcomes were divided into four categories. C1 reactions consisted of metabolites predicted  
7 by PROXIMAL that are already in iML1515 but catalyzed by different enzymes than the ones  
8 already listed in the model. These reactions reflect promiscuous activity that enabled the same  
9 biotransformation catalyzed by a different gene in the model. C2 reactions already existed in  
10 EcoCyc and/or KEGG but not in iML1515. This reflected a curation problem where some  
11 reactions were not included in the iML1515 model. C3 reactions were not in EcoCyc but  
12 documented in KEGG for other organisms. C4 reactions did not exist in either EcoCyc nor in  
13 KEGG. These reactions were thus novel reactions that have not been reported in the literature.

## 14 15 **Declarations**

### 16 **Ethics approval and consent to participate**

17 Not applicable

### 18 **Consent for publication**

19 Not applicable

### 20 **Availability of data and material**

21 The *E. coli* iML1615 model is available on the BiGG database and can be found at  
22 <http://bigg.ucsd.edu/models/iML1515>. Data from ECMDB can be directly downloaded from the

1 ECMDDB website. A full list of derivatives that were predicted by PROXIMAL and had a  
2 chemical ID in PubChem can be found in **Supplementary File 1**.

### 3 **Competing interests**

4 Not applicable

### 5 **Funding**

6 This work is funded under NSF grant #1421972 and NIH grants 1DP2HD091798 and  
7 1R03CA211839-01.

### 8 **Author's contribution**

9 SH conceived the EMMA concept. SA developed the EMMA workflow. EC curated the results.  
10 VP verified the analysis. NN and SH supervised the work done through the development of the  
11 workflow and data curation. Manuscript was written by SA and EC, reviewed by VP, and revised  
12 by NN and SH.

### 13 **Acknowledgments**

14 Not applicable

15

### 16 **References**

- 17 1. Lee SK, Chou H, Ham TS, Lee TS, Keasling JD. Metabolic engineering of microorganisms for  
18 biofuels production: from bugs to synthetic biology to fuels. *Current opinion in biotechnology*.  
19 2008;19(6):556-63.
- 20 2. Trantas EA, Koffas MA, Xu P, Ververidis F. When plants produce not enough or at all: metabolic  
21 engineering of flavonoids in microbial hosts. *Frontiers in plant science*. 2015;6:7.
- 22 3. George KW, Alonso-Gutierrez J, Keasling JD, Lee TS. Isoprenoid drugs, biofuels, and  
23 chemicals—artemisinin, farnesene, and beyond. *Biotechnology of Isoprenoids: Springer*; 2015. p. 355-  
24 89.
- 25 4. Singh R, White D, Demirel Y, Kelly R, Noll K, Blum P. Uncoupling fermentative synthesis of  
26 molecular hydrogen from biomass formation in *Thermotoga maritima*. *Appl Environ Microbiol*.  
27 2018;84(17):e00998-18.
- 28 5. Singh R, Tevatia R, White D, Demirel Y, Blum P. Comparative kinetic modeling of growth and  
29 molecular hydrogen overproduction by engineered strains of *Thermotoga maritima*. *International Journal*  
30 *of Hydrogen Energy*. 2019.

- 1 6. Du J, Shao Z, Zhao H. Engineering microbial factories for synthesis of value-added products.  
2 *Journal of industrial microbiology & biotechnology*. 2011;38(8):873-90.
- 3 7. Furusawa C, Horinouchi T, Hirasawa T, Shimizu H. Systems metabolic engineering: the creation  
4 of microbial cell factories by rational metabolic design and evolution. *Future Trends in Biotechnology*:  
5 Springer; 2012. p. 1-23.
- 6 8. Davy AM, Kildegaard HF, Andersen MR. Cell factory engineering. *Cell Systems*. 2017;4(3):262-  
7 75.
- 8 9. Lee S, Mattanovich D, Villaverde A. Systems metabolic engineering, industrial biotechnology  
9 and microbial cell factories. *BioMed Central*; 2012.
- 10 10. Burgard AP, Pharkya P, Maranas CD. Optknock: a bilevel programming framework for  
11 identifying gene knockout strategies for microbial strain optimization. *Biotechnology and bioengineering*.  
12 2003;84(6):647-57.
- 13 11. Ranganathan S, Suthers PF, Maranas CD. OptForce: an optimization procedure for identifying all  
14 genetic manipulations leading to targeted overproductions. *PLoS computational biology*.  
15 2010;6(4):e1000744.
- 16 12. Yousofshahi M, Lee K, Hassoun S. Probabilistic pathway construction. *Metabolic engineering*.  
17 2011;13(4):435-44.
- 18 13. Wu G, Yan Q, Jones JA, Tang YJ, Fong SS, Koffas MA. Metabolic burden: cornerstones in  
19 synthetic biology and metabolic engineering applications. *Trends in biotechnology*. 2016;34(8):652-64.
- 20 14. Gerstl MP, Ruckerbauer DE, Mattanovich D, Jungreuthmayer C, Zanghellini J. Metabolomics  
21 integrated elementary flux mode analysis in large metabolic networks. *Scientific reports*. 2015;5:8930.
- 22 15. Kim TY, Sohn SB, Kim YB, Kim WJ, Lee SY. Recent advances in reconstruction and  
23 applications of genome-scale metabolic models. *Current opinion in biotechnology*. 2012;23(4):617-23.
- 24 16. Saha R, Chowdhury A, Maranas CD. Recent advances in the reconstruction of metabolic models  
25 and integration of omics data. *Current opinion in biotechnology*. 2014;29:39-45.
- 26 17. Mukherjee S, Stamatis D, Bertsch J, Ovchinnikova G, Katta HY, Mojica A, et al. Genomes  
27 OnLine database (GOLD) v. 7: updates and new features. *Nucleic Acids Research*. 2018.
- 28 18. Kanehisa M, Goto S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic acids research*.  
29 2000;28(1):27-30.
- 30 19. Karp PD, Billington R, Caspi R, Fulcher CA, Latendresse M, Kothari A, et al. The BioCyc  
31 collection of microbial genomes and metabolic pathways. *Brief Bioinformatics*. 2017.
- 32 20. King ZA, Lu J, Dräger A, Miller P, Federowicz S, Lerman JA, et al. BiGG Models: A platform  
33 for integrating, standardizing and sharing genome-scale models. *Nucleic acids research*.  
34 2015;44(D1):D515-D22.
- 35 21. Sorokina M, Stam M, Médigue C, Lespinet O, Vallenet D. Profiling the orphan enzymes. *Biology*  
36 *direct*. 2014;9(1):10.
- 37 22. Raushel FM. Finding homes for orphan enzymes. *Perspectives in Science*. 2016;9:3-7.
- 38 23. Notebaart RA, Szappanos B, Kintsjes B, Pál F, Györkei Á, Bogos B, et al. Network-level  
39 architecture and the evolutionary potential of underground metabolism. *Proceedings of the National*  
40 *Academy of Sciences*. 2014;111(32):11762-7.
- 41 24. Notebaart RA, Kintsjes B, Feist AM, Papp B. Underground metabolism: network-level  
42 perspective and biotechnological potential. *Current Opinion in Biotechnology*. 2018;49:108-14.
- 43 25. Rosenberg J, Commichau FM. Harnessing Underground Metabolism for Pathway Development.  
44 *Trends in biotechnology*. 2018.
- 45 26. Hult K, Berglund P. Enzyme promiscuity: mechanism and applications. *Trends in biotechnology*.  
46 2007;25(5):231-8.
- 47 27. Khersonsky O, Roodveldt C, Tawfik DS. Enzyme promiscuity: evolutionary and mechanistic  
48 aspects. *Current opinion in chemical biology*. 2006;10(5):498-508.
- 49 28. Tawfik OK, S D. Enzyme promiscuity: a mechanistic and evolutionary perspective. *Annual*  
50 *review of biochemistry*. 2010;79:471-505.

- 1 29. Nobeli I, Favia AD, Thornton JM. Protein promiscuity and its implications for biotechnology.  
2 Nature biotechnology. 2009;27(2):157.
- 3 30. D'Ari R, Casadesús J. Underground metabolism. BioEssays. 1998;20(2):181-6.
- 4 31. Liechti G, Singh R, Rossi PL, Gray MD, Adams NE, Maurelli AT. Chlamydia trachomatis dapF  
5 encodes a bifunctional enzyme capable of both D-glutamate racemase and diaminopimelate epimerase  
6 activities. MBio. 2018;9(2):e00204-18.
- 7 32. Carbonell P, Faulon J-L. Molecular signatures-based prediction of enzyme promiscuity.  
8 Bioinformatics. 2010;26(16):2012-9.
- 9 33. Jeffryes JG, Colastani RL, Elbadawi-Sidhu M, Kind T, Niehaus TD, Broadbelt LJ, et al. MINEs:  
10 open access databases of computationally predicted enzyme promiscuity products for untargeted  
11 metabolomics. Journal of cheminformatics. 2015;7(1):44.
- 12 34. Hadadi N, Hafner J, Shajkofci A, Zisaki A, Hatzimanikatis V. ATLAS of biochemistry: a  
13 repository of all possible biochemical reactions for synthetic biology and metabolic engineering studies.  
14 ACS synthetic biology. 2016;5(10):1155-66.
- 15 35. Arora B, Mukherjee J, Gupta MN. Enzyme promiscuity: using the dark side of enzyme specificity  
16 in white biotechnology. Sustainable Chemical Processes. 2014;2(1):25.
- 17 36. Poppe L, Paizs C, Kovács K, Irimie F-D, Vértessy B. Preparation of unnatural amino acids with  
18 ammonia-lyases and 2, 3-aminomutases. Unnatural Amino Acids: Springer; 2012. p. 3-19.
- 19 37. Atsumi S, Hanai T, Liao JC. Non-fermentative pathways for synthesis of branched-chain higher  
20 alcohols as biofuels. nature. 2008;451(7174):86.
- 21 38. Song CW, Kim JW, Cho IJ, Lee SY. Metabolic engineering of Escherichia coli for the production  
22 of 3-hydroxypropionic acid and malonic acid through  $\beta$ -alanine route. ACS synthetic biology.  
23 2016;5(11):1256-63.
- 24 39. Yousofshahi M, Manteiga S, Wu C, Lee K, Hassoun S. PROXIMAL: a method for Prediction of  
25 Xenobiotic Metabolism. BMC systems biology. 2015;9(1):94.
- 26 40. Moriya Y, Shigemizu D, Hattori M, Tokimatsu T, Kotera M, Goto S, et al. PathPred: an enzyme-  
27 catalyzed metabolic pathway prediction server. Nucleic acids research. 2010;38(suppl\_2):W138-W43.
- 28 41. Monk JM, Lloyd CJ, Brunk E, Mih N, Sastry A, King Z, et al. iML1515, a knowledgebase that  
29 computes Escherichia coli traits. Nature biotechnology. 2017;35(10):904.
- 30 42. Guo AC, Jewison T, Wilson M, Liu Y, Knox C, Djoumbou Y, et al. ECMDB: the E. coli  
31 Metabolome Database. Nucleic acids research. 2012;41(D1):D625-D30.
- 32 43. Sajed T, Marcu A, Ramirez M, Pon A, Guo AC, Knox C, et al. ECMDB 2.0: A richer resource  
33 for understanding the biochemistry of E. coli. Nucleic acids research. 2015;44(D1):D495-D501.
- 34 44. Keseler IM, Collado-Vides J, Gama-Castro S, Ingraham J, Paley S, Paulsen IT, et al. EcoCyc: a  
35 comprehensive database resource for Escherichia coli. Nucleic acids research. 2005;33(suppl\_1):D334-  
36 D7.
- 37 45. Tepper N, Noor E, Amador-Noguez D, Haraldsdóttir HS, Milo R, Rabinowitz J, et al. Steady-  
38 state metabolite concentrations reflect a balance between maximizing enzyme efficiency and minimizing  
39 total metabolite load. PloS one. 2013;8(9):e75370.
- 40 46. Quinlan JR. Simplifying decision trees. International journal of man-machine studies.  
41 1987;27(3):221-34.
- 42 47. Coote J, Hassall H. The role of imidazol-5-yl-lactate-nicotinamide-adenine dinucleotide  
43 phosphate oxidoreductase and histidine-2-oxoglutarate aminotransferase in the degradation of imidazol-5-  
44 yl-lactate by Pseudomonas acidovorans. Biochemical Journal. 1969;111(2):237.
- 45 48. Mühlenweg A, Melzer M, Li S-M, Heide L. 4-Hydroxybenzoate 3-geranyltransferase from  
46 Lithospermum erythrorhizon: purification of a plant membrane-bound prenyltransferase. Planta.  
47 1998;205(3):407-13.
- 48 49. Suda S, Lawton EM, Wistuba D, Cotter PD, Hill C, Ross RP. Homologues and bioengineered  
49 derivatives of LtnJ vary in ability to form D-alanine in the lantibiotic lacticin 3147. Journal of  
50 bacteriology. 2012;194(3):708-14.

- 1 50. Häusler E, Petersen M, Alfermann AW. Hydroxyphenylpyruvate Reductase From Cell  
2 Suspension Cultures Of *Coleus Blumei* Benth. *Zeitschrift für Naturforschung C*. 1991;46(5-6):371-6.
- 3 51. Nagayama H, Muramatsu M, Shimura K. Enzymatic formation of aminomalonic acid from  
4 ketomalonic acid. *Nature*. 1958;181(4606):417.
- 5 52. Rakus JF, Kalyanaraman C, Fedorov AA, Fedorov EV, Mills-Groninger FP, Toro R, et al.  
6 Computation-facilitated assignment of the function in the enolase superfamily: a regiochemically distinct  
7 galactarate dehydratase from *Oceanobacillus iheyensis*. *Biochemistry*. 2009;48(48):11546-58.
- 8 53. Hassanpour N. *Computational Methods to Advance Directed Evolution of Enzymes and*  
9 *Metabolomics Data Analysis*: Tufts University; 2018.
- 10 54. Hattori M, Okuno Y, Goto S, Kanehisa M. Development of a chemical structure comparison  
11 method for integrated analysis of chemical and genomic information in the metabolic pathways. *Journal*  
12 *of the American Chemical Society*. 2003;125(39):11853-65.
- 13 55. Misra RV, Horler RS, Reindl W, Goryanin II, Thomas GH. Echo BASE: an integrated post-  
14 genomic database for *Escherichia coli*. *Nucleic acids research*. 2005;33(suppl\_1):D329-D33.
- 15 56. Apweiler R, Bairoch A, Wu CH, Barker WC, Boeckmann B, Ferro S, et al. UniProt: the universal  
16 protein knowledgebase. *Nucleic acids research*. 2004;32(suppl\_1):D115-D9.
- 17 57. Consortium U. UniProt: a hub for protein information. *Nucleic acids research*.  
18 2014;43(D1):D204-D12.
- 19 58. Jewison T, Knox C, Neveu V, Djoumbou Y, Guo AC, Lee J, et al. YMDB: the yeast metabolome  
20 database. *Nucleic acids research*. 2011;40(D1):D815-D20.
- 21 59. Sundararaj S, Guo A, Habibi Nazhad B, Rouani M, Stothard P, Ellison M, et al. The CyberCell  
22 Database (CCDB): a comprehensive, self-updating, relational database to coordinate and facilitate in  
23 silico modeling of *Escherichia coli*. *Nucleic acids research*. 2004;32(suppl\_1):D293-D5.
- 24 60. O'Boyle NM, Morley C, Hutchison GR. Pybel: a Python wrapper for the OpenBabel  
25 cheminformatics toolkit. *Chemistry Central Journal*. 2008;2(1):5.
- 26 61. O'Boyle NM, Banck M, James CA, Morley C, Vandermeersch T, Hutchison GR. Open Babel: An  
27 open chemical toolbox. *Journal of cheminformatics*. 2011;3(1):33.

28

## 1 **List of Figures**

2 **Fig. 1:** Decision tree for classifying reactions identified based on enzyme promiscuity. When  
3 analyzing the iML1515 *E. coli* model, reaction categories C1, C2, C3, and C4 had 4, 10, 5, and 4  
4 predicted reactions, respectively.

5  
6 **Fig. 2:** The set of four reactions belonging to Category 1 (C1). Reactions in C1 are predicted to  
7 be catalyzed by enzymes different than those in iML1515. Each of the four panels is divided into  
8 three sections I) the balanced reaction developed by our workflow indicating the reactants,  
9 products, and the promiscuous enzyme, II) the RDM pattern showing the Reaction Center (R) in  
10 red where the biotransformation occurs, and III) the native reaction catalyzed by the potentially  
11 promiscuous enzyme, as catalogued in KEGG.

12  
13 **Fig. 3:** The set of ten reactions belonging to Category 2 (C2). Reactions in C2 are associated  
14 with derivatives not present in iML1515 but are associated with *E. coli* in KEGG and/or EcoCyc.  
15 Each of the ten panels is divided into two sections I) the balanced reaction developed by our  
16 workflow, that is also documented in KEGG, indicating the reactants, products, and the  
17 promiscuous enzyme, and II) the RDM pattern showing the Reaction Center (R) in red where the  
18 biotransformation occurs.

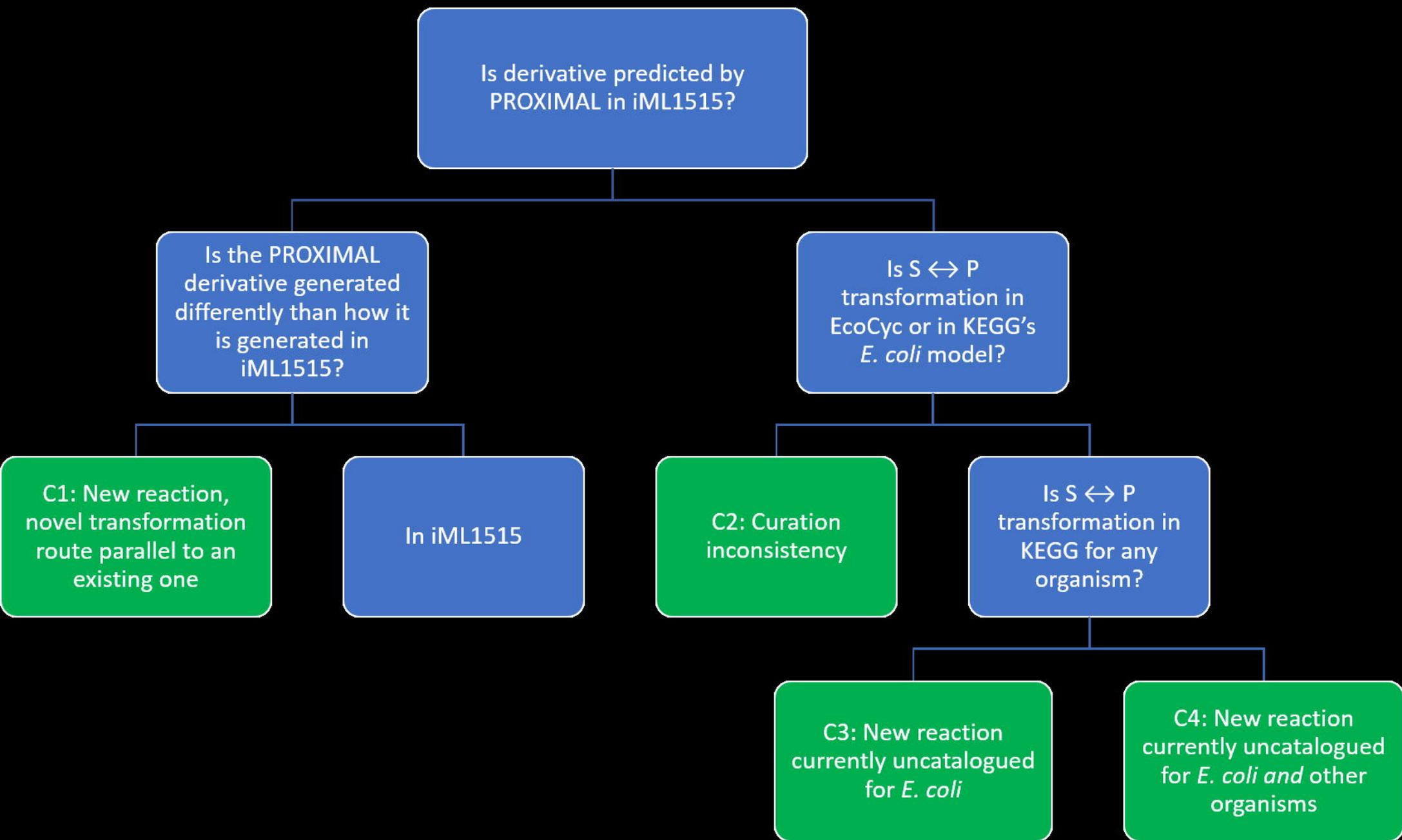
19  
20 **Fig. 4:** The set of five reactions belonging to Category 3 (C3). C3 reactions and derivatives are  
21 neither present in iML1515 nor associated with *E. coli* in KEGG and EcoCyc. However,  
22 according to KEGG, the reactions occur in other organisms. Each of the five panels is divided  
23 into three sections I) the balanced reaction developed by our workflow indicating the reactants,

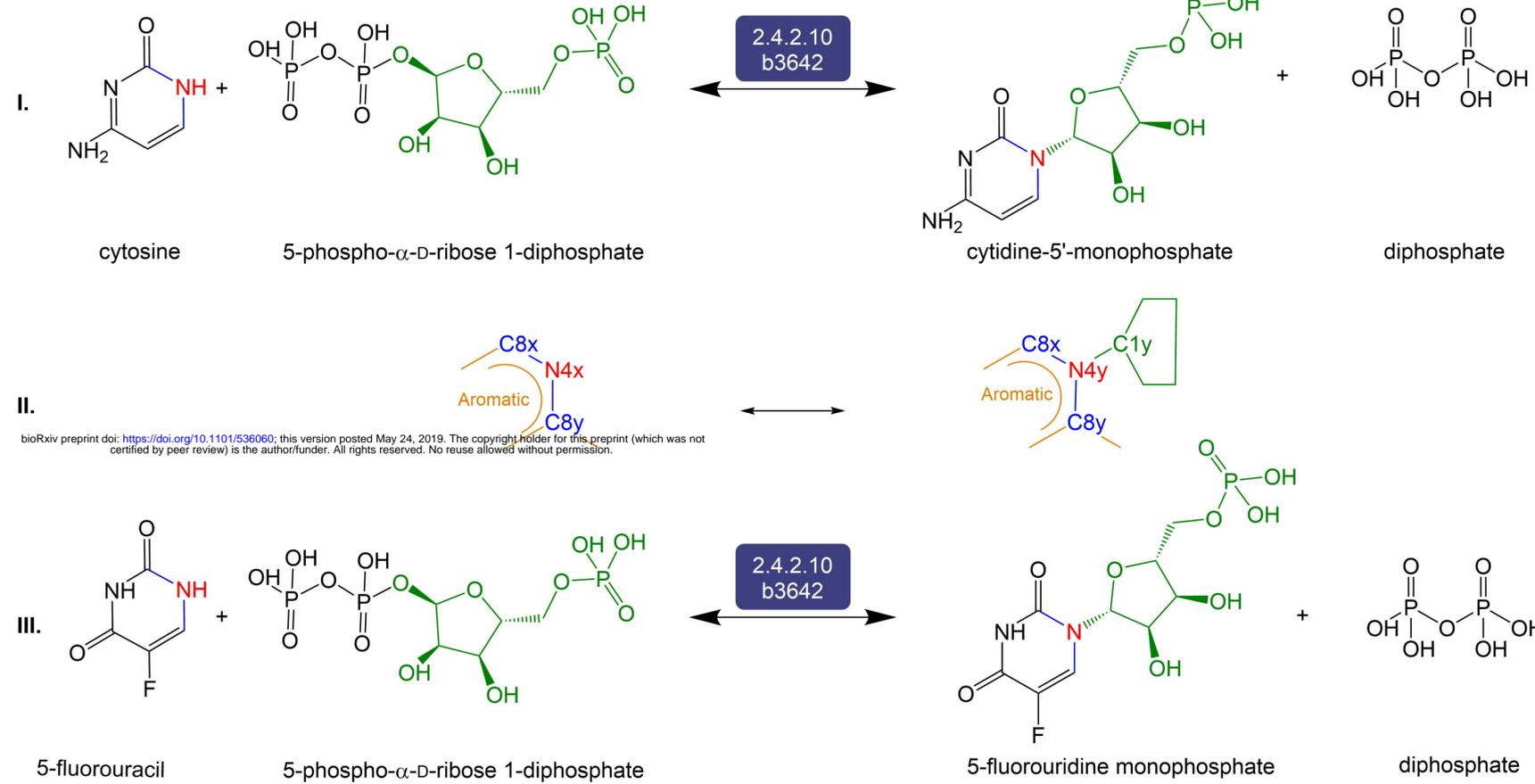
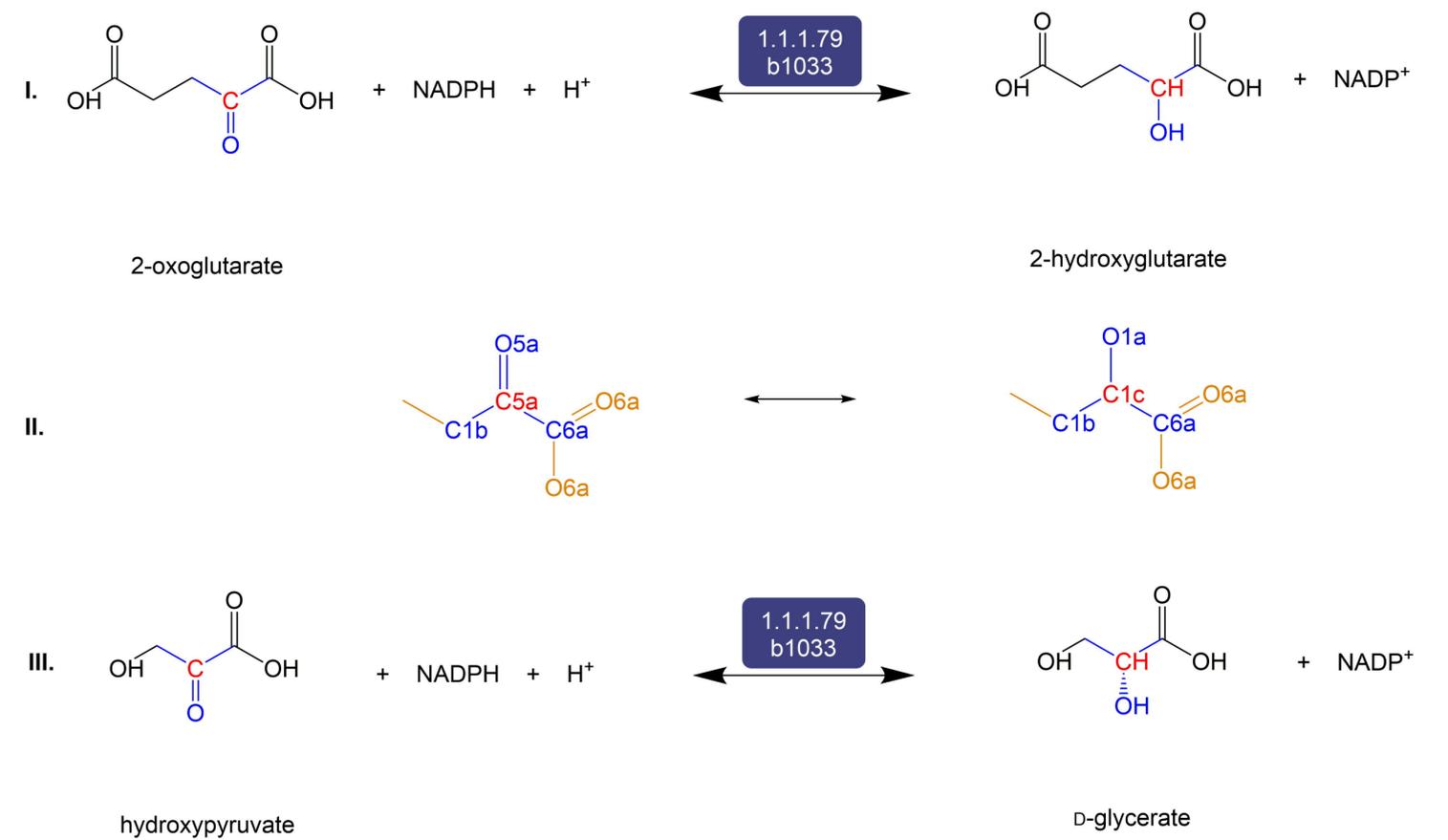
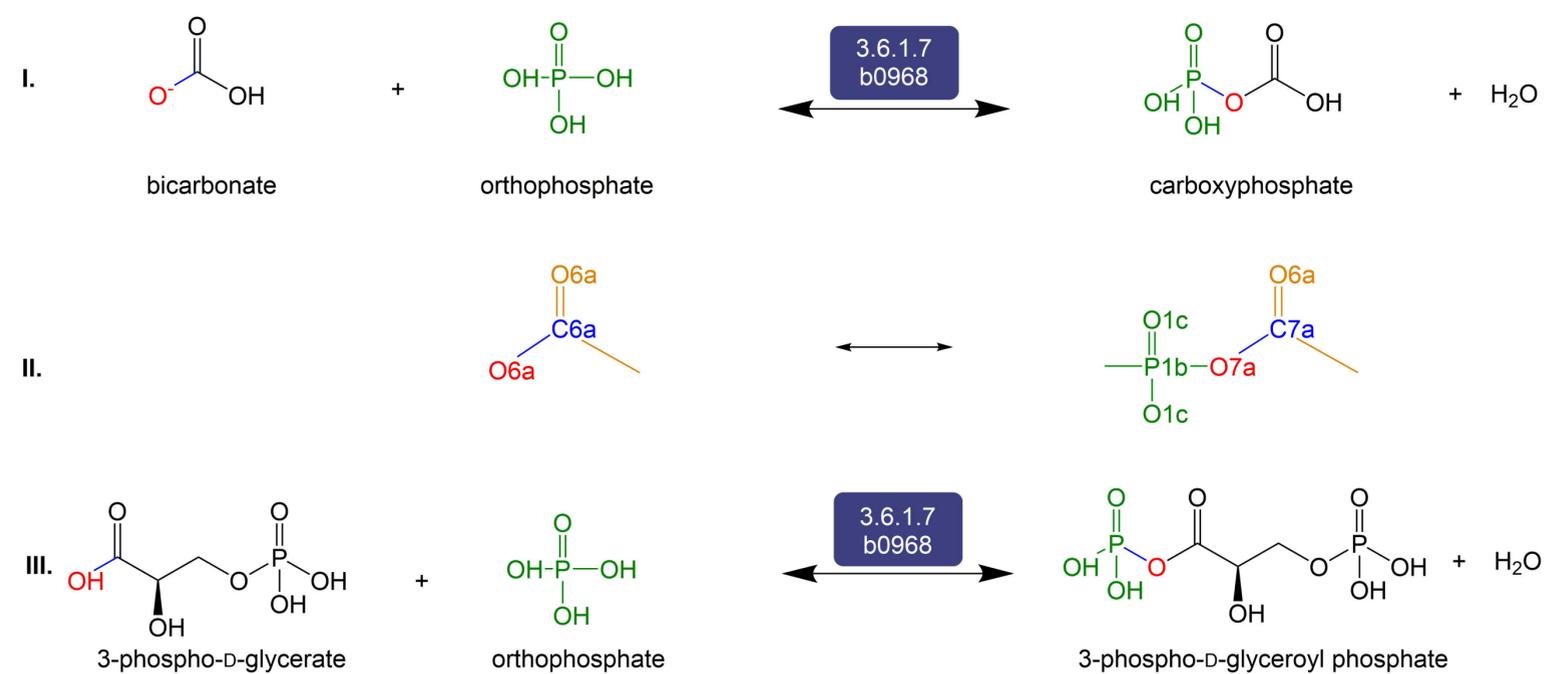
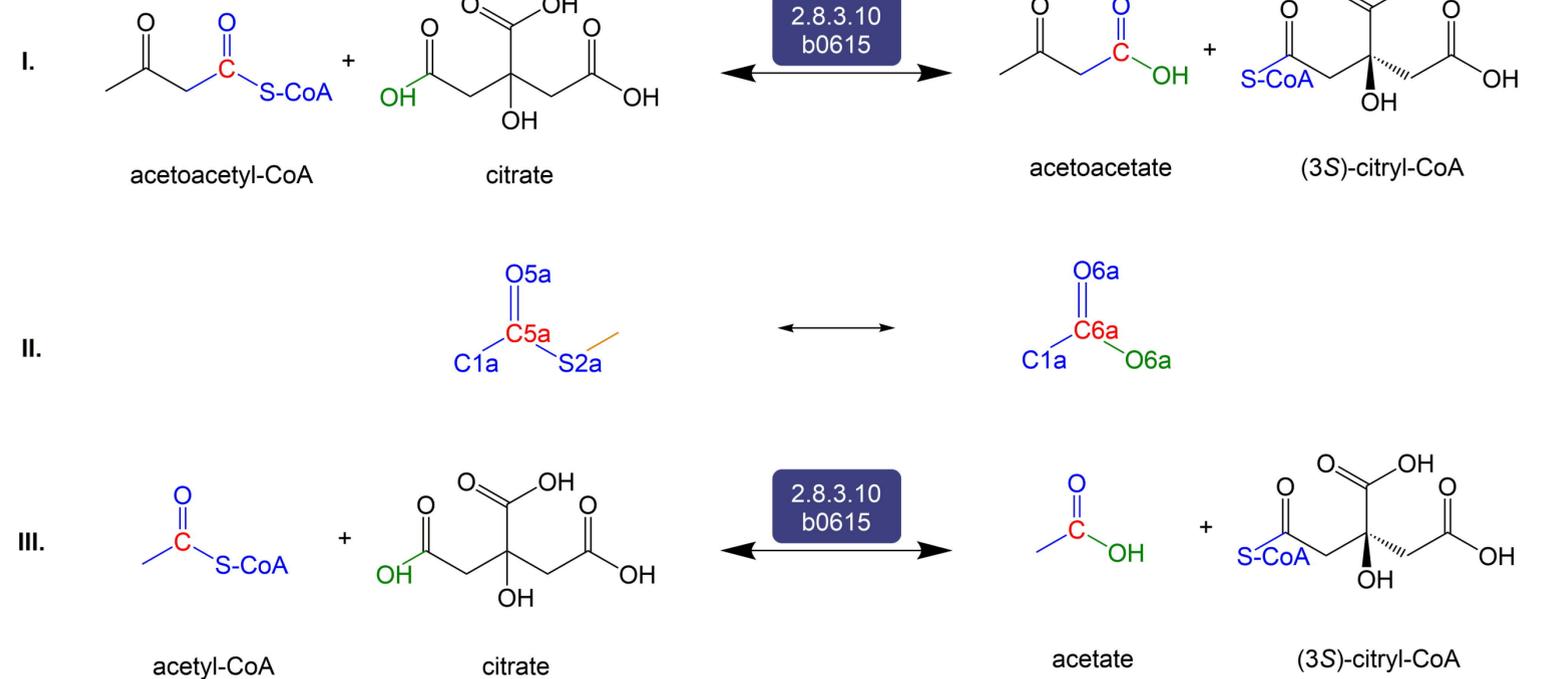
1 products, and the promiscuous enzyme, II) the RDM pattern showing the Reaction Center (R) in  
2 red, and III) the native reaction catalyzed by the potentially promiscuous enzyme, as catalogued  
3 in KEGG.

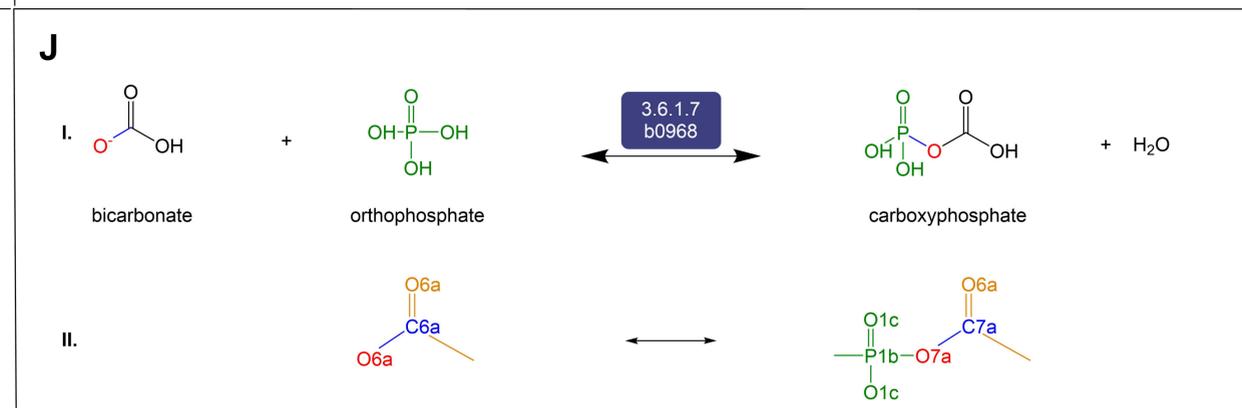
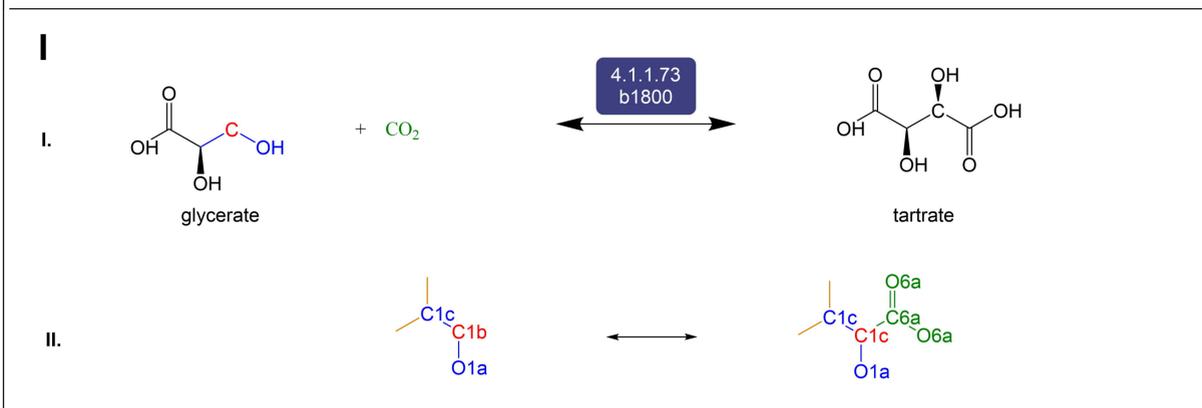
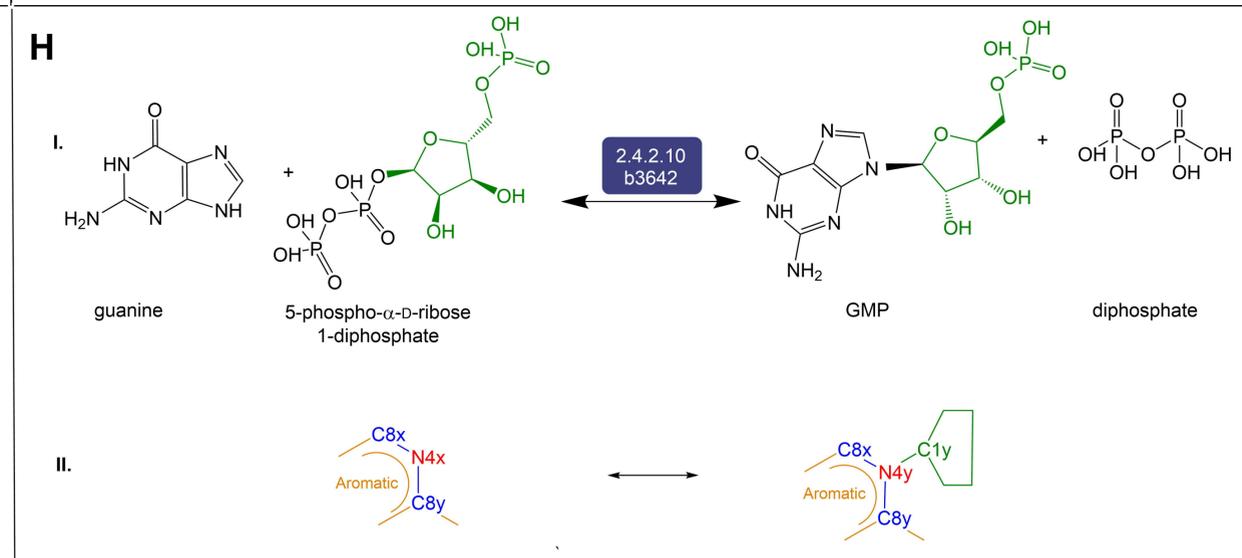
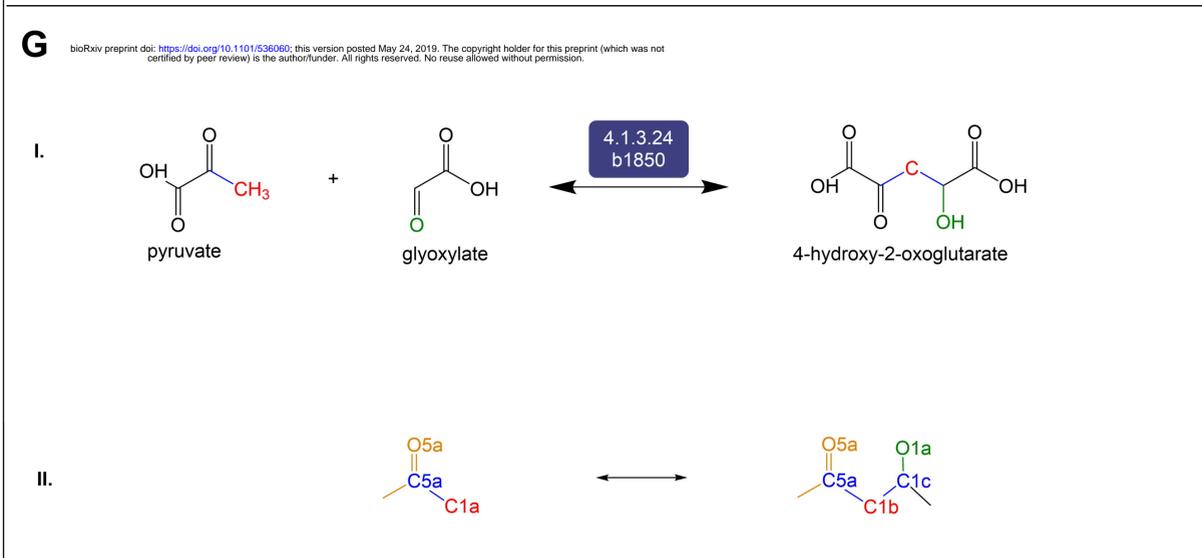
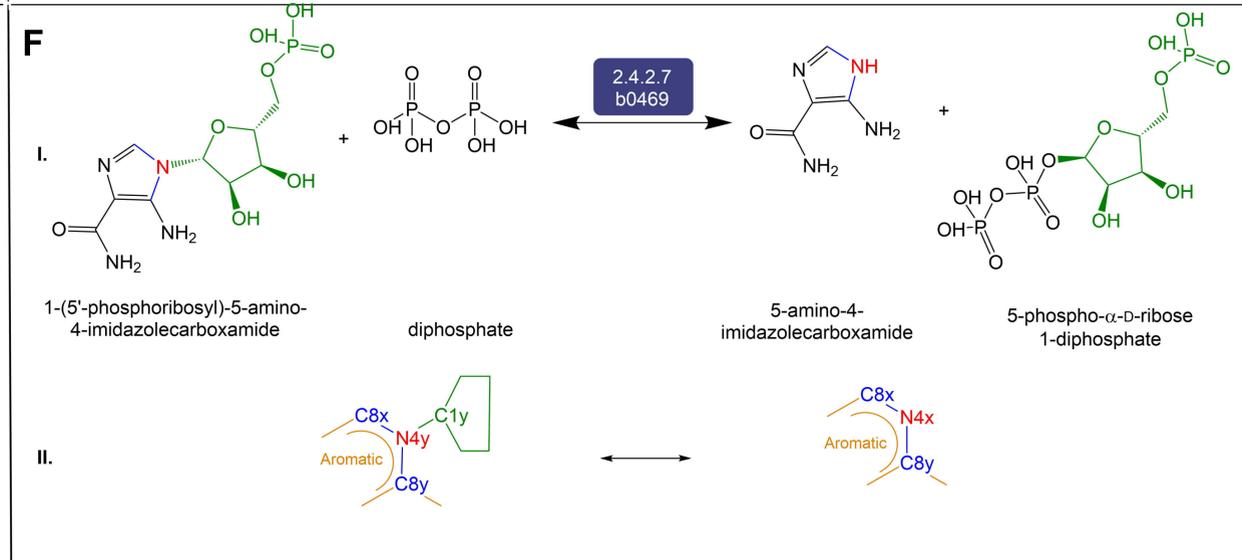
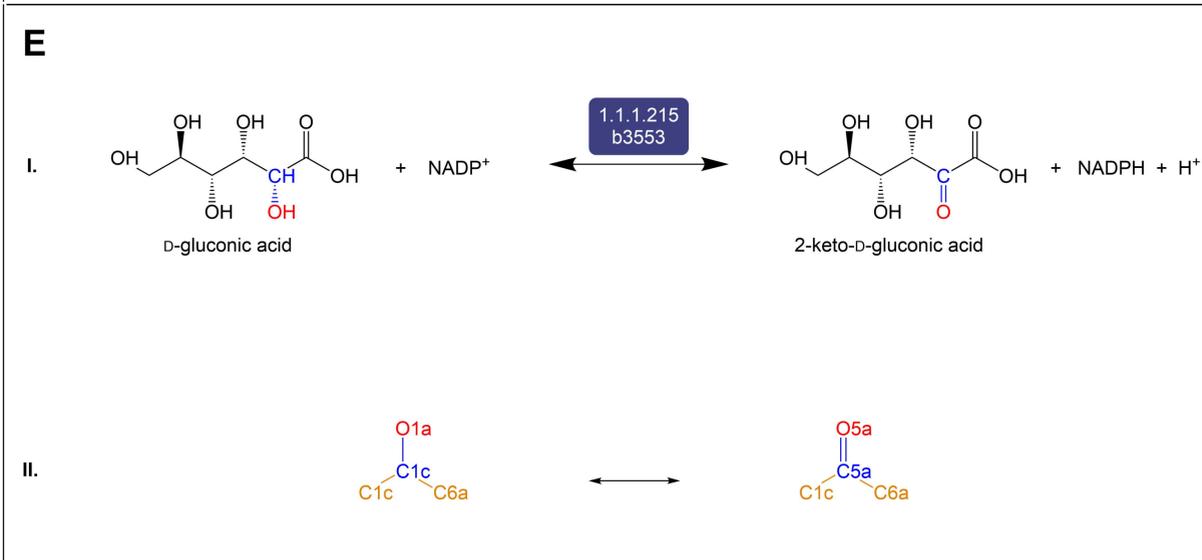
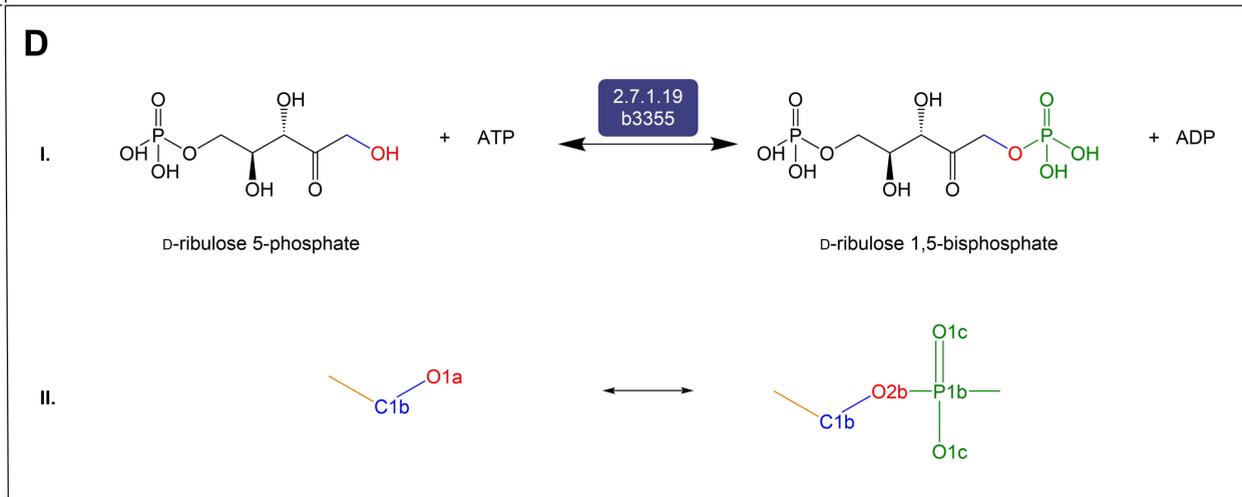
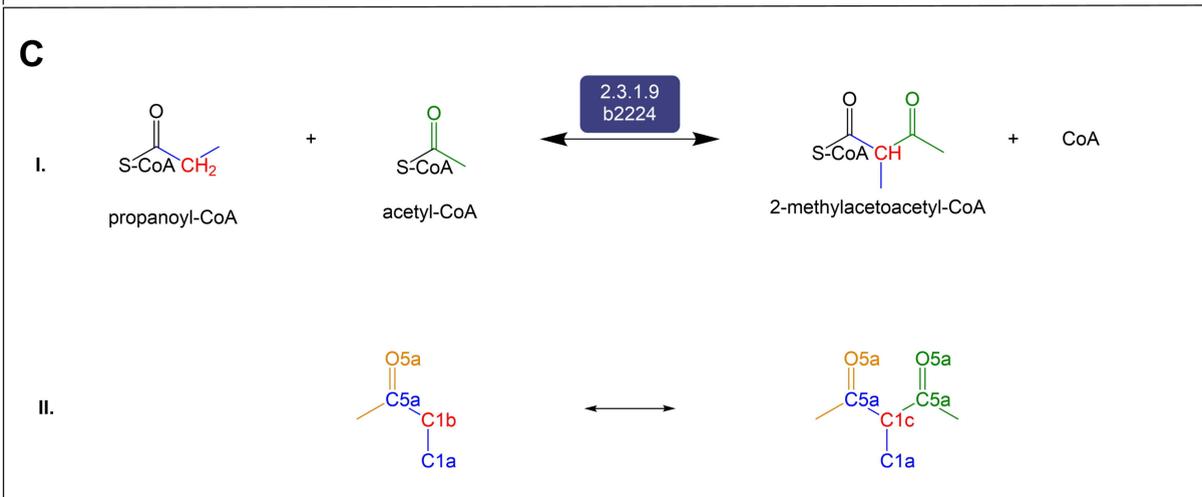
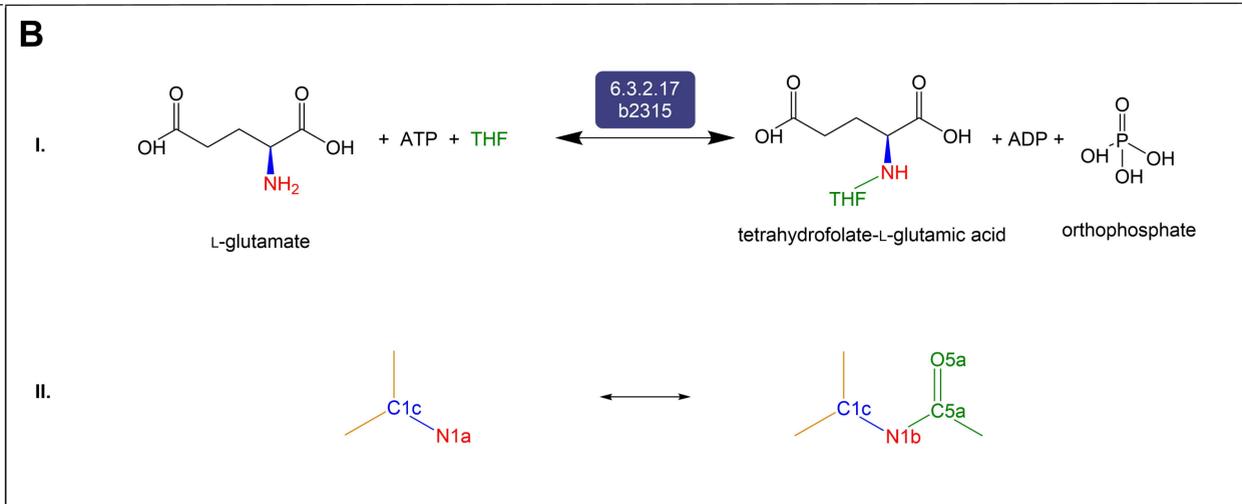
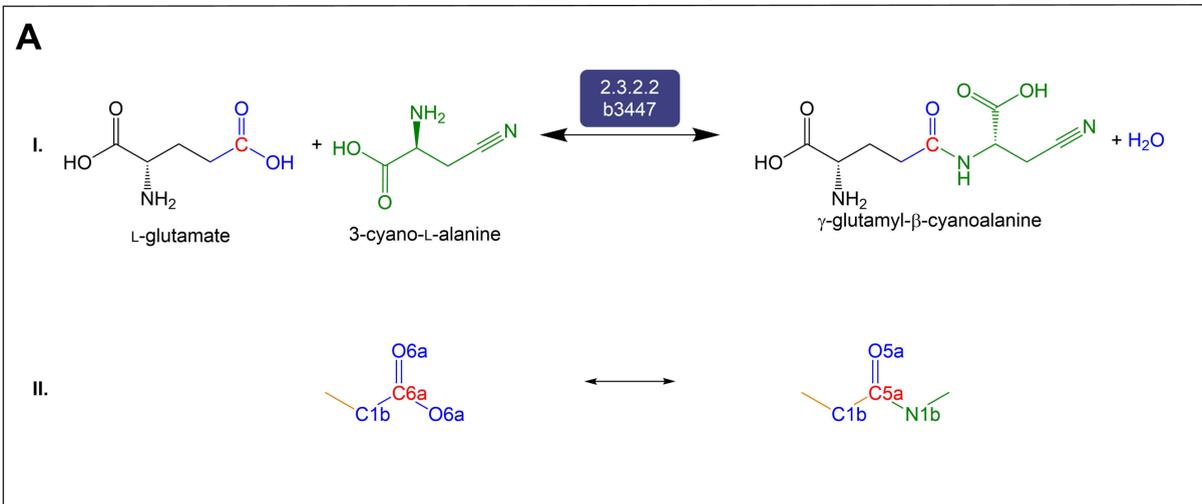
4  
5 **Fig. 5:** The set of four reactions belonging to Category 4 (C4). C4 reactions and derivatives are  
6 neither present in iML1515 nor associated with any other organism in KEGG or EcoCyc. Each  
7 of the four panels is divided into three sections I) the balanced reaction developed by our  
8 workflow indicating the reactants, products, and the promiscuous enzyme, II) the RDM pattern  
9 showing the Reaction Center (R) in red, and III) the native reaction catalyzed by the potentially  
10 promiscuous enzyme, as catalogued in KEGG.

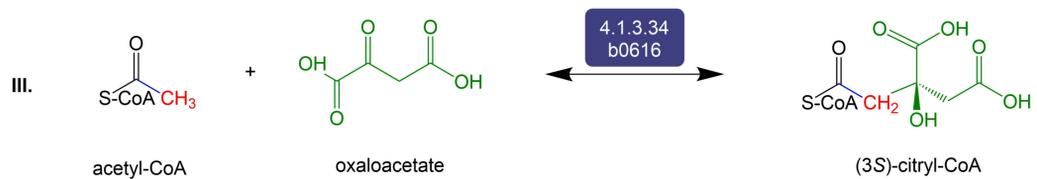
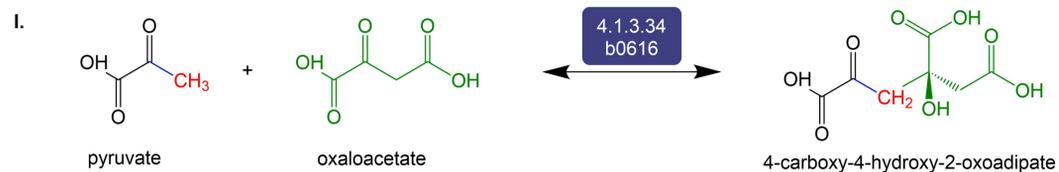
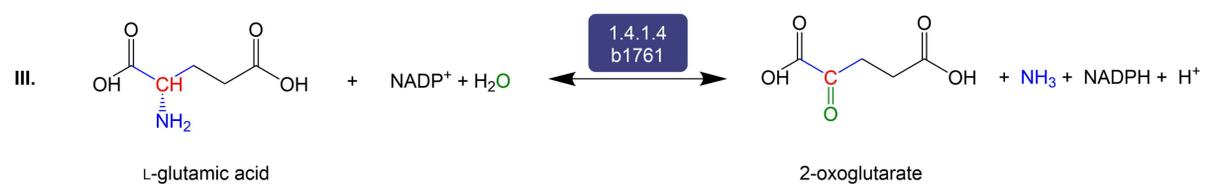
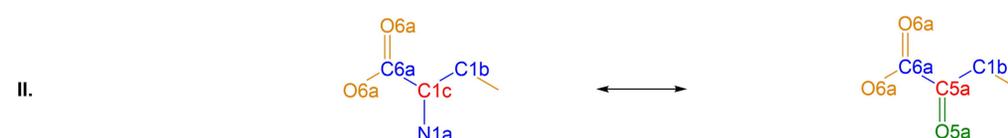
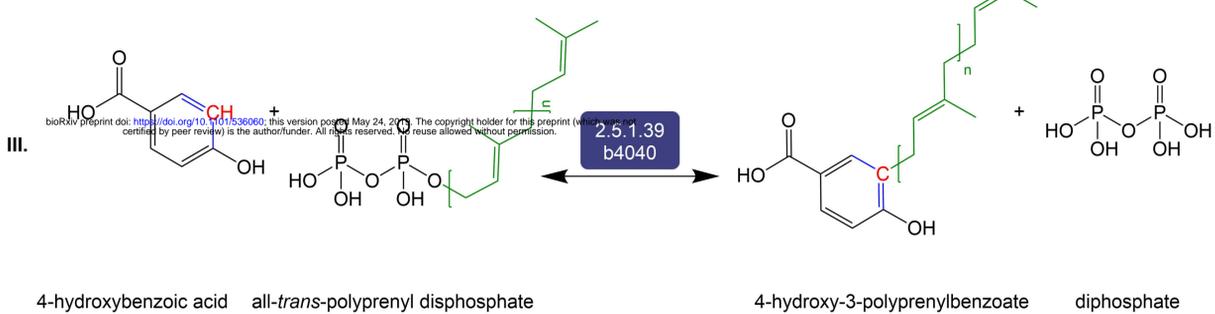
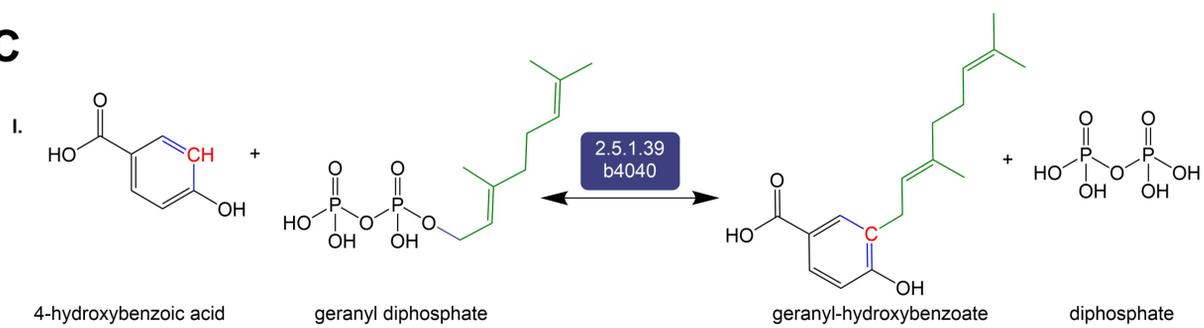
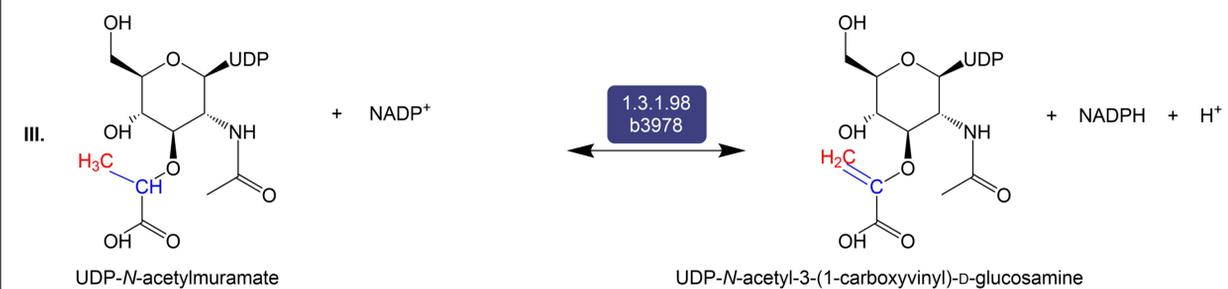
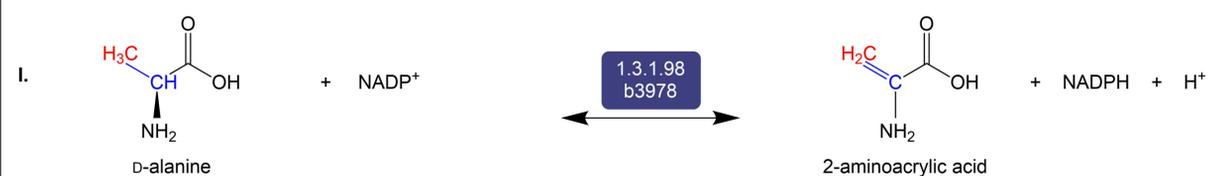
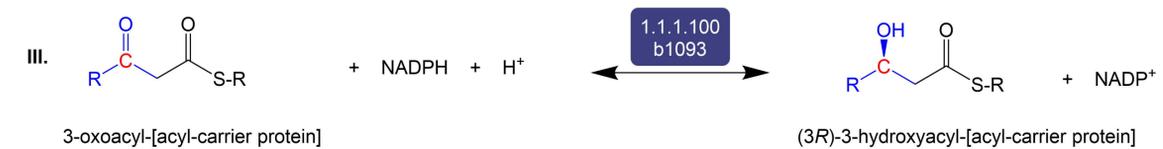
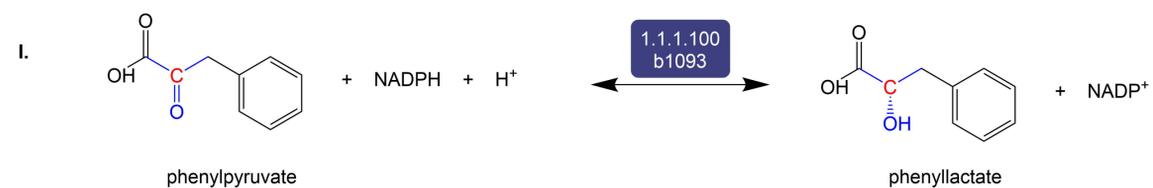
11  
12 **Fig. 6:** Main steps of EMMA workflow customized to extend the *E. coli* iML1515 model with  
13 predicted reactions. Step 1: Predict promiscuous transformations and derivatives using  
14 PROXIMAL. Step 2: Compare derivatives with measured metabolic dataset(s). Step 3: Curation  
15 and stoichiometric balancing of reactions.

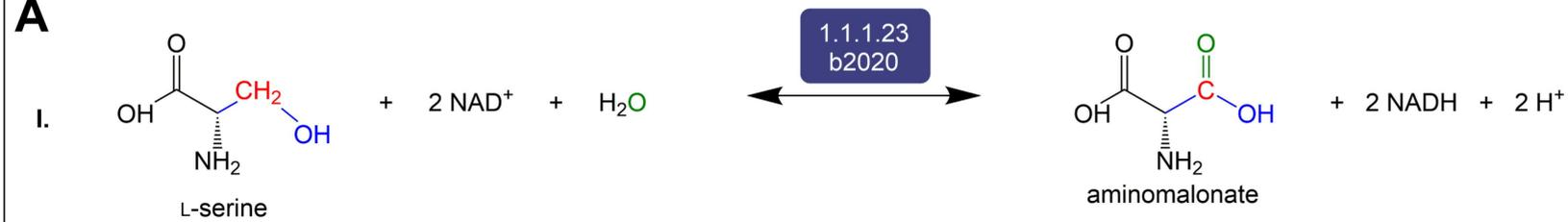
16



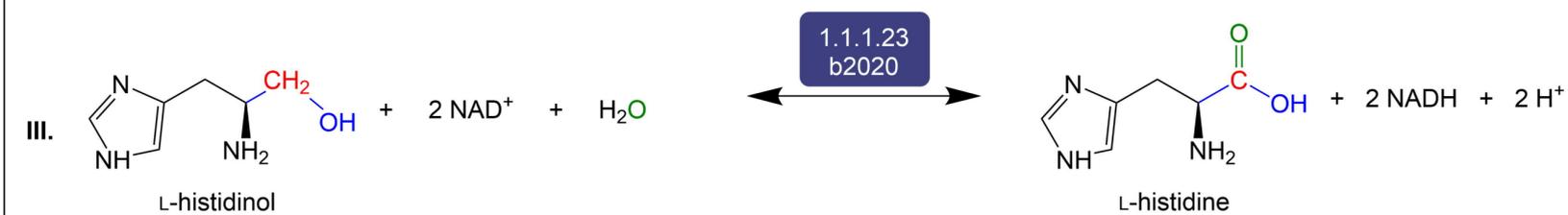
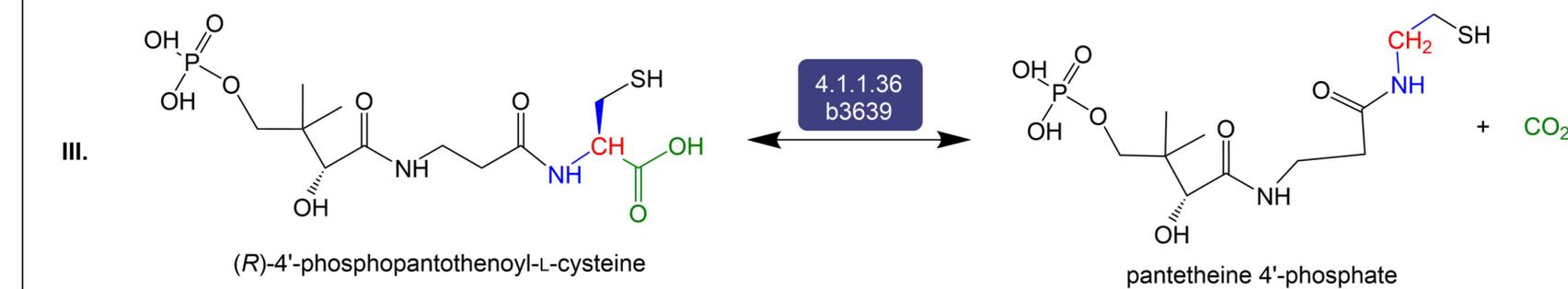
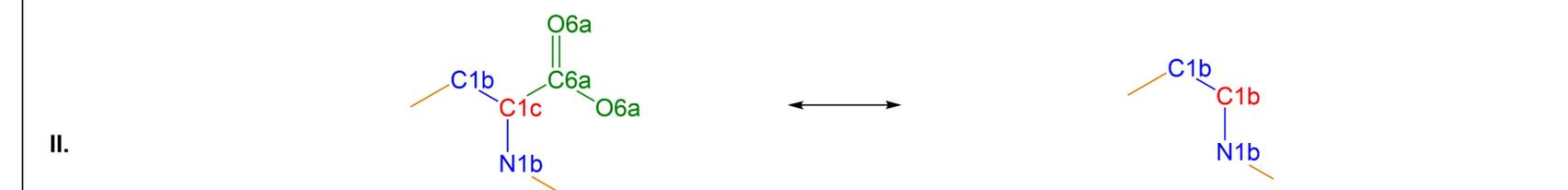
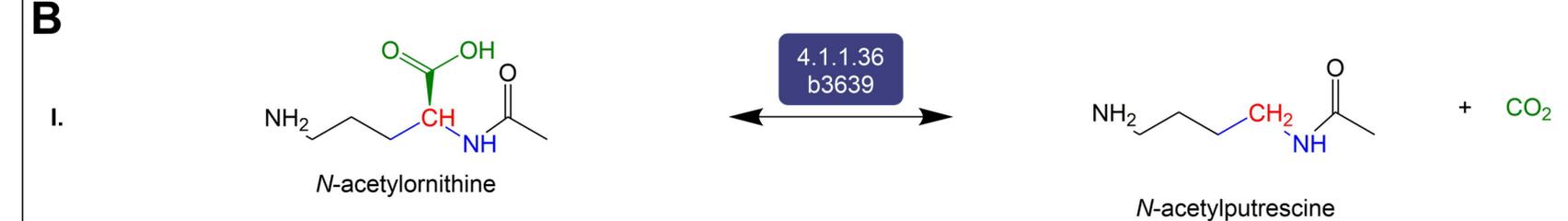
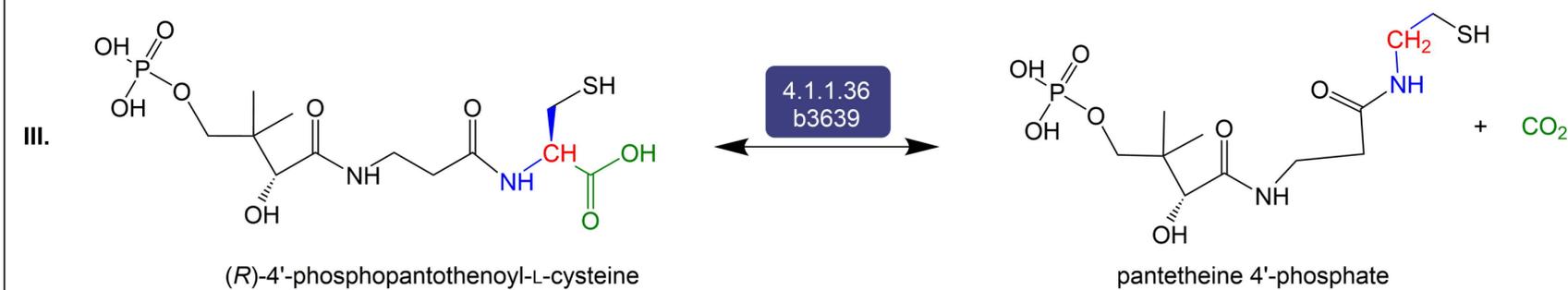
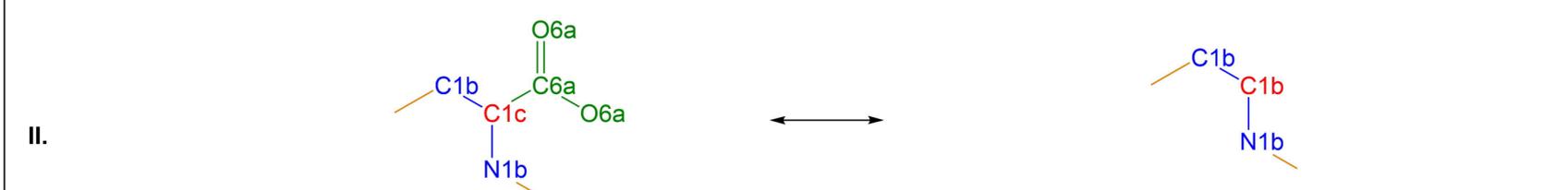
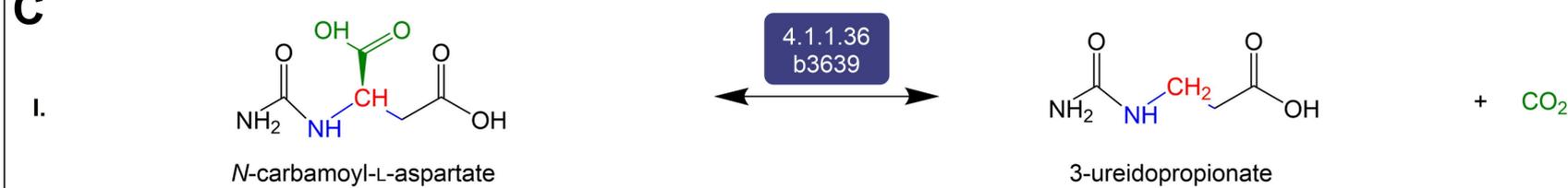
**A****B****C****D**



**A****B****C****D****E**

**A**

bioRxiv preprint doi: <https://doi.org/10.1101/536060>; this version posted May 24, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

**B****C****D**