
A compact model of *Escherichia coli* core and biosynthetic metabolism

Marco Corrao^{1,*}, Hai He², Wolfram Liebermeister³, Elad Noor⁴, Arren Bar-Even^{5,†}

¹ Department of Engineering Science, University of Oxford, Parks Road, OX1 3PJ Oxford, UK

² Max Planck Institute for Terrestrial Microbiology, Karl-von-Frisch-Str. 10, 35043 Marburg, Germany

³ Université Paris-Saclay, INRAE, MaIAGE, 78350 Jouy-en-Josas, France

⁴ Department of Plant and Environmental Sciences, Weizmann Institute of Science, 7610001 Rehovot, Israel

⁵ Max Planck Institute of Molecular Plant Physiology, Am Mühlenberg 1, 14476 Potsdam-Golm, Germany

* Corresponding author: Marco Corrao; email: marco.corrao@eng.ox.ac.uk

† Deceased in 2020

Abstract

Metabolic models condense biochemical knowledge about organisms in a structured and standardised way. As large-scale network reconstructions are readily available for many organisms, genome-scale models are being widely used among modellers and engineers. However, these large models can be difficult to analyse and visualise and occasionally generate predictions that are hard to interpret or even biologically unrealistic. Of the thousands of enzymatic reactions in a typical bacterial metabolism, only a few hundred form the metabolic pathways essential to produce energy carriers and biosynthetic precursors. These pathways carry relatively high flux, are central to maintaining and reproducing the cell, and provide precursors and energy to engineered metabolic pathways. Focusing on these central metabolic subsystems, we present *iCH360*, a manually curated medium-scale model of energy and biosynthesis metabolism for the well-studied bacterium *Escherichia coli* K-12 MG1655. The model is a sub-network of the most recent genome-scale reconstruction, *iML1515*, and comes with an updated layer of database annotations and with a range of metabolic maps for visualisation. We enriched the stoichiometric network with extensive biological information and quantitative data, enhancing the scope and applicability of the model. In addition, we assess the properties of this model in comparison to its genome-scale parent and demonstrate the use of the network and supporting data in various scenarios, including enzyme-constrained flux balance analysis, elementary flux mode analysis, and thermodynamic analysis. Overall, we believe this model holds the potential to become a reference medium-scale metabolic model for *E. coli*.

Keywords: Metabolic network, functional annotation, gene-reaction mapping, Elementary Flux Mode, biochemical thermodynamics

Abbreviations: EFM: Elementary Flux Mode; FBA: Flux Balance Analysis; GEM: Genome-scale models; MDF: Max-Min Driving Force; PMO: Probabilistic Metabolic Optimisation; RMSE: Root Mean Squared Error; 3GP: glycerol 3-phosphate; 3hmrsACP: (R)-3-Hydroxytetradecanoyl-ACP; 3MOB: 3-Methyl-2-oxobutanoate; Ala: L-alanine; Asp: L-aspartate; ATP: adenosine triphosphate; F6P: fructose 6-phosphate; G3P: glyceraldehyde 3-phosphate; HdeACP: cis-hexadec-9-enoyl-ACP; palmACP: palmitoyl-ACP; PEP: phosphoenolpyruvate; Ru5P: D-ribulose 5-phosphate.

1 Introduction

Metabolic models are a valuable tool for biologists and biotechnologists who want to elucidate and engineer cell metabolism [1–3]. In their simplest form, such models encode basic biochemical knowledge, such as network structure, reaction stoichiometries, or known reaction directionalities, in a structured and standardised format. However, the scope of these models is often wider, including information on catalysing enzymes and kinetic parameters, as well as annotations that link model elements to external databases. Rapid development of high-throughput experimental and computational pipelines has led to genome-scale metabolic network reconstructions now existing for a wide range of microorganisms [2, 4]. One of them, *Escherichia coli*, has been the most studied prokaryotic organism and, as such, its metabolism has been the subject of extensive modelling efforts spanning over three decades [5–8]. In particular, for the common laboratory strain *E. coli* K-12 MG1655, the most recent genome-scale reconstruction, *iML1515*, accounts for 2712 enzyme-catalysed reactions mapped in detail to 1515 genes [8].

Genome-scale metabolic network models (GEMs) provide a comprehensive picture of cell metabolism, and constraint-based modelling algorithms that use these models have shown remarkable predictive power, for example when predicting gene essentiality in bacteria [9]. However, working with such large models comes with some disadvantages. In the absence of sufficient constraining or parameterisation, simulations based on large networks can easily lead to biologically unrealistic solutions. For example, when designing and testing gene knockout strategies, genome-scale networks often wrongly predict unphysiological metabolic bypasses that must be manually inspected and filtered out [10, 11] (see Supplementary Table 1 for some examples). Another issue is that, owing to their size and complexity, analysis of genome-scale networks is often limited to relatively simple modelling frameworks, such as flux balance analysis (FBA), that can only answer a limited range of questions. More complex methods, including sampling of metabolic flux distributions, elementary flux mode (EFM) analysis [12], or kinetic modelling can be used to gain additional insight into the governing principles and constraints of microbial metabolism, but are difficult to apply to large models. Finally, genome-scale models are often hard to visualise comprehensively, which can make the interpretation of computed flux distributions cumbersome and unintuitive.

For all these reasons, small-scale models of *E. coli* metabolism are commonly used instead, both for strain design and for the development of novel modelling frameworks. Among these, the *E. coli core model* (ECC) developed by Orth et al. [13] has been widely used in the literature. Although popular as an educational and benchmark tool, ECC is limited in scope: it lacks, among others, most biosynthesis pathways, which would be relevant to many metabolic engineering applications. This limitation was addressed by Hädicke and Klamt [14], who constructed a medium-scale model, *E. coli core 2* (ECC2), as a subnetwork of *iJO1366*, the most up-to-date GEM available at the time [7]. ECC2 was obtained through an algorithmic reduction [15] that iteratively prunes reactions from a template model while retaining user-specified structural and phenotypic features, such as the ability to grow under a defined set of conditions. However, to enforce the desired phenotypes, the algorithm relied only on steady-state stoichiometric modelling and did not account for other important factors, such as thermodynamics, kinetics, or regulatory effects, which are relevant under physiological conditions. Therefore, while the resulting submodel satisfies the stoichiometric constraints imposed by construction, further manual curation is often needed depending on the application at hand.

Here, we introduce *iCH360* (named, according to convention, by the initials of the authors followed by the number of genes covered by the model), a manually curated “Goldilocks-sized” model of *E. coli* K-12 MG1655 energy and biosynthesis metabolism. The model was derived from the most recent genome-scale reconstruction (*iML1515* [8]) and includes all pathways required for energy production and for the biosynthesis of the main biomass building blocks, such as amino acids, nucleotides, and fatty acids, while the conversion of these precursors into more complex biomass components is described by an effective biomass-producing reaction. We extended the coverage of annotations that point to external databases from *iML1515*, and built custom metabolic maps to facilitate visualisation of the model and its subsystems [16]. We complemented the stoichiometric network structure with a curated layer of biological information on catalytic function, protein complex composition, and small molecule regulation. Finally, we enriched the model with useful quantitative data, including thermodynamic and kinetic constants. Thanks to all these extra layers of information, our model can support a wide range of modelling methods beyond simple stoichiometric ones.

In the following, we present the model and demonstrate several use cases across various modelling scenarios, including enzyme-allocation predictions, EFM analysis, and thermodynamic analysis. The model is freely available in the standard formats SBML and JSON and can be used directly with popular metabolic modelling tools such as the COBRA toolbox [17].

Metabolic model *iCH360*

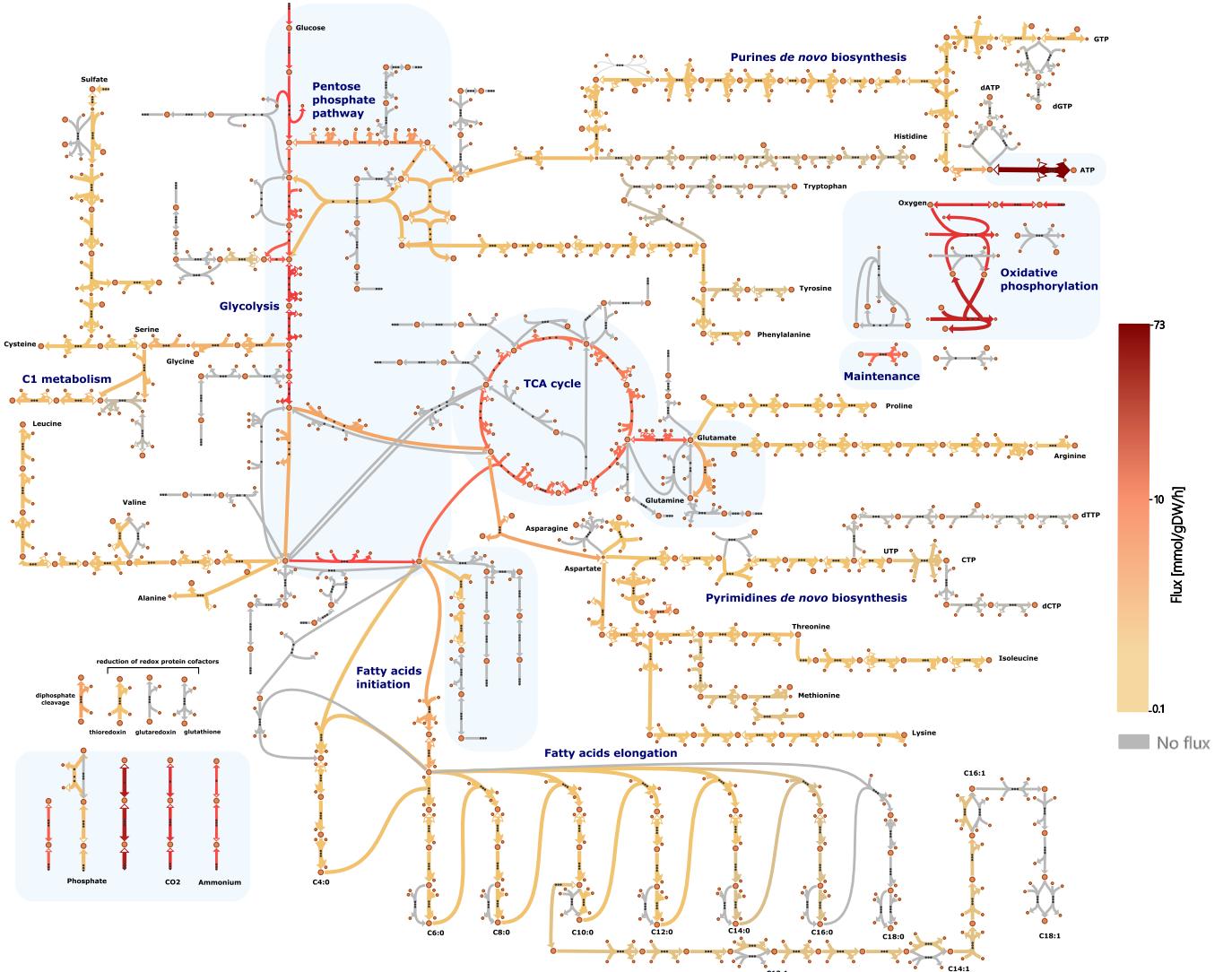


Figure 1: Metabolic map of the *iCH360* model, built with the metabolic visualisation tool Escher [18] and showing the metabolic subsystems included in the model. Shaded areas denote metabolic subsystems already present in the ECC model [13]. Reaction and metabolite names were omitted from the plot for clarity. Overlaid onto the map is a flux distribution for aerobic growth on glucose, computed via parsimonious FBA.

2 Results

2.1 A compact model of *Escherichia coli* energy and biosynthesis metabolism

To assemble the *iCH360* model, we started from the core metabolic reactions present in ECC and extended them with a curated set of pathways required for the biosynthesis of the main biomass building blocks, including the twenty amino acids, the four nucleotides, and both saturated and unsaturated fatty acids (Fig. 1, Table 1, and Supplementary Figures S1-S4). On the other hand, we specifically omitted from our model the pathways required for the biosynthesis of complex biomass components and polymers, most degradation pathways, the pathways involved in the *de novo* biosynthesis of cofactors, and the reactions involved in the uptake of metals and ions. In addition, while not performing a comprehensive review of *iML1515*, we applied a small number of corrections to the original reactions based on the knowledge from the literature (Supplementary Information A.1).

The final assembled model (Figure 1) comprises 304 compartment-specific metabolites (254 chemically unique compounds) and 323 metabolic reactions mapped to 360 genes, thus qualifying as a medium-scale model ranging in between ECC and *iML1515* (Supplementary Figure S5). Although similar in scale, our model and ECC2 present a fundamental structural difference. ECC2 was built by systematically removing reactions from its genome-scale parent (*iJO1366*) [14]. Thus, its metabolic space spans the production of all compounds consumed in the *iJO1336* biomass reaction. In contrast, the metabolic space of

Table 1: The main metabolic subsystems covered by *iCH360*.

| Subsystem | Description | Metabolic map |
|-----------------------------|--|-------------------------|
| Carbon uptake and transport | Reactions required for the uptake and assimilation of the following carbon sources: glucose, fructose, ribose, xylose, lactate, acetate, gluconate, pyruvate, glycerol, glyc erate, succinate, 2-ketoglutarate | Figure 1 |
| Central carbon metabolism | Glycolysis, pentose phosphate pathway, pyruvate fermentation, TCA cycle, oxidative phosphorylation | Figure 1 |
| Amino acids biosynthesis | Biosynthesis of all 20 amino acids from core metabolism precursors | Supplementary Figure S1 |
| Nucleotide biosynthesis | Biosynthesis of purine and pyrimidine nucleotides (and deoxyribonucleotides) from core and amino acid metabolism | Supplementary Figure S2 |
| Fatty-acids biosynthesis | Biosynthesis of saturated and unsaturated fatty acids present in <i>E. coli</i> from acetyl-CoA | Supplementary Figure S3 |
| C1 metabolism | One-carbon metabolism | Supplementary Figure S4 |

Table 2: Biosynthesis pathways outside of the *iCH360* model that were considered to construct a biomass reaction equivalent to the biomass reaction in the *iML1515* model. The right column shows the main precursors present in *iCH360*. Note that only the main precursors are shown here, but the equivalent biomass reaction computed also accounts for any net production or consumption of metabolites in the reduced model.

| Pathway | Precursor in <i>iCH360</i> |
|--|--------------------------------------|
| Biosynthesis of phosphatidylethanolamine (C16:0 and C16:1) | 3GP, PalmACP (C16:0), HdeACP (C16:1) |
| Biosynthesis of KDO2-Lipid-A | F6P, Ru5P, 3hmrsACP |
| Murein Biosynthesis | G3P, PEP, F6P, Ala |
| NAD/NADP de novo biosynthesis | Asp, DHAP |
| FAD de novo biosynthesis | GTP, Ru5P |
| CoA de novo biosynthesis | 3MOB, Asp |
| Active transport of ions | ATP |

iCH360 only reaches biomass building blocks, without explicitly considering peripheral pathways such as the assembly of cell-membrane components, *de novo* synthesis of cofactors, and active transport of ions in the cell. To make the model comparable to its parent model *iML1515*, we constructed an equivalent biomass reaction, in which all biomass requirements not present in our model are summarised by an equivalent metabolic cost in terms of precursors, based on manually curated bioproduction routes (Table 2, Supplementary Information A.2, and Supplementary File S1).

2.2 Range of metabolic conversions described by production envelopes

To check how its significantly smaller size and complexity affect the solution space of our model, we first looked at the maximum achievable biomass production flux under a range of growth conditions (Supplementary Figure S6). Across most conditions considered, the model achieves biomass fluxes comparable to those of its genome-scale parent. The main differences exist in anaerobic growth on fumarate, alpha-ketoglutarate (AKG), and glycerol, where our model predicts zero growth, while the GEM achieves some (albeit small) biomass production rate. In practice, fermentation in these scenarios is biologically unrealistic (e.g. [19]). The reduced model is thus a good basis for metabolic simulation frameworks with few constraints such as FBA, since the reduction procedure, at least in this test, limits the original solution space of the GEM to more physiologically relevant regions.

To further investigate the metabolic capabilities of our model, we looked at production envelopes (projections of the model solution space onto a smaller set of dimensions, also known as phenotypic phase planes) describing production rates for biomass and a range of metabolites (Methods). Figure 2 shows the resulting envelopes for aerobic growth on glucose. The reduced model has production envelopes comparable to *iML1515*, except for the production of acetate, where the genome-scale model can

achieve considerably higher yields, both aerobically and anaerobically. In order to understand the cause of these differences, we investigated optimal acetate production routes in both models using Flux Balance Analysis. Under aerobic conditions, we found that the differences can be traced to different production abilities for acetyl-CoA, a precursor of acetate. In *iCH360*, acetyl-CoA is produced entirely from pyruvate oxidation via either the Pyruvate Dehydrogenase (PDH) or the Pyruvate-Formate-Lyase (PFL) reactions (Supplementary Figure S7A, left). On the other hand, *iML1515* can produce acetyl-CoA via a number of additional pathways not included in our model (Supplementary Figure S7A, right). These include: the degradation of threonine either directly to acetyl-CoA (THR_D, GLYAT), or indirectly via acetaldehyde (THRA, THR_D, ATHRDH_r, THRA2), which is then converted to acetyl-CoA; the degradation of 2-deoxy-D-ribose 5-phosphate into acetaldehyde (DRPA), followed by conversion into acetyl-CoA; the degradation of autoinducer 2 into acetyl-CoA (AI2K, PAI2I, PAI2T). Indeed, we found that, in the region of the production envelope where the two models diverge significantly, up to 50% of acetyl-CoA production in the genome-scale model is accounted for by these degradation pathways. This is in disagreement with the understanding of the existing literature that under aerobic growth on glucose, most of the acetyl-CoA is produced from oxidation of pyruvate [20, 21]. Confirming these findings, the simultaneous deletion of four reactions blocking these degradation pathways (DRPA, PAI2T, THRA, THR_D) in *iML1515* brings aerobic acetate production to levels comparable to *iCH360* (Supplementary Figure S7C).

In the anaerobic scenario, we found that the differences between the two models are exacerbated by the ability of the genome-scale model to produce more pyruvate (which in turn results in higher acetyl-CoA production) than *iCH360*. Investigating the cause of this difference, we found that, under anaerobic conditions, the *iML1515* solution involves the uptake of external CO₂ and its use as a sink for electrons produced in glycolysis (Supplementary Figure S7C), which is thermodynamically unrealistic under ambient CO₂ levels. Blocking CO₂ uptake reduces the maximal anaerobic pyruvate yield in the genome-scale model, but does not fully close the gap with the production capabilities of our model (Supplementary Figure S7D), implying the existence of additional routes for anaerobic production of pyruvate that are not included in *iCH360*. A similar pattern was found when comparing production envelopes between the two models under different growth conditions (Supplementary Figure S8).

Overall, the higher maximal yields achievable by *iML1515* for acetyl-CoA (aerobically) and pyruvate (anaerobically) appear unrealistic. Hence, in agreement with what was previously done with ECC2 [14], we chose not to include in *iCH360* any additional reactions that would allow for higher acetate yields.

Production envelopes generated with *iCH360* were also comparable, albeit not identical, with those computed on other existing medium-scale models, namely ECC and ECC2 (Supplementary Figures S9). Particularly, the maximal yields achievable for each product were nearly identical for the products considered, both aerobically and anaerobically. Some differences between the envelopes can be observed for solutions with higher predicted growth rates. However, since the three models are equipped with different biomass reactions sourced from different parental models (and hence predict different biomass yields), these differences are expected.

2.3 Connecting reactions to their catalysing enzymes and the enzymes' protein components

Metabolic models often contain annotations that connect model elements to entries in biological databases, such as BioCyc [22], KEGG [23] and MetaNetX [24]. However, even for the subset of reactions included in *iCH360*, the annotations present in *iML1515* were incomplete and, in part, outdated. To fill these gaps, we extended and corrected the original annotations through a mixture of automated querying and manual curation (Figure 3A). Notably, the annotations pointing to the BioCyc knowledgebase [22], are nearly complete: Out of 321 enzymatic reactions in the model, 317 are mapped to BioCyc with a single ID (for the remaining four unannotated reactions, involved in the biosynthesis of unsaturated fatty acids, a match in the database could not be found for the specific use of NADPH as a redox cofactor). Further, nearly all of these BioCyc annotations (315 / 317) are in the ECOLI namespace and therefore point to the organism-specific EcoCyc database, a widely used and extensive reference for *E. coli* molecular biology [25, 26]. The remaining two reactions map instead to the broader MetaCyc database, also part of the the BioCyc ecosystem, via the META namespace. Additionally, we found 134 deprecated annotations pointing to the MetaNetX database, which were consequently updated with the most up-to-date IDs.

Using the extensive mapping to the EcoCyc database, we parsed, assembled, and manually curated a knowledge graph that enhances the stoichiometric model with a detailed layer of information about enzymes, polypeptides, and genes related to the network (Methods). This data structure takes the form of a weighted graph, where nodes represent biological entities (reactions, proteins, genes, or compounds) and edges represent (potentially quantitative) functional relationships between them (Methods, Figure 3B, Supplementary Tables 2 and 3), such as catalysis, regulation, protein modification, and protein subunit composition. This graph contains the information collected in a unified form, allowing users to perform a number of tasks that occur in metabolic modelling applications. For example, by explicitly mapping reactions to their catalysing enzymes rather than to the

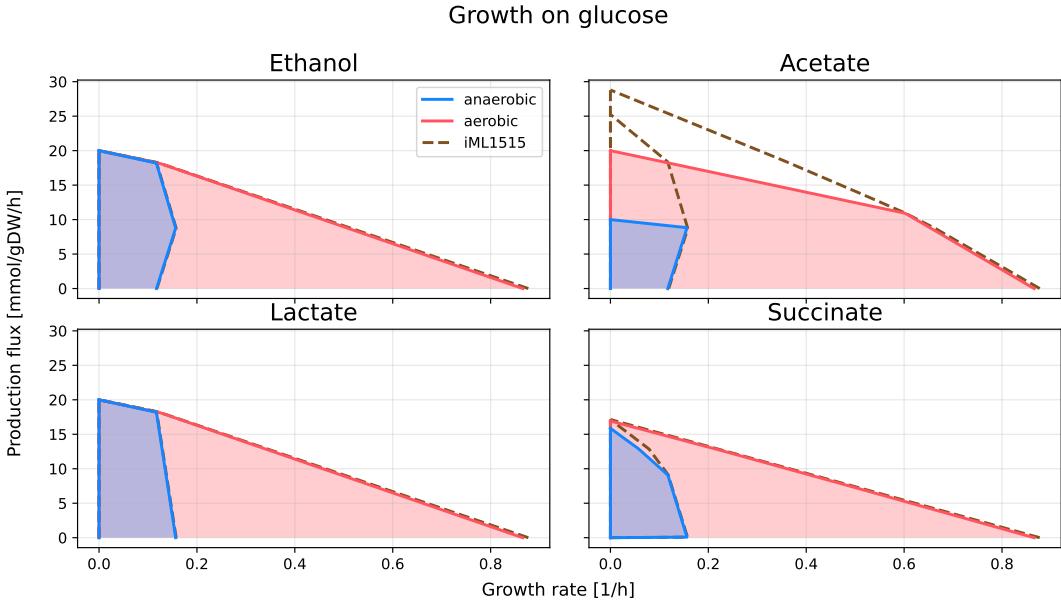


Figure 2: The *iCH360* model shows similar, but more precise metabolic capabilities than *iML1515*. Considering glucose as a feedstock and studying ethanol, lactate, and succinate production, a production envelope analysis yields similar results in the two models (note that the dashed line representing the production envelope of *iML1515* is sometimes hidden behind the coloured lines). Growth rate and production fluxes were computed by limiting the glucose uptake rate to a maximum of 10 mmol/gDW/h. In the scenario of acetate production (top right panel) *iCH360* avoids an unrealistically high production flux [14] as predicted by *iML1515*.

associated genes (which may additionally include protein activators or cofactors), it simplifies the definition of meaningful enzyme capacity constraints in the model. Similarly, since the graph topology implicitly defines associations between reactions and genes, boolean gene-protein-reaction (GPR) rules needed for *in silico* knockout studies can be generated that account for catalytic and noncatalytic requirements for each reaction (Supplementary Information A.3). Crucially, these GPR rules can be regenerated as needed whenever the graph is updated with new nodes or edges. Finally, we used it to estimate the abundance of protein complexes included in the model from the measured polypeptide abundances using a simple automated procedure (Supplementary Information A.4, Results section 2.4).

Through this annotation graph, 318 metabolic reactions in the model are linked to 289 catalysing enzymes, with more than 25% of the reactions being catalysed by multiple enzymes (isozymes). Such enzymatic redundancy plays an important role, for example, when designing metabolic engineering strategies to prevent flux through a pathway. However, the different isoenzymes of a reaction need not all be equally important, and treating them as completely equivalent may generate inaccuracies in some phenotypic predictions. For example, while phosphofructokinase activity in *E. coli* (reaction PFK) is known to be carried by two isozymes, *pfkA* and *pfkB*, the latter is known to be responsible only for minor activity under normal physiological conditions [28]. If these differences are not taken into account, metabolic models can overpredict redundancy in the network, for example, when predicting phenotypes after gene knockouts [8].

To address this issue and make the model usable in a wider range of modelling scenarios, we classified the catalytic edges in the graph as either *primary* or *secondary catalysis* (Methods, Figure 3C). The catalytic relationship between a reaction and an enzyme was annotated as secondary whenever the enzyme, according to experiments, accounts only for negligible activity for the reaction in the wild-type strain. As experimental evidence we considered *in vitro* and *in vivo* complementation studies. Through this curation process, a total of 72 catalytic edges were functionally annotated as secondary.

To test how this functional annotation can support phenotypic predictions, we investigated the disruption of each essential reaction in the model, resulting from the knockout of each of its associated genes (Methods). We classified the outcome of each possible knockout as: (i) complete disruption, if a reaction loses all of its catalytic edges; (ii) full primary disruption, if a reaction loses all of its primary catalysis edges, but secondary ones are left; (iii) partial primary disruption, if the reaction loses some, but not all, of its primary catalysis edges, or (iv) secondary disruption, if some or all secondary edges are disrupted, but none of the primary ones. We classified the effect of *in silico* knockouts for aerobic growth on 9 different carbon sources and compared the results against a large data set of mutant fitness data, obtained via competitive fitness assays, from Price et al. (2018) [27].

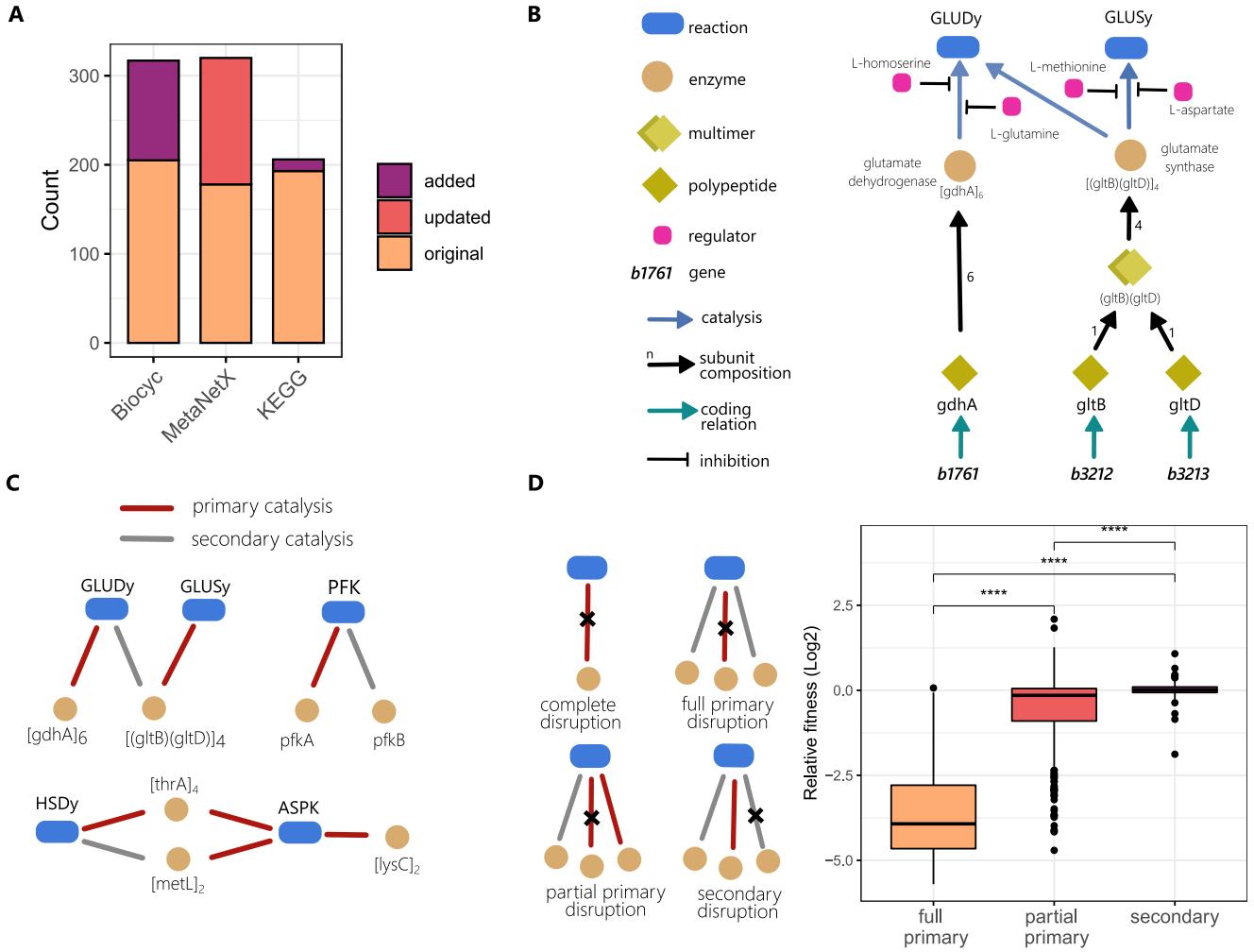


Figure 3: Layers of annotation and biological knowledge supporting the stoichiometric model in *iCH360*. **A:** Annotations for the model reactions point to the BioCyc, metaNetX, and KEGG databases. Bars show the numbers of annotations, highlighting the share of annotations that were added to or updated from the parent model *iML1515*. **B:** Some of the biological knowledge parsed from EcoCyc (and manually curated) included in the model-supporting functional annotation graph. The graph captures catalytic relationships between reactions and enzymes, protein subunit compositions, protein-gene mappings, and small-molecule regulation interactions, among others. Shown here as an example are the branches of the graph corresponding to the Glutamate Dehydrogenase (GLUDy) and Glutamate Synthase (GLUSy) reactions. **C:** Examples of catalytic relationships functionally annotated as either primary or secondary in the graph. Note that all catalytic relationships were classified as primary by default, unless sufficient evidence was found to annotate them as secondary. **D:** Functional annotation of catalytic edges as primary or secondary supports can be used to improve phenotypic predictions. Left: Classification of catalytic edge disruptions in the network resulting from simulated knockout of genes associated with essential reactions in the model across 9 growth conditions (see text for a description of each disruption class). Right: comparison of predicted disruption outcomes against a large data set of mutant fitness data [27] shows that the different types of disruption tend to lead to significantly different fitness changes.

Based on this analysis, we found that disruptions of the primary edges as a result of a knockout were significantly associated with greater fitness losses than disruptions of secondary edges (Wilcoxon rank-sum test, $p < 10^{-6}$, Figure 3D). In addition, primary disruptions that were only partial were associated with more contained fitness losses. Finally, minor fitness gains could be identified when mutants were associated with secondary and partial primary disruptions, but not when complete primary disruptions occurred. We did not find any significant differences between complete disruptions and full primary disruptions (Supplementary Figure S10), supporting the idea that catalytic relationships annotated as secondary are unlikely to be strong enough to stand-in for disrupted primary ones under normal physiological conditions.

2.4 Enzyme-constrained Flux Balance Analysis with EC-*iCH360*

To make *iCH360* applicable for enzyme-constrained flux simulations, we constructed a version of the model containing all necessary extra information, which we denote as EC-*iCH360*. We constructed EC-*iCH360* in the sMOMENT [29] format (Methods). The sMOMENT framework is inherently simple and generates the same solution space as more complex model formats, such as GECKO, unless reaction-specific capacity constraints are specified in the latter [29, 30]. Since it requires

unique reaction-enzyme mappings, we used the knowledge graph to remove all secondary catalytic relationships in the model (Section 2.3).

We first parametrised the model by defining a flux-specific enzyme cost for each reaction (in units of grams of enzyme per unit flux), using as values the estimated *in vivo* turnover numbers from Heckmann et al. (2020) [31]. Using this enzyme-constrained model, we then predicted enzyme abundances and compared them to experimental enzyme abundances, estimated from a data set of proteomic measurements for aerobic growth on eight different carbon sources [32] (Methods). This analysis led to predictions with root mean squared error (RMSE, computed for \log_{10} -transformed enzyme abundances) ranging from 0.53 to 0.62 (Supplementary Figure S11). To assess the nature of residuals between measurements and predictions, we investigated the geometric mean of enzyme abundances across all eight conditions (Supplementary Figure S11, bottom right panel). If the mismatch between measurements and predictions were due to each enzyme operating at different saturation levels in each condition, one would expect that averaging would reduce these differences. However, averaging the predicted and measured enzyme abundances, respectively, across conditions did not significantly improve the RMSE, indicating that the abundances of individual enzymes were systematically over- or under-predicted between conditions.

To increase the predictive capacity of the model, we adjusted the turnover numbers by fitting them to experimental measurements through a custom heuristic (Methods, Supplementary information A.5). By *simultaneously* fitting all available conditions with a single set of parameters, we ensured that our adjustment procedure is robust to condition-specific biases. Furthermore, by introducing regularisation within our adjustment scheme, we penalised large deviations of parameters from the original data set, increasing the robustness of the procedure to overfitting. Our adjusted parameter set shows a mean absolute deviation (computed for \log_{10} -transformed turnover values) from the original parameter set of ≈ 0.22 (Supplementary Figure S12) and results in significantly better agreement with experimental measurements across conditions (Figure 4A). Further, mean enzyme abundances across conditions are very well predicted with the adjusted parameter set (Figure 4B), implying that residuals between measurements and predictions are now to be attributed to variability in enzyme saturation across conditions (which simple frameworks such as enzyme-constrained FBA, which use a constant enzyme cost per unit flux, cannot account for), rather than systematic over- or under-predictions. Thus, each adjusted turnover parameter can be thought of as a “typical” apparent k_{cat} value, incorporating average saturation trends for an enzyme across growth conditions. To further confirm this aspect, we performed a leave-one-out cross-validation analysis (Methods), where each condition was in turn excluded from the model fitting dataset, but used to evaluate predictions. Results (Supplementary Figure S14) show that the predictions of enzyme abundances in a given condition are considerably improved even when data from that condition is not used for fitting, confirming that our parameter fitting heuristics is capturing global trends and not condition-specific effects. Training the model on the full dataset, as we did to compute the final parameter set, further improves predictions, even if to a lesser extent. This can be explained by noting that, by design, our procedure cannot adjust the turnover parameters for enzymes associated with zero flux in the reference flux distribution used for fitting (Supplementary information A.5). Hence, including all conditions in the training set allows for the turnover number of highly condition-specific enzymes (those associated to nonzero flux in only one of the reference flux distributions used by the procedure) to also be adjusted, further improving the overall prediction metrics.

2.5 Elementary Flux Modes in the reduced model variant $i\text{CH360}_{\text{red}}$

Despite the small size of $i\text{CH360}$, we found the explicit enumeration of its elementary flux modes (EFMs) to be intractable. This is not necessarily surprising, since EFM count is crucially dependent on the topology of a metabolic network rather than its sheer size. Metabolic networks can possess different types of redundancy, such as the presence of alternative pathways for the production of the same metabolite, the use of alternative cofactors for the same catalytic step in a pathway, or the presence of alternative transporters for the uptake/excretion of a compound. Although knowledge about these redundancies is often valuable, including them in the model can increase the number of EFMs exponentially, hampering or even preventing EFM-based analyses.

To address this issue, we identified and removed a small set of alternative metabolic routes in $i\text{CH360}$, using available information from the literature, whenever possible, to ensure that the most physiologically relevant alternative was maintained (Supplementary Table 4). This resulted in a metabolic submodel of $i\text{CH360}$, a model variant that we denote by $i\text{CH360}_{\text{red}}$. $i\text{CH360}_{\text{red}}$ contains 305 metabolic reactions (18 less than $i\text{CH360}$) and shows the same production envelopes as its parent model for a number of metabolites of interest (Supplementary Figure S15). While the number of EFMs in $i\text{CH360}_{\text{red}}$ is still relatively large (≈ 13.5 millions for aerobic growth on glucose, see Supplementary Table 5), it is not prohibitive for most types of EFM-based analysis, and their explicit enumeration does not require high-performance computing (Methods).

We used the EFMs of $i\text{CH360}_{\text{red}}$ to study the possible combinations of achievable growth rates and yields in the network [12]. To this end, we considered growth on glucose as a scenario and computed, for each EFM from $i\text{CH360}_{\text{red}}$, its yield,

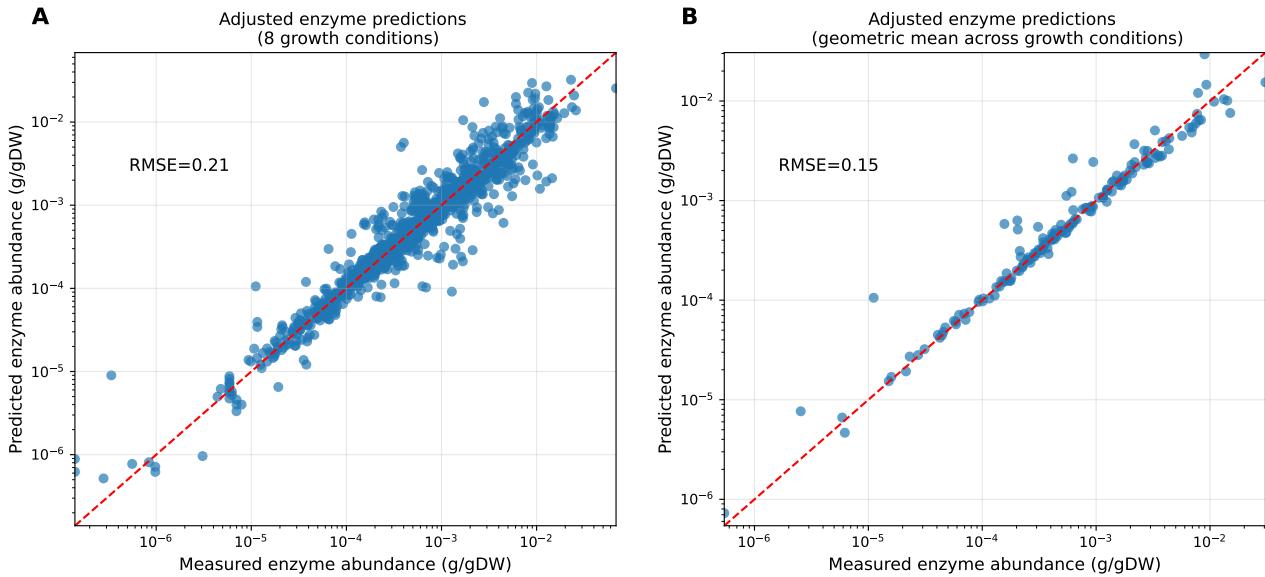


Figure 4: Enzyme allocation predictions obtained with the model variant EC-*i*CH360 after adjusting the turnover parameters. **A:** Predicted vs measured enzyme abundances for aerobic growth on eight different carbon sources. Each data point represents an enzyme-condition pair. A total of 325 data points corresponding to zero predictions (enzymes associated with zero-flux in the enzyme-constrained FBA solution for a given condition) were omitted from the plot. **B:** Geometric mean across conditions of predicted vs measured enzyme abundances. For each enzyme, the geometric mean was computed across the conditions with non-zero predicted abundance. A total of 27 data points, corresponding to enzymes with zero predictions across all conditions, were omitted from the plot.

computed as the ratio of biomass flux and glucose uptake, and its achievable cell growth rate, which we estimated based on the enzyme costs defined for the enzyme-constrained model (Methods). Based on this analysis, we identified a front of Pareto-optimal EFM_s, along which any increase in the growth rate will necessarily lead to a reduction in yield (Figure 5). Along the Pareto front, we observe a transition from a purely respiratory mode at maximum yield (Supplementary Figure S16) to a mixed respiratory-fermentative mode at maximum growth (Supplementary Figure S17). Quantitatively, the extent of this trade-off was rather modest: the EFM with maximal yield reaches almost the maximal growth rate, so only minor gains in growth rate can be achieved by using other, fermentative modes along the Pareto front. However, it is worth noting that this analysis was performed using a simple capacity-based enzyme cost function, which ignores metabolite concentrations by assuming a constant enzyme cost per unit flux for each reaction. Repeating the analysis with a more complete enzyme-cost function, such as one that accounts for variable thermodynamic driving force and enzyme saturation, could help elucidate the nature of this trade-off [33]. To demonstrate that the shape and size of the Pareto front depend strongly on growth conditions, we also simulated an environment with very low oxygen levels. We implemented this by increasing the flux-specific enzyme cost of the oxygen-dependent reaction in the respiratory chain, equivalent to assuming a lower enzyme efficiency due to a lower oxygen level (Methods). Results (Supplementary Figure S18) show a much broader front of Pareto optimal EFM_s, indicating that the nature of the observed trade-off is indeed condition-dependent. Notably, a similar dependence of the Pareto front on extracellular oxygen availability has been previously observed in a small-scale model of *E. coli* core metabolism [12].

2.6 Saturation FBA and modelling of overflow metabolism

In order to study the effect of external conditions on optimal metabolic strategies in more detail, we used another framework which does not require an enumeration of EFM_s and allows for additional flux bounds, for example, to impose a minimum ATP consumption rate for cell maintenance. The saturation FBA (satFBA) framework [34] is a variant of enzyme-constrained FBA, wherein a fixed enzyme cost per flux is assumed for all metabolic reactions in a model, except for the substrate transporter, for which a complete kinetic law is used (Methods). Since the external substrate concentration is a simple parameter, screening this concentration is equivalent to screening the values of the transporter efficiency. Here we used satFBA to simulate how the growth-maximising solution of the network varies in response to changing extracellular glucose concentration. By solving the satFBA problem for a range of glucose concentrations, we predicted the dependence of the cell's growth rate on substrate concentration, resulting in the typically observed Monod curve (Figure 6A). If the problem does not contain any further flux bounds (so that the magnitude of fluxes at the optimum is solely limited by the maximum enzyme availability) the solution of satFBA problems will be an elementary flux mode [34, 35]. Hence, in this case we can use satFBA to explore how a cell should switch across elementary modes as a function of the growth environment.

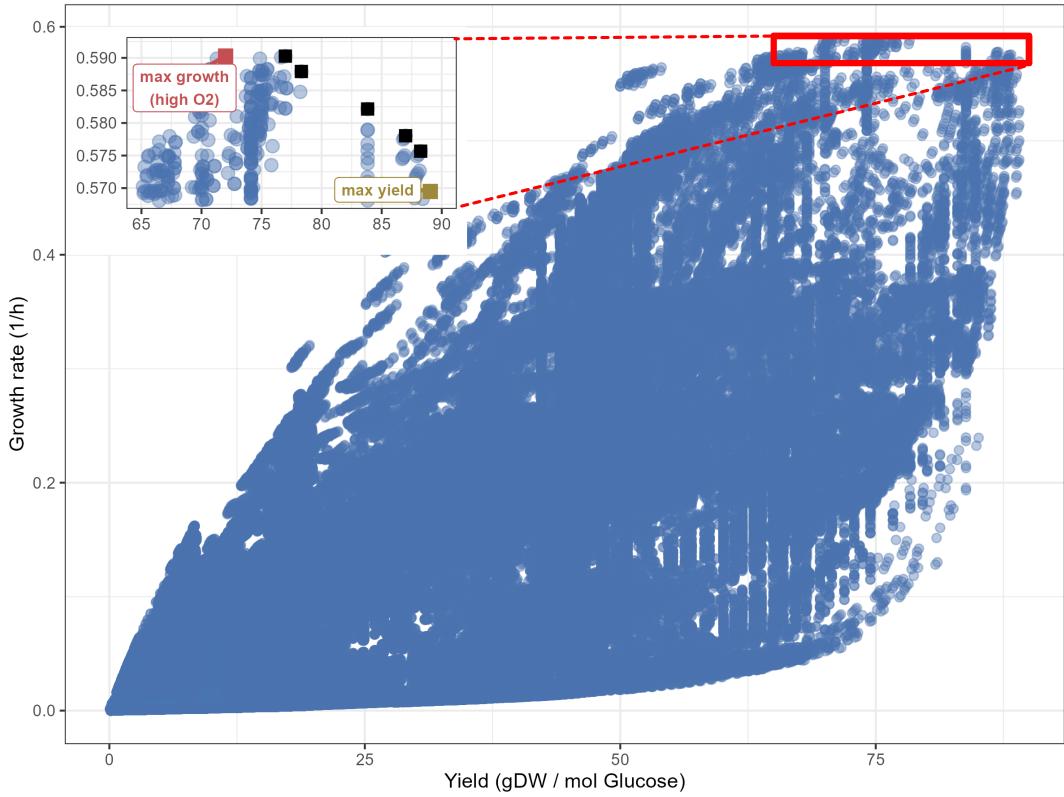


Figure 5: Growth rates and biomass yields achieved by different elementary flux modes of *iCH360_{red}* for growth on glucose. The inset on the top left (corresponding to the area of the plot enclosed by the red rectangle) highlights the front of Pareto-optimal EFM_s (squares), with the maximum-growth and maximum yield modes laying at the extremes of the front. The growth rate of each mode was estimated by assuming that the metabolic enzymes in the model occupy, by mass, a constant fraction of the cell's dry weight (see Methods).

At low glucose concentrations, the glucose transporter operates at low saturation and glucose uptake is enzymatically expensive, leading to a high-yield, purely respiratory metabolic mode at the optimum (Figure 6B and C). As substrate availability is increased, the cost of substrate uptake decreases and higher growth rates are achieved by switching to lower yield, acetate-secreting modes [34]. Since the yield of a flux distribution is, by definition, constant along an elementary flux mode, the yield varies in step-like manner as a function of external glucose concentration (Figure 6C, inset), where each jump represents a change of optimal mode. The satFBA formalism can also be used with additional flux bounds. For example, if a positive lower bound on ATP hydrolysis is added as a maintenance requirement, optimal solutions to the satFBA problem will no longer be elementary modes, and the yield of the optimal solution no longer follows a piecewise constant profile (Supplementary Figure S19C).

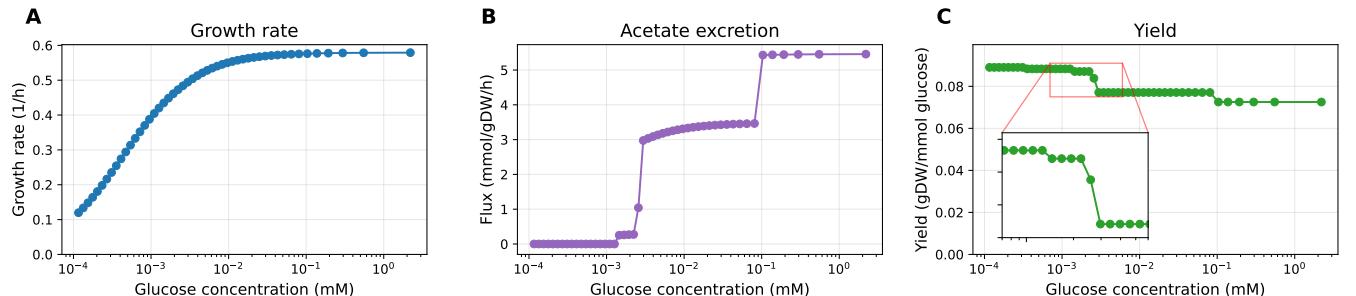


Figure 6: Saturation FBA enables the exploration of the optimal switching across elementary modes as a function of the growth environment. **A:** satFBA predictions for the growth rate as a function of external glucose concentration, showing a typical Monod curve. Note that satFBA computes the cell's growth rate by assuming a fixed total enzyme mass budget while varying the saturation of the substrate transporter as a function of external substrate concentration. Importantly, although the curve is continuous and smooth, it comprises many smaller sections, each dominated by a different elementary mode. **B:** satFBA predictions for the acetate excretion flux, showing progressively higher use of fermentative metabolism in the optimal solution as external glucose availability increases. **C:** The yield of the optimal satFBA solution progressively decreases in a step-like manner as external glucose availability increases. Each jump represents a switch in the optimal elementary flux mode.

2.7 Equilibrium constants, thermodynamic forces, and thermodynamically feasible states

Living systems operate outside of thermodynamic equilibrium, and thermodynamics places strong constraints on the operation of metabolic systems. In any metabolic state, the flux directions must follow the signs of thermodynamic forces, which depend on metabolite concentrations and equilibrium constants. To provide these constants as parts of our model, we used the component contribution framework [36] to estimate the standard Gibbs free energy ($\Delta_r G^\circ$) of each reaction (Methods). These estimates account for compartment-specific chemical environment parameters, such as pH, pMg, and ionic strength, and were corrected to account for protons and charge translocation in multi-compartment reactions.

The resulting parameter set covers the vast majority of the model (over 97% of metabolic reactions, with the remaining being not covered by the component contribution database used here) and accounts for the uncertainty in the estimates through a multivariate covariance matrix. Accounting for the correlations between different $\Delta_r G^\circ$ values becomes important when imposing thermodynamic constraints on the model. For example, the fatty acid biosynthesis subsystem in the model consists of repeated elongation cycles, where a short sequence of chemical transformations is repeatedly performed on a growing carbon chain. As a result, even if the $\Delta_r G^\circ$ for each reaction in the pathway is known with some uncertainty, this uncertainty is tightly correlated across the reactions, which constrains the set of achievable thermodynamic states in the network.

Using this set of thermodynamic constants, we first tested whether some typical flux distributions obtained from the model are thermodynamically feasible. To this end, we considered flux distributions generated by parsimonious Flux Balance Analysis (pFBA) across 12 growth conditions and computed their max-min driving force (MDF) [37], accounting for uncertainty in the estimates (Methods). We found a positive MDF for each of the flux distributions, indicating that all pFBA solutions tested are thermodynamically feasible. Notably, we found that the computed MDF values cluster very clearly in three groups, corresponding to aerobic growth on glycolytic substrates (high MDF), aerobic growth on gluconeogenic substrates (medium MDF), and anaerobic conditions (low MDF, Figure 7A).

Having confirmed that our reference FBA-derived flux distributions are thermodynamically feasible, we then employed an alternative flux prediction method that ensures thermodynamic feasibility by construction. The probabilistic metabolic optimisation (PMO) framework [38] uses a mixed-integer quadratic programming approach to compute a set of fluxes, metabolite concentrations, and reaction driving forces which is probabilistically most in agreement with experimentally measured metabolite concentrations (Methods). Using this framework, we computed a maximum-likelihood thermodynamic state for the model, simulating aerobic growth on glucose. In the thermodynamic state computed by PMO, we found that all metabolite concentrations lie within physiologically reasonable ranges (1 μM – 1 mM). In addition, all “anomalous concentrations” identified by the framework (metabolite concentrations lying more than one standard deviation away from the mean experimental value) had been identified and explained previously [38], and can potentially be addressed by lumping together those reactions for which substrate channelling is known to happen, as reported in literature.

We used the PMO-derived thermodynamic state to identify candidate bottlenecks, in terms of enzyme demand, across the network. To this end, we first computed the flux-force efficacy of each reaction in the thermodynamic state [37] (Methods). The flux-force efficacy is a unitless quantity, ranging from 0 to 1, denoting the ratio between net flux (forward minus backward flux) and total flux (forward plus backward flux) of a reaction. Reactions operating at low flux-force efficacy have a lower net flux due to two reasons: (1) the forward flux is lower in absolute terms, and (2) the higher backward flux is counter-productive and subtracts from the forward flux. Therefore, to achieve a given required net flux, the cell has to invest more resources in maintaining a higher enzyme level. To identify potential thermodynamic bottlenecks that lead to high enzyme costs in specific reactions, we therefore screened all reactions for their predicted flux-force efficacy and flux (Figure 7B) and identified those predicted to operate at low efficacy, while carrying significant flux (Figure 7B, labelled points).

To assess to which extent the thermodynamic operating states predicted by PMO are predictive of enzyme investment, we split reactions into two groups based on whether their predicted flux-force efficacy sits above or below 50%, and compared the distributions of measured enzyme abundances between the two groups (Methods, Supplementary Figure S20). While there are many determinants of enzyme abundance beyond thermodynamics, including enzyme turnover, affinity, or regulation, we observed a significant difference between the two distributions ($p < 0.001$, two-sided Wilcoxon rank-sum test), with the low efficacy having approximately a 3-fold higher median enzyme abundance than the high efficacy group. Establishing a causal link can be very difficult, but we can speculate that this negative correlation might be explained by metabolism evolving to compensate for the low efficacy of some reactions with a higher expression level of their catalysing enzyme [37, 39].

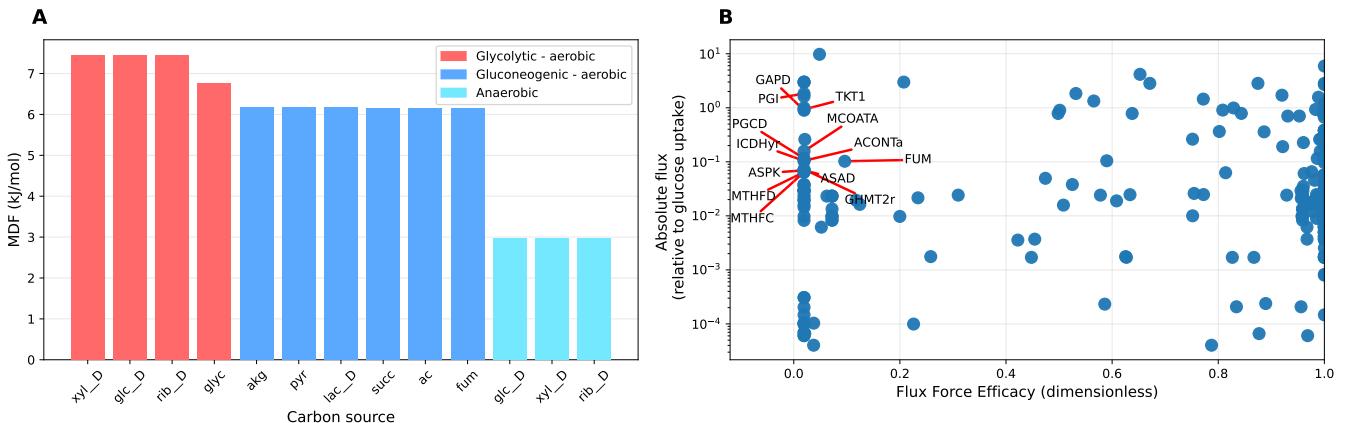


Figure 7: Thermodynamic analysis of the model via the curated thermodynamic parameter set. **A** Probabilistic max-min driving force (MDF) analysis of flux distributions obtained by parsimonious flux balance analysis for a total of 12 growth conditions. All flux distributions tested have a positive MDF, implying they are thermodynamically feasible under reasonable physiological metabolite concentration ranges. The computed MDF values cluster in three groups, corresponding to glycolytic aerobic, gluconeogenic aerobic, and anaerobic growth conditions. **B**: Fluxes relative to the glucose uptake flux (EX_glc_D.e) and flux-force efficacies computed by probabilistic metabolic optimisation (PMO). The labelled data points represent examples of reactions (excluding transport and spontaneous reactions) with low predicted flux-force efficacy (here, below 20%), but carrying high relative flux in the optimal solution (here, more than 5% of the glucose uptake flux). xyl_D: D-xylose; glc_D: D-glucose; rib_D: D-ribose; glyc: glycerol; akg: alpha-keto-glutarate; ac: acetate; pyr: pyruvate; lac_D: D-lactate; succ: succinate; fum: fumarate.

3 Discussion

Here we presented *iCH360*, a medium-scale metabolic model of *E. coli* covering central and biosynthesis metabolism, together with the associated data and metabolic maps and results from several analysed use cases. Similarly to previously constructed *core* models [13, 14], this model trades metabolic coverage for usability, interpretability, and ease of visualisation. It is well suited whenever a relatively small, highly curated network is desired, when computationally demanding analyses are to be performed, or as an educational tool in the field of metabolic modelling. When comparing some key properties of this model with those of its parent genome-scale model, we observed only small differences in the achievable biomass and product yields across a range of growth conditions, validating that, despite its contained size, the model captures the most salient metabolic features of the genome-scale network.

Further, we showed that the use of a well-curated, smaller-scale model can, in some instances, even correct unrealistic phenotypes predicted by its genome-scale parent. These unrealistic predictions from *iML1515* are not the result of “errors” in the metabolic model. Rather, they are the result of applying simple stoichiometric methods, such as FBA, to a large network with many degrees of freedom. It is possible that the inclusion of additional constraints, such as thermodynamic feasibility under physiological conditions or proteome allocation bounds, would automatically render such solutions infeasible. However, these constraints, and the parameters required to implement them, are not always readily available. By assembling a smaller model and curating it with expert knowledge, we filter out many of these behaviours *by construction*, providing users with a versatile and interpretable tool to investigate central and biosynthetic metabolism in *E. coli*. Clearly, this comes at the cost of limited applicability in other scenarios, such as those where the metabolic subsystems not included in *iCH360* (e.g. degradation pathways) are crucial to explaining or modelling a given phenotype. In these cases, the use of a genome-scale model (or an *ad hoc* reduction thereof) would still be an invaluable tool.

To make this model easily usable in a variety of applications, we enriched the stoichiometric network structure with a curated layer of biological knowledge in the form of a knowledge graph. This graph encodes information about biological entities in the network in a structured, ready-to-use format, including catalytic relationships between reactions and enzymes, the stoichiometric composition of protein complexes, and small-molecule regulation interactions. In addition, we mapped to the model a range of quantitative parameters, including *in vivo* turnover number estimates and thermodynamic constants, extending the use of the model beyond a simple stoichiometric analysis. A summary of the biological knowledge captured by *iCH360* is shown in Table 3.

Due to its medium size and the high level of curation, *iCH360* lends itself to a wide range of modelling methods. Here, we demonstrated some representative examples. These include the calculation of production envelopes, the modelling of metabolic proteome allocation via the enzyme-constraint model variant *EC-iCH360*, enumerating and analysing elementary flux modes in the network via the reduced model variant *iCH360red*, and performing thermodynamic-based analysis using

Table 3: A summary of knowledge captured by the *i*CH360 model, as well as example simulations and analyses shown in this article.

| Model structure | Notes | Data source |
|--|---|-----------------------|
| Network (reaction stoichiometries) | Selected reactions from <i>i</i> ML1515, hand-curated | [8] |
| Annotations to external databases | Parsed from <i>i</i> ML1515, extended, and updated | [8], manual curation |
| Network graphics | Escher maps of the full model and its subsystems | |
| Biological knowledge supporting the stoichiometric model | Catalytic relationships, protein complex composition, small-molecule regulations, and others | [25] |
| Physico-chemical parameters mapped to model | | |
| Thermodynamic constants (ΔG^0) | Account for compartment-specific chemical environment (pH, pMg, ionic strength, and potential) and include corrections for reactions occurring across compartments. | [36], manual curation |
| k_{app}^{max} values | <i>in vivo</i> estimates of catalytic turnover numbers | [31] |
| Typical k_{app} values | Adjusted estimates of turnover numbers, fitted to proteomic data, accounting for typical saturation levels across growth conditions | [31] [32] |
| Protein molecular masses | Molecular masses for all proteins/protein complexes covered by the model | [25] |
| Cell-state data mapped to the model | | |
| Protein abundances | Abundance across different growth conditions for proteins/protein complexes covered by the model, estimated from experimentally measured polypeptide abundances | [32] |
| Metabolite concentrations | Measured metabolite concentrations across growth conditions | [40] |
| Metabolic fluxes | Measured metabolite fluxes for aerobic growth on glucose | [31], [41] |
| Example applications shown | | |
| Production envelope analysis | See Figures 2 and S8 | |
| Enzyme-constrained FBA | With enzyme-constrained version of the model, EC- <i>i</i> CH360, constructed in sMOMENT format [29] | |
| EFM analysis | With reduced model variant <i>i</i> CH360red. | |
| Saturation FBA | Predicts metabolic fluxes as a function of external substrate concentration | |
| Max-Min Driving force | Formulation from Noor et al. (2014) [37], extended to account for correlated uncertainty in the thermodynamic constants estimates. | |
| Probabilistic Metabolic Optimisation | Prediction of thermodynamic states (metabolite concentrations and relative fluxes) maximally consistent with measured metabolite concentrations [38] | |

the set of thermodynamic constants provided. Nevertheless, we believe that many other analyses are possible. For example, alternative definitions of elementary pathways, such as elementary conversion modes, would be valuable to explore as a more tractable alternative to elementary flux modes [42]. Similarly, our analyses of metabolic enzyme cost presented here have relied on simple, capacity-based definitions of enzyme cost, where a constant cost per unit flux is assumed for each reaction regardless of the condition studied. However, with additional kinetic parametrisation of the model, more complete enzyme cost functions, explicitly accounting for condition-specific metabolite concentrations and, consequently, enzyme costs, could be used to generate more realistic estimates of metabolic tradeoffs predicted by the model [33].

Indeed, while at this stage the parameterisation of the model is limited to turnover numbers and thermodynamic constants, we anticipate that additional parameter sets can easily be mapped to the model. Facilitated by extensive annotations present in the model and by the recent development of machine learning-enabled kinetic constant estimators [43–46], a complete kinetic parameterisation of *i*CH360 is thus a valuable potential future development. In addition, probabilistic estimates of the kinetic parameters [46, 47] can be combined with our existing thermodynamic parameterisation, making it possible to account for (potentially correlated) parameter uncertainty throughout the kinetic modelling process.

In light of the above results, we believe that *i*CH360 has the potential to become a reference metabolic model for *E. coli*.

4 Methods

4.1 Model assembly and curation

All relevant pathways included in the model (the core metabolism reactions from ECC [13] and the biosynthesis pathways shown in Supplementary Figures S1-S4) were assembled and curated based on information available in the EcoCyc [25] and KEGG [23] databases. The respective reactions were then extracted from *iML1515* and parsed into a new model. To compute an equivalent biomass reaction, we first collected all the pathways required for the production of the components present in the *iML1515* biomass reaction (with the exception of a small number of compounds present with very small stoichiometry, which were neglected from this analysis for simplicity), but not in our model. We used the "core" biomass reaction from *iML1515* (BIOMASS_Ec_iML1515_core_75p37M), rather than alternative "WT" reaction (BIOMASS_Ec_iML1515_WT_75p37M), since its smaller number of requirements made it easier to manually curate the pathways for their production. These additional pathways (available in the repository supporting this manuscript) were manually curated based on the available literature and database annotations to ensure they represent the biologically most relevant bio-production route for each biomass component. By adding these pathways to *iCH360*, we obtained an extended model that was able to predict growth rates directly through the original *iML1515* biomass reaction. The equivalent biomass function was then computed based on a reference flux distribution computed on this extended model, as explained in Supplementary Information A.2. Both growth-associated (the stoichiometry of intracellular ATP in the biomass reaction) and non-growth-associated (the lower bound on the maintenance reaction, ATPM) energy requirements were directly inherited from *iML1515*. Model assembly, manipulation and validation were performed using the COBRA Toolbox [17]. The extension of database cross-annotation for the model reactions was performed through a mixture of automated database query and manual curation.

4.2 Network graphics

iCH360 can be visualised through a series of custom-built maps using the metabolic visualisation tool Escher [18] (see Figure 1). There are three main ways to visualise the model or solutions thereof. First, a complete map of the model, including all of its reactions, can be used (Figure 1). In order to provide a more compact representation of the network, a compressed second variant of the same map was constructed (Supplementary Figure S21). Here, long biosynthetic linear pathways were lumped into single pseudo-reactions, which only show the net production or consumption of metabolites by the pathway while omitting intermediates. Finally, individual maps are provided for each of the main subsystems in the model.

4.3 Production Envelope Analysis

All production envelopes shown in the main text and Supplementary Information were generated using the built-in production envelopes tools from the `cobra.flux_analysis.phenotype_phase_plane` module in the COBRApy package [17]. Briefly, the algorithm first computes the maximum and minimum production rates of the metabolite of interest, given the existing constraints in the model (including the specified bounds on the uptake of the carbon source and oxygen). The interval between the maximum and minimum achievable production rates is then discretised into an equally spaced grid of points. For each point in the interval, the production rate is fixed and the model's objective (here, the growth rate) is sequentially maximised and minimised, thus generating the boundary of the production envelope. All prediction envelopes were computed by specifying an upper bound on the uptake of the carbon source of 10 mmol/gDW/h, and blocking oxygen uptake for the anaerobic scenario. For comparisons with ECC and ECC2, the maintenance requirement (lower bound on the ATPM reaction) of these two models was set to the same value used in *iML1515* and *iCH360* (6.86 mmol/gDW/h).

4.4 Knowledge graph for linking reactions to enzymes and proteins

Information about the enzymes and proteins behind the stoichiometric model was collected in a knowledge graph. To build this graph, all available data on reaction-protein association and subunit composition were retrieved by automatically querying the BioCyc database through its REST-based data retrieval API (<https://biocyc.org/web-services.shtml>). This information was then extended and curated based on a comparison with existing *iML1515* GPR annotations. The resulting data were used to generate a directed graph in which nodes represent biological entities (such as reactions, proteins, genes, and metabolites) and edges represent functional dependencies across them, including catalysis, subunit composition, post-translational modifications, and others. A complete list of node and edge types in the graph is provided in Supplementary Tables 2 and 3, respectively. All polypeptide nodes were annotated with their molecular mass (parsed from EcoCyc), enabling recursive computation of the molecular mass of *any* protein node in the graph (see Supplementary Information A.3). All manipulation

and analysis of the graph data structure were performed using the NetworkX Python package [48], and the final data structure is provided to the user in Cytoscape and GML formats.

4.5 Primary and secondary catalytic edges

Catalytic edges in the graph, i.e. edges connecting reaction and enzyme nodes, were manually annotated as either primary or secondary, based on available evidence for the activity of each enzyme in the model with respect to its associated reactions (Supplementary File S2). More specifically, a catalytic edge between a reaction-enzyme pair was labelled as secondary whenever the enzyme was shown in the literature to account for only minor catalytic activity for a reaction when compared with another isozyme. In this case, references to the relevant literature were included as metadata for that edge. Whenever sufficient information was not available, all isozyme edges for a reaction were conservatively treated as primary.

4.6 Catalytic disruption analysis

For the catalytic disruption test, we identified condition-specific essential reactions in the following way: in the model, for each condition considered, we determined the reactions whose knockout led to the inability to produce biomass. A small number of known false positives, which are essential in *i*CH360 only due to its lack of certain reactions or pathways, were excluded from the analysis. For each essential reaction, we considered all associated genes and investigated the removal of each of them individually from the graph. The effect of the simulated knockout was propagated across the graph by removing all nodes from the graph for which the gene is required according to the nodes' Boolean GPR rules (see Supplementary Information A.3). Finally, the result of each simulated knockout was catalogued according to the reaction-level disruption it caused (see the main text). Whenever multiple reactions were disrupted by the knockout of a gene, the strongest disruption among them was assigned to the knockout, in the following precedence order: complete disruption, full primary disruption, partial primary disruption, secondary disruption. The analysis was repeated for a total of 9 growth conditions, and each condition-gene pair, labelled with the assigned disruption type, was compared with experimentally measured mutant relative fitness (averaged across available replicates for that condition-gene combination), using the data set from Price et al. (2018) [27].

4.7 Construction of the enzyme-constrained metabolic model

In the sMOMENT formalism, the positive and negative fluxes in each reaction are formally described, respectively, as positive fluxes in separate “forward” and “backward” versions of the reaction. To construct the enzyme-constrained model in the sMOMENT format, reversible reactions in the model were duplicated in the model to separately represent fluxes in forward or backward direction. Direction-specific turnover number estimates were parsed from Heckmann et al. (2020) [31] and a default value of 65 s^{-1} was used for transporters as in [49]. To account for the fact that [31] reports values as turnover numbers *per polypeptide*, the values were multiplied by the number of polypeptide subunits in each enzyme. For each reaction-enzyme pair, an enzyme cost per unit flux (in units of $\text{g} \cdot \text{h} \cdot \text{mmol}^{-1}$) was then defined as:

$$a_i = \frac{M_i}{k_{\text{cat},i} \sigma} \quad (1)$$

where a_i is the enzyme cost of reaction-enzyme pair i , $k_{\text{cat},i}$ is the turnover rate estimate for the pair (here, in units of h^{-1}), M_i is the molecular mass of the enzyme involved (in kDa), and σ is a unitless condition-specific scaling factor (typically interpreted as an average enzyme saturation value). Then, a unique enzyme was assigned to each reaction. To this end, secondary catalytic relationships were first discarded and, for reactions with multiple annotated primary isoenzymes, the enzyme with the highest measured abundance in the integrated PAX Database [50] was heuristically chosen. Based on these costs per unit flux, an enzyme capacity constraint was introduced as:

$$\sum_i a_i v_i \leq e_{\text{tot}} \quad (2)$$

where v_i denotes the flux in reaction i (in $\text{mmol h}^{-1} \text{gDW}^{-1}$) and e_{tot} is a parameter denoting the total amount of enzyme mass (in g/gDW) that can be allocated to the flux mode. The constraint is enforced by augmenting the stoichiometric matrix of the model with an additional enzyme supply pseudoreaction (upper-bounded by e_{tot}) and an enzyme pool pseudometabolite consumed in each reaction with stoichiometry a_i [29].

4.8 Adjustment of turnover numbers across conditions

Turnover numbers from Heckmann et al. (2020) [31] were adjusted based on the nonlinear programming (NLP) formulation detailed in Supplementary Information A.5). The reference flux distributions required by the procedure were obtained using

the original (unadjusted) set of turnover numbers and the bounds on allowable adjustments (\mathbf{u}_{\min} and \mathbf{u}_{\max} in Eq. 34) set to ± 2 (corresponding to a maximal 100-fold increase or reduction of each parameter from the original value). The linear program used to obtain reference flux distributions was formulated and solved with GUROBI [51], while the nonlinear program used for turnover adjustment was formulated and solved with the open-source optimisation package CasADi [52]. To investigate the effect of the ridge regularisation hyperparameter (ρ in equation (34)), we solved the adjustment problem for a broad range of values of this parameter and, each time, we computed the RMSE between measurements and predictions of enzyme abundance, as well as the mean absolute deviation of between original and adjusted turnover values, both computed for \log_{10} -transformed data (Supplementary Figure S13). Based on this information, a value of $\rho = 1$, after which any further decrease in the amount of regularisation results in marginal reduction of the RMSE, was heuristically chosen to compute the final set of adjusted turnover numbers. Upon reparametrisation into *apparent* turnover numbers, (see Supplementary Information A.5.5), this set of adjusted parameters was used to parametrise the enzyme-constrained model variant EC-*i*CH360. For the leave-one-out cross validation, the adjustment procedure was run multiple times, each time excluding abundance data from a condition from the training dataset. The resulting adjusted parameter set was then used to generate enzyme abundance predictions for the condition left out, which were then compared against measured values.

4.9 Enzyme allocation predictions

To validate enzyme allocation predictions against experimental values, we first retrieved measured polypeptide abundances for each growth condition from Schmidt et al. (2016) [32] and imputed missing values using, whenever available, abundance values from the PAX Database [50]. Then we estimated enzyme counts across conditions from polypeptide counts using non-negative least-squares estimation (Supplementary Information A.4). Next, we converted them into mass abundances (in units of g gDW^{-1}) based on the molecular mass of each complex (see Supplementary Information A.3) and assuming a cell dry mass of $2.8 \times 10^{-13}\text{g}$ (BIONUMBER ID 103904 [53]). Enzyme allocation predictions for each condition were then computed via EC-*i*CH360 by fixing the growth rate to the experimentally measured one from Schmidt et al. (2016) [32] and minimising the total enzyme cost, initially using a value of 1 for the average saturation parameter σ (Eq. (1)). For the predictions computed using the turnover number estimates from Heckmann et al. (2020) [31], the average saturation coefficient σ was then estimated from data, for each condition, as:

$$\sigma = \frac{\sum_{i \in \mathcal{M}} e_i}{\bar{e}_{\text{tot}}} \quad (3)$$

where \mathcal{M} denotes the index set of model enzymes for which measurements are available, e_i is the predicted abundance for the i th enzyme, and \bar{e}_{tot} is the total measured model enzyme abundance for that condition. This value of σ was then used to scale the predicted enzyme abundances before comparing them with the ones measured experimentally. Note that this choice of σ ensures that the sum of predicted abundances matches that of measured ones. To compute predictions with the adjusted parameter set, the scaling factors for each condition were obtained as part of the fitting procedure (see Supplementary Information A.5). All reported root mean squared errors were computed on \log_{10} transformed enzyme abundances, excluding enzymes with zero predicted abundance from the computation.

4.10 Enumeration of Elementary Flux Modes

Elementary Flux Modes (EFMs) for the submodel *i*CH360red were enumerated, for each growth condition, using EFMtools [54]. Filtered modes (Supplementary Table 5) were defined as those supporting nonzero biomass flux and, for aerobic modes, nonzero oxygen uptake. In addition, in the aerobic case, filtered modes exclude those carrying flux in either of three reactions – Pyruvate-Formate Lyase (PFL), Fumarate Reductase (FRD2), and the menaquinone-dependent Dihydroorotate dehydrogenase (DHORD5) – known to be physiologically active only under anaerobic conditions [55–57].

4.11 Growth/yield trade-off analysis

To analyse trade-offs between growth rate and yield, we computed the yield and a predicted growth rate of each EFM, using the unit conventions common in stoichiometric metabolic models. The specific growth rate μ (1/h) of a cell or a cell population can be defined as the rate of biomass production per amount of biomass present. In stoichiometric metabolic models, units are chosen in a specific way such that the growth rate is identically given by the rate v_{BM} of the biomass reaction. In FBA models, normal reaction fluxes are given in units of mmol/gDW/h, so the yield of an EFM, computed as the ratio of its biomass flux v_{BM} to its glucose uptake flux, has a unit of gDW/mmol. The biomass flux, whose conventional unit is h^{-1} was directly interpreted as the cell growth rate μ . To determine the cell growth rate allowed by a flux distribution, we first computed its absolute enzyme

cost

$$c_{\text{enz}} = \sum_i a_i v_i \quad (4)$$

where c_{enz} is the enzyme cost of the flux distribution (an enzyme mass, measured in g/gDW), v_i is the flux of reaction i , and a_i is the enzyme cost per unit flux in reaction i , computed as per equation (1) using the set of adjusted turnover numbers. Since an elementary flux mode can be scaled arbitrarily, both v_{BM} and c_{enz} depend on the particular choice of scaling of the mode (though their ratio, $v_{\text{BM}}/c_{\text{enz}}$, does not). In order to obtain an estimate for the achievable growth rate of a mode μ that is a unique property of each EFM, independent of its scaling, we thus normalised all modes to the same total enzyme cost f_{enz} (in g/gDW), and looked at the resulting flux through the biomass reaction. More formally, the achievable growth rate μ (in h⁻¹) for an (arbitrarily scaled) EFM with biomass flux v_{BM} and enzyme cost c_{enz} was computed as:

$$\mu = f_{\text{enz}} \frac{v_{\text{BM}}}{c_{\text{enz}}} \quad (5)$$

where f_{enz} denotes the total mass of enzyme available for the flux distribution, relative to the total dry mass of the cell. For simplicity, we approximate f_{enz} by a constant value of $f_{\text{enz}} = 0.285$ g/gDW, which we obtained by taking the minimum enzyme investment required by the enzyme-constrained model to support the experimentally measured growth rate reported in the proteomic data set [32] for the condition of interest in this analysis (aerobic growth on glucose).

To simulate low-oxygen conditions, the cost per unit flux of all oxygen-consuming reactions (CYTBO3_4pp, CYTBDpp, CYTBD2pp) was increased by a 1000-fold, mimicking the physiological state in which these reactions operate at low saturation with oxygen.

4.12 Saturation FBA analysis

Saturation FBA calculations were performed by optimising biomass production in the enzyme-constrained model, setting the saturation coefficient (σ in Eq. (1) to the value fitted as part of the turnover number adjusting procedure (Section 4.8) for all reactions except the glucose transporter (GLCptspp). The saturation of the glucose transporter, σ_{up} , was computed as a function of external glucose concentration and assuming irreversible Michaelis-Menten kinetics, so that:

$$\sigma_{\text{up}} = \frac{[\text{Glc}]}{K_m + [\text{Glc}]} \quad (6)$$

where [Glc] is the external glucose concentration (in mM) and a value of 0.116 mM was used for the Michaelis constant K_m [12].

4.13 Component contribution estimates of thermodynamic constants

Estimates of the free energies of reactions, and their uncertainties, were obtained using the component contribution framework previously described in [36]. Several reactions in the model involve protein side groups, such as the acyl-carrying protein (ACP), or cofactors, such as glutaredoxin, for which a decomposition in terms of chemical groups is not available. As a result, thermodynamic constants for these reactions cannot be *directly* estimated through database-based implementations of the component contribution method, such as eQuilibrator [58]. However, since these non-decomposable protein groups are conserved in all reactions within our model, their net contribution to the reaction thermodynamics is, at least from a group contribution perspective, null. If we were only interested in computing standard free energies of reaction, $\Delta_r G^\circ$, we could simply treat protein groups as non-decomposable “black-box” units and add them to the group incidence matrix of the eQuilibrator database (see Section S1.1 in [58]), enabling us to construct a group decomposition for compounds that contain them. However, for the computation of *transformed* standard free energies of reaction, $\Delta_r G'^\circ$, an exact chemical definition of each metabolite is required.

To address this issue, protein groups were replaced by an appropriate chemical moiety that best approximates the metabolite’s chemical environment. Specifically, ACP-groups were replaced by a phosphopantetheine group, the natural prosthetic group of acyl-carrier proteins, with a methyl group at the attachment site to the protein scaffold. Similarly, glutaredoxin was replaced by its Cys-Pro-Tyr-Cys active site, with the two cysteines being either free (for the reduced form of the cofactor) or linked by a disulfide bridge (for the oxidised state). International Chemical Identifiers (InChI) were constructed for these “replacement” metabolites (Supplementary File S3) and used to extend the default compound cache of the eQuilibrator database. This custom-extended eQuilibrator compound cache is available in the repository supporting this manuscript.

Corrections for reactions that occur in different compartments were calculated as described in [58], using compartment-specific pH, pMg, ionic strength, and potentials from Gollub et al. (2023) [46].

4.14 Max-min driving force computation

The max-min driving force (MDF) of each reference flux distribution was calculated by extending the original formulation described in [37] to account for correlated uncertainty in the estimates of the thermodynamic constants. Let $\Delta G_r'^\circ \in \mathbb{R}^N$ be a random vector representing the (uncertain) standard Gibbs free energy of reaction for the reactions in the network. Importantly, this vector includes only balanced metabolic reactions and excludes pseudoreactions such as exchange reactions. To describe our uncertain knowledge, we assume that the vector $\Delta G_r'^\circ$ follows a multivariate normal distribution:

$$\Delta G_r'^\circ \sim \mathcal{N}(\bar{\Delta G}_r'^\circ, \Sigma) \quad (7)$$

where $\bar{\Delta G}_r'^\circ$ and Σ are the mean vector and covariance matrix of the estimates obtained through the component contribution methods. This random vector can equivalently be expressed as:

$$\Delta G_r'^\circ = \bar{\Delta G}_r'^\circ + \mathbf{Q} \mathbf{z} \quad (8)$$

where \mathbf{z} is a standard normal random vector in \mathbb{R}^q , with $q = \text{rank}(\Sigma)$, and $\mathbf{Q} \in \mathbb{R}^{N \times q}$ is a square root of the covariance matrix, i.e. it satisfies:

$$\Sigma = \mathbf{Q} \mathbf{Q}^T. \quad (9)$$

In order to integrate this probabilistic description within a typical constrained-optimisation formulation, we define the set \mathcal{D}_α as the α -level confidence region around the mean of $\Delta G_r'^\circ$. Using Eq. (7), and noting that the squared norm of a standard normal random variable is known to be Chi-squared distributed, we can represent this set as:

$$\mathcal{D}_\alpha = \{ \mathbf{x} \in \mathbb{R}^N \mid \mathbf{x} = \bar{\Delta G}_r'^\circ + \mathbf{Q} \mathbf{m}, \quad \|\mathbf{m}\|_2^2 \leq \chi_{q;\alpha}^2 \} \quad (10)$$

where $\mathbf{m} \in \mathbb{R}^q$ is a vector of free parameters and $\chi_q(\cdot)$ denotes the quantile function (inverse cumulative distribution function) of a chi-squared distribution with q degrees of freedom. We can thus account for the uncertainty in thermodynamic estimates by treating the free energies of reaction as decision variables, rather than known parameters, and constraining their value to belong to \mathcal{D}_α . For a given reference flux distribution \mathbf{v} , this results in the following quadratically-constrained program (QCP):

$$\begin{aligned} & \max_{\mathbf{m}, \mathbf{c}, b} \quad b \\ \text{s.t.} \quad & \Delta_r G'^\circ = \bar{\Delta G}_r'^\circ + \mathbf{Q} \mathbf{m} \\ & \Delta_r G' = \Delta_r G'^\circ + RT \mathbf{S}^T \mathbf{c} \\ & -\text{sign}(v_i) \Delta_r G'_i > b, \quad \text{if } v_i \neq 0 \\ & \|\mathbf{m}\|_2^2 \leq \chi_{q;\alpha}^2 \\ & \mathbf{c}_{\min} \leq \mathbf{c} \leq \mathbf{c}_{\max} \end{aligned} \quad (11)$$

where b is the min driving force (or the MDF after optimisation), $\mathbf{c} \in \mathbb{R}^m$ is a vector of log-metabolite concentrations, $\mathbf{c}_{\min}, \mathbf{c}_{\max} \in \mathbb{R}^m$ are lower and upper bounds on these log-concentrations, $\mathbf{S} \in \mathbb{R}^{m \times N}$ is the stoichiometric matrix of the model, R is the ideal gas constant, and T is the temperature used for the computation of the free energy estimates. For our MDF calculations, we used a confidence level of 90% and set the bounds on the concentration of all metabolites to the physiologically plausible range of (1μM, 10mM). The above quadratically constrained program was formulated and solved using the GUROBI package [51].

4.15 Probabilistic thermodynamic analysis

Probabilistic metabolic optimisation (PMO) of the model was performed using the PTA Python package [38], providing the software with the curated thermodynamic estimates generated for this model. The default values available through the package for growth in M9 medium with glucose (which include measurements from Gerosa et al. (2015) [40]) were used as priors for the concentration of metabolites, and the growth rate was bounded from below by the reported value in [40]. Furthermore, for this analysis, which requires that each flux in the solution have a well-defined directionality, two transhydrogenase reactions (NADH17pp and THD2pp) were allowed to operate reversibly [38].

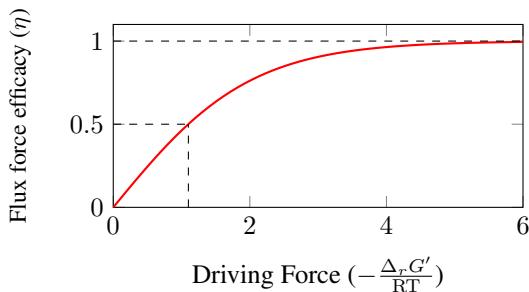


Figure 8: The flux force efficacy (η) as a function of the (scaled) negative Gibbs free energy of reaction, $-\frac{\Delta_r G'}{RT}$. The efficacy of the flux force corresponds to the ratio between the net flux (forward minus backward flux) and the total flux (forward plus backward flux) of a reaction, which approaches 1 for reactions operating far from chemical equilibrium ($\Delta_r G' \ll 0$).

To analyse the thermodynamic state computed by PTA, we compute the flux-force efficacy η for each reaction as (see Figure 8):

$$\eta = \frac{e^{-\frac{\Delta_r G'}{RT}} - 1}{e^{-\frac{\Delta_r G'}{RT}} + 1} = \tanh\left(-\frac{1}{2} \frac{\Delta_r G'}{RT}\right) \quad (12)$$

where $\Delta_r G'$ is the Gibbs free energy of reaction in the thermodynamic state, R is the universal gas constant and T is the temperature considered for the analysis (310.15 K). To compare the PTA-predicted flux force efficacies with experimental measurements of enzyme abundance, we selected reactions that carried at least 2.5% of the glucose uptake flux and pooled them into two groups, corresponding to $\eta > 0.5$ (72 reactions) and $\eta < 0.5$ (31 reactions). $\eta = 0.5$ corresponds to $\Delta_r G' \approx -1.1 RT = -2.8$ kJ/mol.

5 Author contributions

A.B.-E. conceived the project and oversaw its initial stage. M.C and H.H. assembled and curated the model. M.C. wrote software, performed simulations and analysed results. W.L., E.N. and H.H. oversaw the project and contributed to the analysis and interpretation of results. M.C., H.H., W.L., and E.N. wrote the manuscript.

6 Competing interests

The authors declare that they comply with the PCI rule of having no financial conflicts of interest with respect to the content of the article.

7 Code availability

The model, together with all relevant data, is available on Github at <https://github.com/marco-corrao/iCH360>. The code and files required to reproduce all analyses in this manuscript are available on Github at https://github.com/marco-corrao/iCH360_paper and on Zenodo at <https://doi.org/10.5281/zenodo.11092781>. Additional analyses and comparisons with other models are also available in these repositories.

References

- [1] Debolina Sarkar and Costas D. Maranas. “Engineering microbial chemical factories using metabolic models”. In: *BMC Chemical Engineering* (2019). DOI: [10.1186/s42480-019-0021-9](https://doi.org/10.1186/s42480-019-0021-9).
- [2] Changdai Gu, Gi Bae Kim, Won Jun Kim, Hyun Uk Kim, and Sang Yup Lee. “Current status and applications of genome-scale metabolic models”. In: *Genome Biology* (2019). DOI: [10.1186/s13059-019-1730-3](https://doi.org/10.1186/s13059-019-1730-3).
- [3] Jiangong Lu, Xinyu Bi, Yanfeng Liu, Xueqin Lv, Jianghua Li, Guocheng Du, and Long Liu. “In silico cell factory design driven by comprehensive genome-scale metabolic models: development and challenges”. In: *Systems Microbiology and Biomanufacturing* (2023). DOI: [10.1007/s43393-022-00117-4](https://doi.org/10.1007/s43393-022-00117-4).
- [4] Xin Fang, Colton J. Lloyd, and Bernhard O. Palsson. “Reconstructing organisms in silico: genome-scale models and their emerging applications”. In: *Nature Reviews Microbiology* (2020). DOI: [10.1038/s41579-020-00440-4](https://doi.org/10.1038/s41579-020-00440-4).
- [5] Jennifer L. Reed and Bernhard Ø. Palsson. “Thirteen years of building constraint-based in silicon models of *Escherichia coli*”. In: *Journal of Bacteriology* (2003). DOI: [10.1128/JB.185.9.2692-2699.2003](https://doi.org/10.1128/JB.185.9.2692-2699.2003).

- [6] Adam M Feist, Christopher S Henry, Jennifer L Reed, Markus Krummenacker, Andrew R Joyce, Peter D Karp, Linda J Broadbelt, Vassily Hatzimanikatis, and Bernhard Ø Palsson. “A genome-scale metabolic reconstruction for Escherichia coli K-12 MG1655 that accounts for 1260 ORFs and thermodynamic information”. In: *Molecular Systems Biology* (2007). DOI: 10.1038/msb4100155.
- [7] Jeffrey D Orth, Tom M Conrad, Jessica Na, Joshua A Lerman, Hojung Nam, Adam M Feist, and Bernhard Ø Palsson. “A comprehensive genome-scale reconstruction of Escherichia coli metabolism”. In: *Molecular Systems Biology* (2011). DOI: 10.1038/msb.2011.65.
- [8] Jonathan M. Monk, Colton J. Lloyd, Elizabeth Brunk, Nathan Mih, Anand Sastry, Zachary King, Rikiya Takeuchi, Wataru Nomura, Zhen Zhang, Hirotada Mori, Adam M. Feist, and Bernhard O. Palsson. “iML1515, a knowledgebase that computes Escherichia coli traits”. In: *Nature Biotechnology* (2017). DOI: 10.1038/nbt.3956.
- [9] David B. Bernstein, Batu Akkas, Morgan N. Price, and Adam P. Arkin. “Evaluating E. coli genome-scale metabolic model accuracy with high-throughput mutant fitness data”. In: *Molecular Systems Biology* (2023). DOI: 10.15252/msb.202311566.
- [10] Hai He, Rune Höper, Moritz Dodenhöft, Philippe Marlière, and Arren Bar-Even. “An optimized methanol assimilation pathway relying on promiscuous formaldehyde-condensing aldolases in E. coli”. In: *Metabolic Engineering* (2020). DOI: 10.1016/j.ymben.2020.03.002.
- [11] Ari Satanowski, Beau Dronsella, Elad Noor, Bastian Vögeli, Hai He, Philipp Wichmann, Tobias J. Erb, Steffen N. Lindner, and Arren Bar-Even. “Awakening a latent carbon fixation cycle in Escherichia coli”. In: *Nature Communications* (2020). DOI: 10.1038/s41467-020-19564-5.
- [12] Meike T. Wortel, Elad Noor, Michael Ferris, Frank J. Bruggeman, and Wolfram Liebermeister. “Metabolic enzyme cost explains variable trade-offs between microbial growth rate and yield”. In: *PLOS Computational Biology* (2018). DOI: 10.1371/journal.pcbi.1006010.
- [13] Jeffrey D. Orth, R. M. T. Fleming, and Bernhard Ø Palsson. “Reconstruction and Use of Microbial Metabolic Networks: the Core Escherichia coli Metabolic Model as an Educational Guide”. In: *EcoSal Plus* (2010). DOI: 10.1128/ecosalplus.10.2.1.
- [14] Oliver Hädicke and Steffen Klamt. “EColiCore2: a reference network model of the central metabolism of Escherichia coli and relationships to its genome-scale parent model”. In: *Scientific Reports* (2017). DOI: 10.1038/srep39647.
- [15] Philipp Erdrich, Ralf Steuer, and Steffen Klamt. “An algorithm for the reduction of genome-scale metabolic network models to meaningful core models”. In: *BMC systems biology* (2015). DOI: 10.1186/s12918-015-0191-x.
- [16] Nusrat Jahan, Kazuhiro Maeda, Yu Matsuoka, Yurie Sugimoto, and Hiroyuki Kurata. “Development of an accurate kinetic model for the central carbon metabolism of Escherichia coli”. In: *Microbial Cell Factories* (2016). DOI: 10.1186/s12934-016-0511-x.
- [17] Ali Ebrahim, Joshua A. Lerman, Bernhard O. Palsson, and Daniel R. Hyduke. “COBRApy: COnstraints-Based Reconstruction and Analysis for Python”. In: *BMC Systems Biology* (2013). DOI: 10.1186/1752-0509-7-74.
- [18] Zachary A. King, Andreas Dräger, Ali Ebrahim, Nikolaus Sonnenschein, Nathan E. Lewis, and Bernhard O. Palsson. “Escher: A web application for building, sharing, and embedding data-rich visualizations of biological pathways”. In: *PLOS Computational Biology* (2015). DOI: 10.1371/journal.pcbi.1004321.
- [19] Simon Boecker, Sebastián Espinel-Ríos, Katja Bettenbrock, and Steffen Klamt. “Enabling anaerobic growth of Escherichia coli on glycerol in defined minimal medium using acetate as redox sink”. In: *Metabolic Engineering* (2022). DOI: 10.1016/j.ymben.2022.05.006.
- [20] Anastasia Krivoruchko, Yiming Zhang, Verena Siewers, Yun Chen, and Jens Nielsen. “Microbial acetyl-CoA metabolism and metabolic engineering”. In: *Metabolic Engineering* (2015). DOI: 10.1016/j.ymben.2014.11.009.
- [21] Shasha Zhang, Wei Yang, Hao Chen, Bo Liu, Baixue Lin, and Yong Tao. “Metabolic engineering for efficient supply of acetyl-CoA from different carbon sources in Escherichia coli”. In: *Microbial Cell Factories* (2019). DOI: 10.1186/s12934-019-1177-y.
- [22] Peter D Karp, Richard Billington, Ron Caspi, Carol A Fulcher, Mario Latendresse, Anamika Kothari, Ingrid M Keseler, Markus Krummenacker, Peter E Midford, Quang Ong, Wai Kit Ong, Suzanne M Paley, and Pallavi Subhrawati. “The BioCyc collection of microbial genomes and metabolic pathways”. In: *Briefings in Bioinformatics* (2019). DOI: 10.1093/bib/bbx085.
- [23] Minoru Kanehisa, Miho Furumichi, Mao Tanabe, Yoko Sato, and Kanae Morishima. “KEGG: new perspectives on genomes, pathways, diseases and drugs”. In: *Nucleic Acids Research* (2017). DOI: 10.1093/nar/gkw1092.

- [24] Sébastien Moretti, Van Du T Tran, Florence Mehl, Mark Ibberson, and Marco Pagni. “MetaNetX/MNXref: unified namespace for metabolites and biochemical reactions in the context of metabolic models”. In: *Nucleic Acids Research* (2021). DOI: 10.1093/nar/gkaa992.
- [25] Ingrid M. Keseler et al. “The EcoCyc database: reflecting new knowledge about Escherichia coli K-12”. In: *Nucleic Acids Research* (2017). DOI: 10.1093/nar/gkw1003.
- [26] Peter D. Karp et al. “The EcoCyc Database (2023)”. In: *EcoSal Plus* (2023). DOI: 10.1128/ecosalplus.esp-0002-2023. eprint: <https://journals.asm.org/doi/pdf/10.1128/ecosalplus.esp-0002-2023>.
- [27] Morgan N. Price et al. “Mutant phenotypes for thousands of bacterial genes of unknown function”. In: *Nature* (2018). DOI: 10.1038/s41586-018-0124-0.
- [28] J. Babul. “Phosphofructokinases from Escherichia coli. Purification and characterization of the nonallosteric isozyme.” In: *Journal of Biological Chemistry* (1978). DOI: 10.1016/S0021-9258(17)34726-9.
- [29] Pavlos Stephanos Bekiaris and Steffen Klamt. “Automatic construction of metabolic models with enzyme constraints”. In: *BMC Bioinformatics* (2020). DOI: 10.1186/s12859-019-3329-9.
- [30] Benjamín J Sánchez, Cheng Zhang, Avlant Nilsson, Petri-Jaan Lahtvee, Eduard J Kerkhoven, and Jens Nielsen. “Improving the phenotype predictions of a yeast genome-scale metabolic model by incorporating enzymatic constraints”. In: *Molecular Systems Biology* (2017). DOI: 10.15252/msb.20167411.
- [31] David Heckmann, Anaamika Campeau, Colton J. Lloyd, Patrick V. Phaneuf, Ying Hefner, Marvic Carrillo-Terrazas, Adam M. Feist, David J. Gonzalez, and Bernhard O. Palsson. “Kinetic profiling of metabolic specialists demonstrates stability and consistency of in vivo enzyme turnover numbers”. In: *Proceedings of the National Academy of Sciences* (2020). DOI: 10.1073/pnas.2001562117.
- [32] Alexander Schmidt, Karl Kochanowski, Silke Vedelaar, Erik Ahrné, Benjamin Volkmer, Luciano Callipo, Kévin Knoops, Manuel Bauer, Ruedi Aebersold, and Matthias Heinemann. “The quantitative and condition-dependent Escherichia coli proteome”. In: *Nature Biotechnology* (2016). DOI: 10.1038/nbt.3418.
- [33] Elad Noor, Avi Flamholz, Arren Bar-Even, Dan Davidi, Ron Milo, and Wolfram Liebermeister. “The protein cost of metabolic fluxes: prediction from enzymatic rate laws and cost minimization”. In: *PLOS Computational Biology* (2016). DOI: 10.1371/journal.pcbi.1005167.
- [34] Stefan Müller, Georg Regensburger, and Ralf Steuer. “Resource allocation in metabolic networks: kinetic optimization and approximations by FBA”. In: *Biochemical Society Transactions* (2015). DOI: 10.1042/BST20150156.
- [35] Meike T. Wortel, Han Peters, Josephus Hulshof, Bas Teusink, and Frank J. Bruggeman. “Metabolic states with maximal specific rate carry flux through an elementary flux mode”. In: *The FEBS Journal* (2014). DOI: 10.1111/febs.12722.
- [36] Elad Noor, Hulda S. Haraldsdóttir, Ron Milo, and Ronan M. T. Fleming. “Consistent estimation of Gibbs energy using component Contributions”. In: *PLOS Computational Biology* (2013). DOI: 10.1371/journal.pcbi.1003098.
- [37] Elad Noor, Arren Bar-Even, Avi Flamholz, Ed Reznik, Wolfram Liebermeister, and Ron Milo. “Pathway Thermodynamics Highlights Kinetic Obstacles in Central Metabolism”. In: *PLOS Computational Biology* (2014). DOI: 10.1371/journal.pcbi.1003483.
- [38] Mattia G Gollub, Hans-Michael Kaltenbach, and Jörg Stelling. “Probabilistic thermodynamic analysis of metabolic networks”. In: *Bioinformatics* (2021). DOI: 10.1093/bioinformatics/btab194.
- [39] Daven B. Khana, Annie Jen, Evgenia Shishkova, Eashant Thusoo, Jonathan Williams, Alex Henkel, David M. Stevenson, Joshua J. Coon, and Daniel Amador-Noguez. *Thermodynamics shape the in vivo enzyme burden of glycolytic pathways*. 2025. DOI: 10.1101/2025.01.31.635972.
- [40] Luca Gerosa, Bart R. B. Haverkorn van Rijswijk, Dimitris Christodoulou, Karl Kochanowski, Thomas S. B. Schmidt, Elad Noor, and Uwe Sauer. “Pseudo-transition Analysis Identifies the Key Regulators of Dynamic Metabolic Adaptations from Steady-State Data”. In: *Cell Systems* (2015). DOI: 10.1016/j.cels.2015.09.008.
- [41] Christopher P. Long and Maciek R. Antoniewicz. “Metabolic flux responses to deletion of 20 core enzymes reveal flexibility and limits of *E. coli* metabolism”. In: *Metabolic Engineering* (2019). DOI: 10.1016/j.ymben.2019.08.003.
- [42] Tom J. Clement, Erik B. Baalhuis, Bas Teusink, Frank J. Bruggeman, Robert Planqué, and Daan H. de Groot. “Unlocking Elementary Conversion Modes: ecmtool Unveils All Capabilities of Metabolic Networks”. In: *Patterns* (2021). DOI: 10.1016/j.patter.2020.100177.
- [43] Alexander Kroll, Martin K. M. Engqvist, David Heckmann, and Martin J. Lercher. “Deep learning allows genome-scale prediction of Michaelis constants from structural features”. In: *PLOS Biology* (2021). DOI: 10.1371/journal.pbio.3001402.

- [44] Feiran Li, Le Yuan, Hongzhong Lu, Gang Li, Yu Chen, Martin K. M. Engqvist, Eduard J. Kerkhoven, and Jens Nielsen. “Deep learning-based kcat prediction enables improved enzyme-constrained model reconstruction”. In: *Nature Catalysis* (2022). DOI: 10.1038/s41929-022-00798-z.
- [45] Alexander Kroll, Yvan Rousset, Xiao-Pan Hu, Nina A. Liebrand, and Martin J. Lercher. “Turnover number predictions for kinetically uncharacterized enzymes using machine and deep learning”. In: *Nature Communications* (2023). DOI: 10.1038/s41467-023-39840-4.
- [46] Mattia G. Gollub, Thierry Backes, Hans-Michael Kaltenbach, and Jörg Stelling. ENKIE: A package for predicting enzyme kinetic parameter values and their uncertainties. 2023. DOI: 10.1101/2023.03.08.531697.
- [47] Timo Lubitz and Wolfram Liebermeister. “Parameter balancing: consistent parameter sets for kinetic metabolic models”. In: *Bioinformatics* (2019). DOI: 10.1093/bioinformatics/btz129.
- [48] Aric Hagberg, Pieter J. Swart, and Daniel A. Schult. “Exploring network structure, dynamics, and function using NetworkX. Available at : <https://www.osti.gov/biblio/960616>”. In: *Proceedings of the 7th Python in Science Conference (SciPy 2008)*. 2008.
- [49] David Heckmann, Colton J. Lloyd, Nathan Mih, Yuanchi Ha, Daniel C. Zielinski, Zachary B. Haiman, Abdelmoneim Amer Desouki, Martin J. Lercher, and Bernhard O. Palsson. “Machine learning applied to enzyme turnover numbers reveals protein structural correlates and improves metabolic models”. In: *Nature Communications* (2018). DOI: 10.1038/s41467-018-07652-6.
- [50] Qingyao Huang, Damian Szkłarczyk, Mingcong Wang, Milan Simonovic, and Christian von Mering. “PaxDb 5.0: Curated Protein Quantification Data Suggests Adaptive Proteome Changes in Yeasts”. In: *Molecular & Cellular Proteomics* (2023). DOI: 10.1016/j.mcpro.2023.100640.
- [51] Gurobi Optimization, LLC. *Gurobi Optimizer Reference Manual*. 2023.
- [52] Joel A E Andersson, Joris Gillis, Greg Horn, James B Rawlings, and Moritz Diehl. “CasADi – A software framework for nonlinear optimization and optimal control”. In: *Mathematical Programming Computation* (2019). DOI: 10.1007/s12532-018-0139-4.
- [53] Ron Milo, Paul Jorgensen, Uri Moran, Griffin Weber, and Michael Springer. “BioNumbers—the database of key numbers in molecular and cell biology”. In: *Nucleic Acids Research* (2010). DOI: 10.1093/nar/gkp889.
- [54] Marco Terzer and Jörg Stelling. “Large-scale computation of elementary flux modes with bit pattern trees”. In: *Bioinformatics* (2008). DOI: 10.1093/bioinformatics/btn401.
- [55] Wenhai Zhang, Kenny K. Wong, Richard S. Magliozzo, and John W. Kozarich. “Inactivation of Pyruvate Formate-Lyase by Dioxygen: Defining the Mechanistic Interplay of Glycine 734 and Cysteine 419 by Rapid Freeze-Quench EPR”. In: *Biochemistry* (2001). DOI: 10.1021/bi002589k.
- [56] Gary Cecchini, Imke Schröder, Robert P Gunsalus, and Elena Maklashina. “Succinate dehydrogenase and fumarate reductase from *Escherichia coli*”. In: *Biochimica et Biophysica Acta (BBA) - Bioenergetics* (2002). DOI: 10.1016/S0005-2728(01)00238-9.
- [57] Stuart Andrews, Graeme B. Cox, and Frank Gibson. “The anaerobic oxidation of dihydroorotate by *Escherichia coli* K-12”. In: *Biochimica et Biophysica Acta (BBA) - Bioenergetics* (1977). DOI: 10.1016/0005-2728(77)90197-9.
- [58] Moritz E Beber, Mattia G Gollub, Dana Mozaffari, Kevin M Shebek, Avi I Flamholz, Ron Milo, and Elad Noor. “eQuilibrator 3.0: a database solution for thermodynamic constant estimation”. In: *Nucleic Acids Research* (2022). DOI: 10.1093/nar/gkab1106.
- [59] Tania Bizouarn, Gijs I. van Boxel, Tina Bhakta, and J. Baz Jackson. “Nucleotide binding affinities of the intact proton-translocating transhydrogenase from *Escherichia coli*”. In: *Biochimica et Biophysica Acta (BBA) - Bioenergetics* (2005). DOI: 10.1016/j.bbabi.2005.04.004.
- [60] Lior Zelcbuch, Manuel Razo-Mejia, Elad Herz, Sagit Yahav, Niv Antonovsky, Hagar Kroytoro, Ron Milo, and Arren Bar-Even. “An in vivo metabolic approach for deciphering the product specificity of glycerate kinase proves that both *E. coli*’s glycerate kinases generate 2-phosphoglycerate”. In: *PLOS ONE* (2015). DOI: 10.1371/journal.pone.0122957.
- [61] Oliver Bartsch, Martin Hagemann, and Hermann Bauwe. “Only plant-type (GLYK) glycerate kinases produce d-glycerate 3-phosphate”. In: *FEBS Letters* (2008). DOI: 10.1016/j.febslet.2008.07.038.
- [62] Costas D. Maranas and Ali R. Zomorrodi. *Optimization methods in metabolic networks*. Wiley, 2016.

- [63] H. Yu, X. Li, F. Duchoud, D. S. Chuang, and J. C. Liao. “Augmenting the Calvin-Benson-Bassham cycle by a synthetic malyl-CoA-glycerate carbon fixation pathway”. In: Nature communications (2018). DOI: 10 . 1038 / s41467 - 018 - 04417 - z.
- [64] Helmut Bergler, Sandra Fuchsbichler, Gregor Högenauer, and Friederike Turnowsky. “The enoyl-[acyl-carrier-protein] reductase (FabI) of Escherichia coli, which catalyzes a key regulatory step in fatty acid biosynthesis, accepts NADH and NADPH as cofactors and is inhibited by palmitoyl-CoA”. In: European Journal of Biochemistry (1996). DOI: 10 . 1111 / j . 1432 - 1033 . 1996 . 0689r . x.
- [65] J. T. Tsay, W Oh, T. J. Larson, S Jackowski, and C. O. Rock. “Isolation and characterization of the beta-ketoacyl-acyl carrier protein synthase III gene (fabH) from Escherichia coli K-12.” In: Journal of Biological Chemistry (1992). DOI: 10 . 1016/S0021 - 9258(19)50498 - 7.
- [66] A Sirko, M Zatyka, E Sadowsy, and D Hulanicka. “Sulfate and thiosulfate transport in Escherichia coli K-12: evidence for a functional overlapping of sulfate- and thiosulfate-binding proteins.” In: Journal of Bacteriology (1995).
- [67] Li Zhang, Wangshu Jiang, Jie Nan, Jonas Almqvist, and Yafei Huang. “The Escherichia coli CysZ is a pH dependent sulfate transporter that can be inhibited by sulfite”. In: Biochimica et Biophysica Acta (BBA) - Biomembranes (2014). DOI: 10 . 1016 / j . bbamem . 2014 . 03 . 003.

Appendix A Supplementary Information

A.1 Changes to reactions from the model *iML1515*

After assembling the model as a subset of reactions from the genome-scale model *iML1515*, a few minor corrections were applied to some of the reactions based on evidence gathered from the literature. Note, however, that these changes did not result from an exhaustive process of review of the parent model. The corrections include:

- In *iCH360*, the membrane-bound transhydrogenase reaction (THD2pp) only translocates 1 proton across the periplasmic membrane, as opposed to the 2 protons translocated in the *iML1515* reaction [59]
- The gene-protein rule (GPR) of gene *glxK* (b0514) was reassigned to GLYCK2, which produces 2-phosphoglycerate. Meanwhile, the reaction GLYCK in from *iML1515* (which produces 3-phosphoglycerate) was removed [60, 61]
- The NAPH-dependent homoserine dehydrogenase reaction (HSDy) was made irreversible towards the homoserine production direction [10]. To avoid having a reaction irreversible in the backward direction, substrate and products of the reactions were flipped.
- The succinate transport reaction SUCCt1pp was made irreversible in the export direction, enforcing the use of the more thermodynamically favourable SUCCt2_2pp (which translocates two protons instead of one) for succinate import. To avoid having a reaction irreversible in the backward direction, substrate and products of SUCCt1pp were flipped.

A.2 Computation of an equivalent biomass reaction

In constraint-based models of metabolism, the biomass reaction summarises the production of all molecules that are not explicitly described by the model. These are typically macromolecules such as proteins or polynucleotides. Which metabolites are drained from the model network towards these other parts of metabolism, and in what proportions, depends on what compounds are described by the model. Hence, starting from an existing model, taking only a subset of the reactions but leaving the biomass reaction unchanged would lead to inconsistencies.

To construct a biomass reaction for *iCH360* that corresponds equivalently to the biomass reaction in *iML1515*, the following method was used. First, we collected all the pathways required for the production of the components present in the *iML1515* biomass reaction (with the exception of a small number of compounds present with very small stoichiometry, which were excluded from this analysis for simplicity), but not in our model. These additional pathways (available in the repository supporting this manuscript) were manually curated based on available literature and database annotation to ensure they represent the most biologically-relevant bioproduction route for each biomass component. By adding these pathways to *iCH360*, we obtained an extended model (*iCH360_{ext}*), which was able to predict growth directly through the original biomass reaction.

With this extended model at hand, we can continue as follows. Let N and M denote the number of reactions and metabolites, respectively, in *iCH360*. Further, let $N_{\text{ext}} > N$ denote the number of reactions in the extended model *iCH360_{ext}*. A reference flux distribution $\mathbf{v}_{\text{ref}}^* \in \mathbb{R}^{N_{\text{ext}}}$ is computed on the extended model through FBA. We can express this flux vector as:

$$\mathbf{v}_{\text{ref}}^* = \begin{pmatrix} \mathbf{v}^* \\ \mathbf{v}_+^* \\ v_{\text{BM}}^* \end{pmatrix} \quad (13)$$

where \mathbf{v}^* is the subset of its fluxes corresponding to reactions present in *iCH360*, \mathbf{v}_+^* is the subset of fluxes corresponding to the reactions added to create the extended model *iCH360_{ext}*, and v_{BM}^* is the flux through the *iML1515* biomass reaction. If the fluxes \mathbf{v}^* were to be imposed on *iCH360* (which, at this stage, does not yet contain a biomass reaction), a number of metabolites would necessarily remain unbalanced, that is:

$$\mathbf{S} \mathbf{v}^* \neq \mathbf{0} \quad (14)$$

where $\mathbf{S} \in \mathbb{R}^{M \times N}$ denotes the stoichiometric matrix of *iCH360*. We now seek to define an "equivalent" biomass reaction that, when added to the *iCH360* network, will drain or produce metabolites so as to balance equation (14). Let $\mathbf{r}_{\text{eq}} \in \mathbb{R}^M$ denote the stoichiometry of such biomass reaction. We compute \mathbf{r}_{eq} as the additional column of \mathbf{S} required to balance the system while

achieving the same biomass flux of the reference distribution. That is, \mathbf{r}_{eq} satisfies:

$$(\mathbf{S} \mathbf{r}_{\text{eq}}) \begin{pmatrix} \mathbf{v}^* \\ v_{\text{BM}}^* \end{pmatrix} = \mathbf{0} \quad (15)$$

which can be solved as:

$$\mathbf{r}_{\text{eq}} = -\frac{1}{v_{\text{BM}}^*} \mathbf{S} \mathbf{v}^* \quad (16)$$

The stoichiometry of this equivalent biomass reaction depends on the choice of reference flux distribution $\mathbf{v}_{\text{ref}}^*$ and, generally, may be different for different growth conditions. This results from the fact that, depending on the condition, the extended model may produce the same biomass component through different pathways, which would then be converted by our procedure into different equivalent costs of precursors in the sub-model. Nevertheless, the additional biosynthesis pathways we used to form the extended model do not allow for alternative routes to biomass, hence we found the equivalent biomass reaction to be, in this case, unique across conditions.

A.3 Computation of attributes using the knowledge graph

Using the knowledge graph supporting the stoichiometric model of *iCH360* (Section 3 in the main text), a number of useful properties can be computed based on simple operations. Here, we outline how such an approach was used to i) compute the molecular mass of all enzyme complexes in the model, based on known masses for all polypeptides and ii) construct the boolean rules (GPRs) linking all reactions and proteins in the graph to the model genes. For a description of the different types of nodes and edges mentioned in this section, see Supplementary Tables 2 and 3.

A.3.1 Computation of molecular masses for all protein nodes

In order to compute the molecular masses of all protein nodes in the graph, the protein nodes corresponding to polypeptides were first annotated with their molecular masses, readily available on the EcoCyc database. Then the molecular masses of all other protein nodes are estimated recursively as follows. Let \mathcal{I} denotes the index set of all protein nodes and $\mathcal{C}(i) \in \mathcal{I}$ the index set of protein components of node i , i.e. all nodes connected to node i by a subunit composition relationship. The molecular mass of any protein node i , M_i , is computed as:

$$M_i = \begin{cases} \bar{M}_i & \text{if node } i \text{ is a polypeptide} \\ \sum_{k \in \mathcal{C}(i)} w_{ik} M_k & \text{otherwise} \end{cases} \quad (17)$$

where \bar{M}_i denotes the (known) molecular mass of polypeptide node i and w_{ik} denotes the weight of the edge between i and k .

A.3.2 Computation of gene-protein-reaction rules

Boolean gene-protein-reaction (GPR) rules are a widely used tool defining a map between a genotype (set of active genes) and a phenotype (set of active reactions) in a metabolic model. Conventionally, this is achieved by assigning to each reaction a boolean expression given in terms of genes in the model. In this section, we show how such expressions were computed for the reactions in *iCH360* using the knowledge graph.

Starting from the leaves of the graph (genes), we construct, for each node, a boolean expression describing its state (active/inactive, corresponding to a boolean *true* / *false*) in terms of its children. The exact form of this boolean expression depends on the type of the node (reaction, protein, or logical) and the type of edges connecting it to its neighbours (Figure S22). Particularly:

- A polypeptide node is active if its associated gene is also active (Figure S22A, left).
- A multimeric protein is active if all of its subunits (the child nodes connected to it by a "subunit composition" edge) are active (Figure S22A, middle).
- A modified protein is active if its unmodified form (the child node connected to it by a "protein modification" edge) and its modification requirements (the child nodes connected to it by a "protein modification requirement" edge) are active (Figure S22A, right).
- A logical AND (logical OR) node is active if all (any) of its child nodes are active (Figure S22B).

- Finally, a reaction node is active if any of its catalysing isoenzymes (the child nodes connected to it via a "catalysis" edge) are active and, at the same time, all of its catalytic requirements (the child nodes connected to it via a "non-catalytic requirement" edge) are also active (Figure S22C).

Using these definitions, the boolean expression describing the state of a reaction can be written, ultimately, solely in terms of genes, enabling computation of conventional GPR rules and their incorporation in the standard metabolic model.

A.4 Estimation of enzyme complex abundances

In order to estimate the abundances of all enzymes in the model from proteomics data, we use the model graph (Main Text, Section 4.4) to construct a stoichiometric map between enzymes and polypeptides. This map takes the form of a matrix $\mathbf{E} \in \mathbb{R}^{n \times m}$, where n is the number of enzymes in the model and $m \geq n$ the number of polypeptides, such that \mathbf{E}_{ij} denotes the stoichiometry of polypeptide j in enzyme i . Some polypeptides may be part of additional enzyme complexes that are not part of the model. Using the available annotation to the EcoCyc database, we identified 7 such polypeptides, mapping to 9 out-of-model complexes. If these out-of-model complexes were not taken into account, the abundance of model enzymes to which these polypeptides map would be overestimated. Hence, we constructed a matrix $\hat{\mathbf{E}}$ by augmenting \mathbf{E} with additional rows corresponding to the identified out-of-model complexes.

With this mapping at hand, we assume that polypeptide abundances \mathbf{p} are related to enzyme abundances \mathbf{e} (including the required additional complexes not accounted for in the model) by

$$\mathbf{p} = \hat{\mathbf{E}}^\top \mathbf{e} \quad (18)$$

Hence, given a vector of experimental measurements of polypeptide abundances $\bar{\mathbf{p}}$, we estimate enzyme abundances by solving the nonnegative least square (NNLS) problem:

$$\begin{aligned} \min_{\mathbf{e}} \quad & \|\bar{\mathbf{p}} - \hat{\mathbf{E}}^\top \mathbf{e}\|_2^2 \\ \text{s.t.} \quad & \mathbf{e} \geq 0 \end{aligned} \quad (19)$$

A.5 Adjustment of turnover numbers based on proteomics measurements across conditions

In this section, we outline the procedure used to adjust the turnover numbers in EC-iCH360 by fitting proteomic measurements across conditions (see Section 2.4 in the main text). Briefly, our aim is to adjust the turnover numbers that parametrise the model so that enzyme allocation predictions obtained through the enzyme-constrained formulation of FBA (see Methods in the main text) match more closely experimental measurements of enzyme abundances. By *simultaneously* fitting experimental measurements across many growth conditions, we improve the robustness of the fitting procedure to experimental error and generate a condition-independent set of "typical" apparent turnover numbers that predict average trends of enzyme allocation across conditions. The output of our procedure is a set of typical enzyme efficiencies, one for each enzyme, estimated from proteomic data across conditions, as well as a set of condition-specific scaling factors that account for differences in total measured enzyme abundances between conditions. In section A.5.1 we rigorously define these parameters and state the main assumption underlying our heuristic. In sections A.5.3, A.5.4, and A.5.5 we formulate a two-steps optimisation problem whose solution, upon a suitable reparameterisation, yields data-fitted estimates of the desired parameters.

A.5.1 Preliminaries

Consider an enzyme i in a given metabolic state j . This enzyme catalyses a metabolic flux v_{ij} (in (mol/gDW)/h) given by

$$v_{ij} = \kappa_{ij} c_{ij}^{\text{enz}} \quad (20)$$

where c_{ij}^{enz} is the enzyme concentration and κ_{ij} is the enzyme efficiency (or "apparent turnover number k_{app} "). The efficiency κ_{ij} is a positive rate (here in 1/h, but usually reported in 1/s). Since, by definition, it must be lower than the enzyme's turnover number $k_{\text{cat},i}$, we write it as

$$\kappa_{ij} = \sigma_{ij} k_{\text{cat},i} \quad (21)$$

where $\sigma_{ij} \in [0, 1]$ is a unitless "capacity usage" factor. In enzyme-constrained models, enzymes are often expressed by their mass abundance $e_i = M_i c_i^{\text{enz}}$ (in g/gDW) instead of enzyme concentrations, where M_i is the enzyme molecular mass (in g/mol), so we can write the flux as

$$v_{ij} = \frac{1}{a_{ij}} e_{ij} \quad (22)$$

where a_{ij} is the enzyme cost per catalysed flux, given by the molecular mass M_i divided by the enzyme efficiency κ_{ij} .

In principle, the capacity usage factors σ_{ij} (and therefore efficiencies) of enzymes may freely vary between growth conditions. However, as a heuristic, we here assume that they can be approximated as the product of two factors: an enzyme-specific term and a condition-specific one, that is:

$$\sigma_{ij} \approx \sigma_i \cdot \tau_j \quad (23)$$

Here, the enzyme-specific factor σ_i denotes the "typical" capacity usage factor of our enzyme in the range of conditions studied. The condition-specific term, τ_j , is a unitless scaling factor that simultaneously increases or decreases the efficiencies of all enzymes depending on the cell's growth conditions. By convention, we assume that the values of τ_j are centered around 1. Substituting (23) in (21), we obtain an equivalent expression (under our assumptions) for the apparent turnover number κ_{ij} :

$$\kappa_{ij} = \underbrace{\sigma_i k_{\text{cat},i}}_{\kappa_i} \tau_j \quad (24)$$

where κ_i is a "typical apparent turnover number" that is condition-independent. Practically, the above assumption allows us to simplify the problem by reducing the number of parameters to be fitted from $I \cdot J$ to $I + J$, where I and J denote the total number of enzymes and conditions considered, respectively. Our heuristic assumption corresponds to the idea that "high-quality carbon sources" allow a cell to establish metabolic states in which enzyme efficiencies are generally high, allowing for large fluxes per enzyme abundance in all the reactions, and therefore for high cell growth rates. Probably, this heuristics would fail in some other cases, e.g. cases in which enzymes are specifically perturbed by enzyme inhibitors. But in fact, it turns out that our model, assuming a single "typical apparent turnover number" $k_{\text{app},i}$ for each enzyme, yields very good enzyme allocation predictions. This is what would be expected if our heuristic assumption were correct, and therefore supports our heuristic prediction.

We now describe how the estimates of (enzyme-specific) efficiencies κ_i and (condition specific) scaling factors τ_j for our model were obtained from model simulations and proteomics data.

A.5.2 Overview of the fitting procedure

Fitting the typical turnover parameters to proteomic data is, in general, not simple. Due to the linear programming formulation underlying the enzyme-constrained FBA problem, optimal flux distributions (and, by direct consequence, enzyme allocation predictions) are discontinuous over the turnover parameter space, making derivative-based searches through this space problematic from a numerical perspective. Similarly, the high dimensionality of the parameter space limits the applicability of gradient-free optimisation algorithms. Hence, we restrict ourselves to the (comparatively simpler) problem of adjusting turnover parameters for a fixed set of reference flux distributions across growth conditions, constraining our search to the portion of parameter space in which these reference flux vectors remain optimal for their respective conditions. This simplified fitting procedure thus consists of two steps. In the first part of the procedure, we use an initial parameter set to compute a set of reference flux distributions (one per growth condition) using enzyme-constrained FBA. In the second part, turnover parameters are fitted based on these reference flux distributions and experimental measurements of enzyme abundances.

A.5.3 Obtaining a set of reference flux distributions

Following the sMOMENT formulation of enzyme-constrained FBA [29], we consider a metabolic network with N reactions and M metabolites where all metabolic fluxes are positive (i.e. reversible reactions are split into forward and backwards components) and at most one enzyme is associated with each reaction (see Methods in the main text for more information about how such unique mapping between reactions and catalyzing enzymes was generated in our case). Hence, we assume that the enzyme cost required to sustain flux v_i for a given growth condition j is given by $a_{ij} v_i$, where the cost per unit flux a_{ij} is given by:

$$a_{ij} = \begin{cases} \frac{M_i}{\kappa_{ij}} & \text{if reaction } i \text{ is enzymatic} \\ 0 & \text{otherwise} \end{cases} \quad (25)$$

Here, M_i is the molecular mass of the enzyme associated with the reaction, κ_{ij} is the condition-dependent enzyme efficiency, as defined in 24. Given an initial guess for the value of each κ_{ij} , we compute a reference flux distribution for the j th growth

condition, \mathbf{v}_j^* by fixing the biomass flux, v_{BM} to the experimentally measured rate and minimizing the total enzyme cost:

$$\begin{aligned} \mathbf{v}_j^* &= \arg \min_{\mathbf{v}} \mathbf{a}_j^\top \mathbf{v} \\ \text{s.t. } \mathbf{S} \mathbf{v} &= \mathbf{0} & \text{(a)} \\ \mathbf{B}_j \mathbf{v} &\leq \mathbf{b}_j & \text{(b)} \\ v_{\text{BM}} &= \bar{v}_{\text{BM},j} & \text{(c)} \\ \mathbf{v} &\geq \mathbf{0} \end{aligned} \quad (26)$$

Here, the objective $\mathbf{a}_j^\top \mathbf{v}$ is the total enzyme cost for the j th growth condition, $\bar{v}_{\text{BM},j}$ is the experimentally measured growth rate for the condition, $\mathbf{S} \in \mathbb{R}^{M \times N}$ is the stoichiometric matrix of the network and $\mathbf{B}_j \in \mathbb{R}^{P \times N}$ and $\mathbf{b}_j \in \mathbb{R}^P$ are a matrix and a vector, respectively, encoding any desired upper bound (or positive lower bound) on the fluxes for the growth condition. Noting that constraint 26c can be equivalently cast as a double inequality, we rewrite the problem in the more general form:

$$\begin{aligned} \mathbf{v}_j^* &= \arg \min_{\mathbf{v}} \mathbf{a}_j^\top \mathbf{v} \\ \text{s.t. } \mathbf{S} \mathbf{v} &= \mathbf{0} & \text{(a)} \\ \hat{\mathbf{B}}_j \mathbf{v} &\leq \hat{\mathbf{b}}_j & \text{(b)} \\ \mathbf{v} &\geq \mathbf{0} \end{aligned} \quad (27)$$

where the biomass flux is assumed to be the last component of the flux vector and the augmented matrices

$$\hat{\mathbf{B}}_j \equiv \begin{pmatrix} \mathbf{B}_j \\ [0, \dots, 1] \\ [0, \dots, -1] \end{pmatrix} \quad \hat{\mathbf{b}}_j \equiv \begin{pmatrix} \mathbf{b}_j \\ \bar{v}_{\text{BM},k} \\ -\bar{v}_{\text{BM},k} \end{pmatrix} \quad (28)$$

were introduced.

In order to solve the linear program 26, we shall consider an initial guess for the turnover parameters and assume that, for each condition, all enzymes operate at the same saturation level, so that $\sigma_{ij} \equiv \bar{\sigma}_j$. Since the optimal flux distribution obtained as a solution of problem (27) is unchanged by the choice of $\bar{\sigma}_j$ (as this merely amounts to a scaling of the objective function), we can simply set $\bar{\sigma}_j = 1$ (that is, set $\kappa_{ij} = k_{\text{cat},i}$ for all conditions at this stage).

Solving (27) for each of the J growth conditions available in the experimental dataset, we obtain a set of optimal flux distributions $\mathcal{V}^* = \{\mathbf{v}_1^*, \dots, \mathbf{v}_J^*\}$, which we will use as a reference in the next step.

A.5.4 Estimating typical enzyme efficiencies for the reference set of flux distributions

We now turn to fitting the relevant parameters against experimental measurements of enzyme abundances. For this purpose, we shall express the efficiency of each enzyme-condition pair as:

$$\kappa_{ij} = \theta_j k_i \quad (29)$$

where θ_j is a condition-specific scaling term that simultaneously scales all efficiencies for a given condition, while k_i is an "adjusted" turnover parameter that simultaneously accounts for inaccuracies in the original turnover parameter numbers as well as differences in typical saturation across enzymes (and, hence, it's not formally a turnover number). While the above parameterization differs from the one in equation (24) – the terms have a different interpretation! – it will greatly simplify the notation of the fitting problem, as it allows us to easily distinguish between "global" adjustments (those affecting all enzymes in a condition), which we wish to pick freely, and "local" adjustments (affecting the efficiency of individual enzymes), which instead we wish to regularise. While this factorisation of saturation effects is merely a computational convenience and not necessarily biologically meaningful, we will retrieve the parameters in Eq. (24) from those in Eq. (29) upon a simple reparameterisation, as we detail in Section A.5.5

To formulate our fitting procedure, we will denote with \mathbf{p} a vector of \log_{10} adjusted turnover numbers (i.e. $p_i = \log_{10} k_i$) and with \mathbf{s} a vector of condition-specific \log_{10} -scaling factor (i.e. $s_j = \log_{10} \theta_j$). Further, we denote with $\bar{\mathbf{p}}$ the vector of original log-turnover numbers (i.e. the one used to obtain the reference flux distributions). For the choice of reference flux distribution

computed in the previous step, the abundance of the i th enzyme in condition j , e_{ij} , is then a function of \mathbf{p} and \mathbf{s} :

$$e_{ij}(\mathbf{p}, \mathbf{s}) = \sum_k a_{kj} v_{kj}^* \quad (30)$$

where the summation index k runs across all reactions catalysed by enzyme i and the flux cost of reaction k in condition j , a_{kj} , is computed as in (25). From the formulation of the linear program (27), there must exist a region \mathcal{S}_j of log-turnover parameter space such that $\bar{\mathbf{p}} \in \mathcal{S}_j$ and that, for every $\mathbf{p} \in \mathcal{S}_j$, $\mathbf{v}_j^* \in \mathcal{V}^*$ is the optimal solution of problem (27) for its growth condition. Hence, in this step of the adjustment procedure, we aim to find a set of typical log efficiencies, \mathbf{p}^* and log scaling factors, \mathbf{s}^* by minimising the discrepancy between enzyme abundance predictions and measurements, constraining our search to this region of the parameter space $\mathcal{S} \equiv \bigcap_j \mathcal{S}_j$ where *all* reference flux distributions are optimal for their respective growth condition:

$$\begin{aligned} (\mathbf{p}^*, \mathbf{s}^*) = \arg \min_{\mathbf{p}, \mathbf{s}} \quad & \frac{1}{N_e} \sum_{i,j} [l(e_{ij}) - l(\bar{e}_{ij})]^2 + \frac{\rho}{N_p} \sum_i (p_i - \bar{p}_i)^2 \\ \text{s.t.} \quad & \mathbf{u}_{\min} \leq \mathbf{p} - \bar{\mathbf{p}} \leq \mathbf{u}_{\max} \\ & \mathbf{p} \in \mathcal{S} \end{aligned} \quad (31)$$

where N_e is the number of enzyme-condition pairs, N_p is the number of turnover parameters, \bar{e}_{ij} is the experimental measurement of enzyme i in condition j , \mathbf{u}_{\min} and \mathbf{u}_{\max} are bounds on the allowable adjustment, $\rho > 0$ is a scalar hyperparameter, and the function $l(\cdot)$ is defined as:

$$l(x) = \begin{cases} \log_{10}(x) & x > 0 \\ 0 & x = 0 \end{cases} \quad (32)$$

The objective function of the nonlinear program (31) is a combination of two terms. The first term penalises the mean squared deviation between measurements and predictions of log-enzyme abundance. Note that the above definition of $l(\cdot)$ implies that, for each condition, only enzymes with nonzero predicted abundance are included in this term. The second term is a regularisation expression (whose strength is controlled by the hyperparameter ρ) penalising the mean squared deviation between the adjusted turnover parameters and the original parameter set. The latter term mainly serves two purposes: first, it ensures that, whenever a turnover parameter is "free" in the problem (which can happen if its associated reaction fluxes are 0 across all conditions, or if no experimental measurements are available for its associated enzyme), it will be kept at its original value; secondly, it provides a mean to tune the strength of the adjustment procedure.

In order to define the region \mathcal{S} , we shall exploit the sufficient optimality conditions of LP (27). Introducing, for each growth condition, the vectors of dual variables, $\boldsymbol{\lambda}_j \in \mathbb{R}^M$ and $\boldsymbol{\mu}_j \in \mathbb{R}^P$, corresponding to constraints 27b and 27c, respectively, we note that problem (27) admits (for the j th growth condition) a dual problem in the form:

$$\begin{aligned} \max_{\boldsymbol{\lambda}_j, \boldsymbol{\mu}_j} \quad & -\hat{\mathbf{b}}_j^\top \boldsymbol{\mu}_j \\ \text{s.t.} \quad & \mathbf{S}^\top \boldsymbol{\lambda}_j + \hat{\mathbf{B}}_j^\top \boldsymbol{\mu}_j + \mathbf{a}_j \geq \mathbf{0} \\ & \boldsymbol{\mu}_j \geq \mathbf{0} \end{aligned} \quad (a) \quad (33)$$

A well-known result in linear programming duality theory [62] states that a flux distribution is the optimal solution of the primal problem (27), if, and only if, the dual problem (33) is feasible (dual feasibility) and its optimal objective coincides with the primal optimal objective, that is $\mathbf{a}_j^\top \mathbf{v}_j^* = -\mathbf{b}_j^\top \boldsymbol{\mu}_j$ (strong duality). Taken together, dual feasibility and strong duality thus define the region of optimality in turnover parameter space of each reference flux distribution. Hence, we can integrate the above definition of \mathcal{S} within problem (31) by introducing the two vectors of dual variables ($\boldsymbol{\lambda}_j$ and $\boldsymbol{\mu}_j$) for each condition as additional optimisation variables, and simultaneously enforcing the optimality for each reference distribution. By doing this, we obtain the final formulation of the nonlinear program for turnover number adjustment, which we solved to obtain the results shown in the main text:

$$\begin{aligned} (\mathbf{p}^*, \mathbf{s}^*) = \arg \min_{\mathbf{p}, \mathbf{s}} \quad & \frac{1}{N_e} \sum_{ij} [l(e_{ij}) - l(\bar{e}_{ij})]^2 + \frac{\rho}{N_p} \sum_i (p_i - \bar{p}_i)^2 \\ \text{s.t.} \quad & \mathbf{u}_{\min} \leq \mathbf{p} - \bar{\mathbf{p}} \leq \mathbf{u}_{\max} \\ & \mathbf{S}^\top \boldsymbol{\lambda}_j + \hat{\mathbf{B}}_j^\top \boldsymbol{\mu}_j + \mathbf{a}_j \geq \mathbf{0} \quad j = 1, \dots, J \\ & \boldsymbol{\mu}_j \geq \mathbf{0} \quad j = 1, \dots, J \\ & \mathbf{a}_j^\top \mathbf{v}_j^* = -\hat{\mathbf{b}}_j^\top \boldsymbol{\mu}_j \quad j = 1, \dots, J \end{aligned} \quad (34)$$

A.5.5 Conversion to apparent turnover numbers

The above procedure produces (after conversion back to a linear scale) a set of adjusted turnover parameters, \mathbf{k}^* and scalings, θ^* . In order to retrieve the typical enzyme efficiencies and condition-specific scaling factors introduced in 24, we simply parametrise the solution by factorising the scaling terms θ_j^* as

$$\theta_j^* \equiv \bar{\sigma}^* \tau_j^* \quad (35)$$

Here, $\bar{\sigma}^*$ is the geometric mean of the scalings across conditions, which we interpret as typical enzyme saturation level across conditions, while τ_j^* is a residual scaling factor fluctuating around 1 which is required to account for differences in total measured enzymes between conditions. The typical enzyme efficiencies (κ_i) introduced in 24 are then simply recovered by incorporating the $\bar{\sigma}^*$ constant into the fitted turnover parameters:

$$\kappa_i \equiv \bar{\sigma}^* k_i^* \quad (36)$$

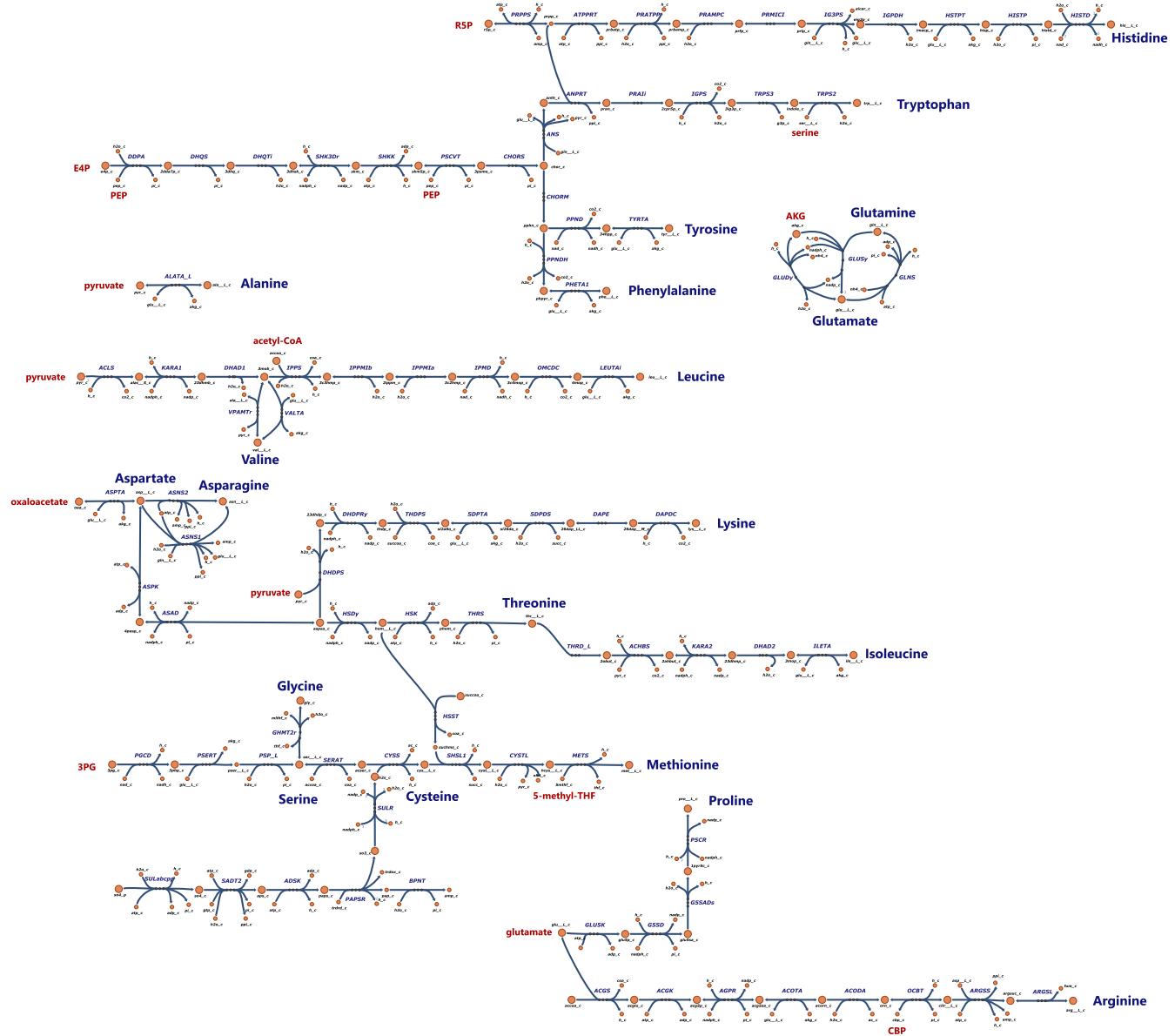
A.5.6 Potential extensions

We conclude this section by noting that the procedure described above assumes that the original parameter set – the set of k_{cat} values used as proxies for apparent k_{cat} values in section A.5.3 – is sufficiently good to produce a realistic flux distribution to use as a reference for the adjustment step. If this is not the case, then a multi-start approach can be implemented, where multiple turnover parameter sets are first generated by perturbing the original parameter set, and each of them is used to generate a separate set of flux distributions. Each reference set is then provided as an input to problem (34), and the solution achieving the lowest objective is chosen in the end. The perturbed parameter set may be generated either randomly (for example, by introducing log-normal noise in the original turnover parameter vector) or systematically. The latter could be achieved, for example, by identifying reactions with zero-predicted flux but high measured abundance of the associated enzyme. By systematically increasing the corresponding turnover number, one can "encourage" these reactions to be included in the reference set, potentially leading to the exploration of more relevant reference distributions than achieved by random perturbations.

Note that, for in this work, we limited ourselves to the original parameter set and thus did not explore this potential heuristic.

Appendix B Supplementary Figures

Biosynthesis of aminoacids

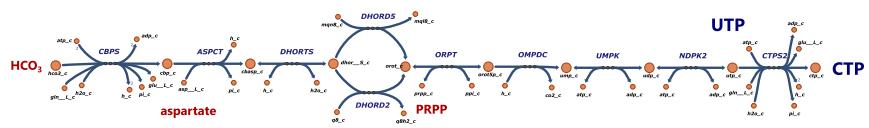


Supplementary Figure S1: Biosynthesis of amino acids (blue labels) from core metabolism precursors (red labels) in iCH360. The Escher maps (including Suppl. Fig. S2 to S4) are available along with the code. R5P: ribose-5-phosphate; E4P: erythrose 4-phosphate; PEP: phosphoenolpyruvate; 3PG: 3-phosphoglycerate; AKG: alpha-ketoglutarate; CBP: carboamyl phosphate.

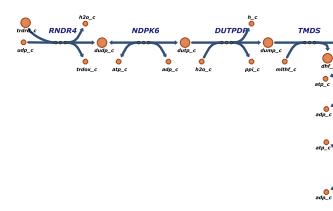
Biosynthesis of nucleotides

Pyrimidine deoxyribonucleotides

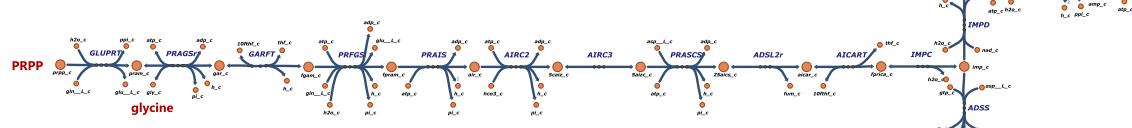
Pyrimidine ribonucleotides



Pyrimidine deoxyribonucleotides



Purine ribonucleotides

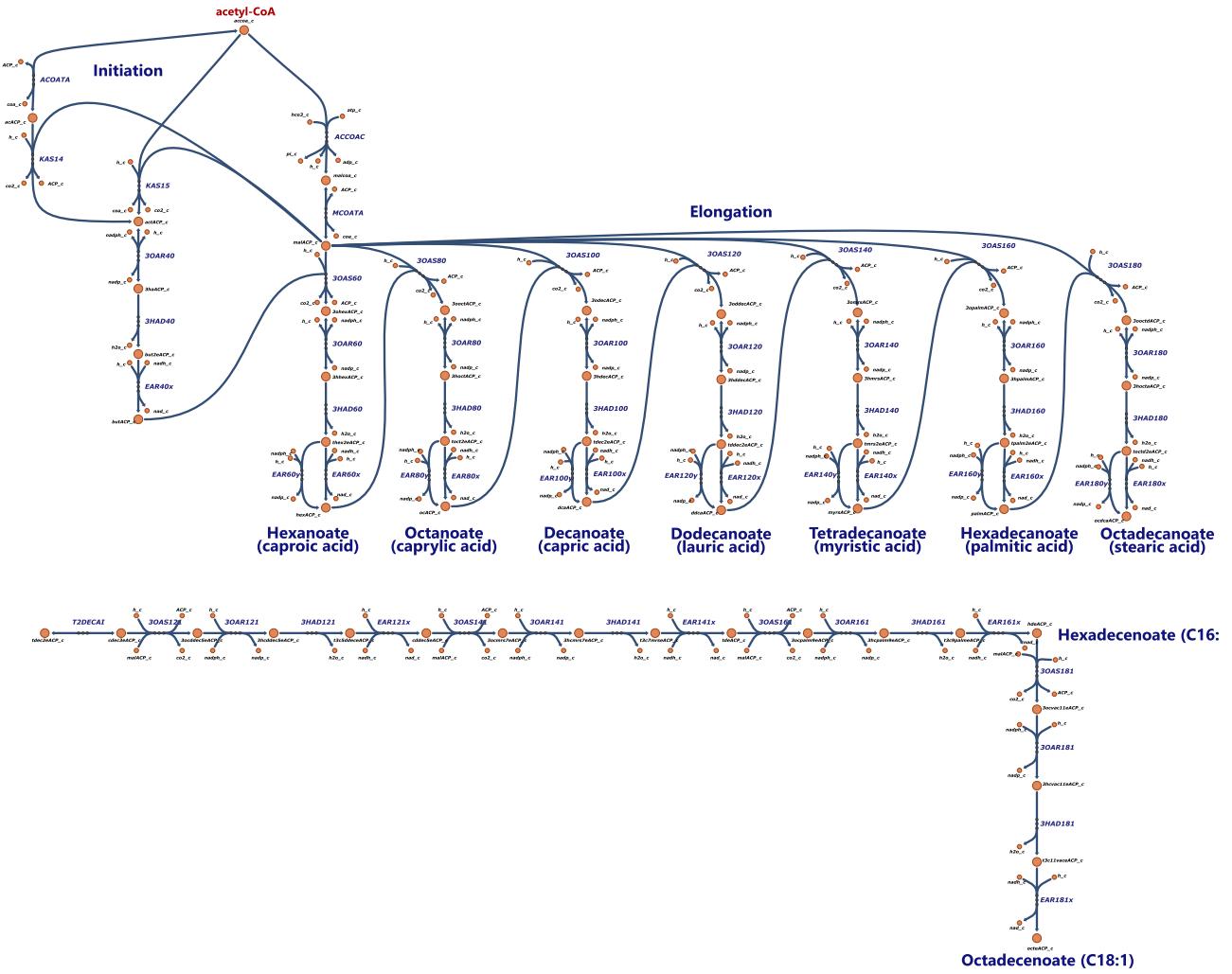


Purine deoxyribonucleotides



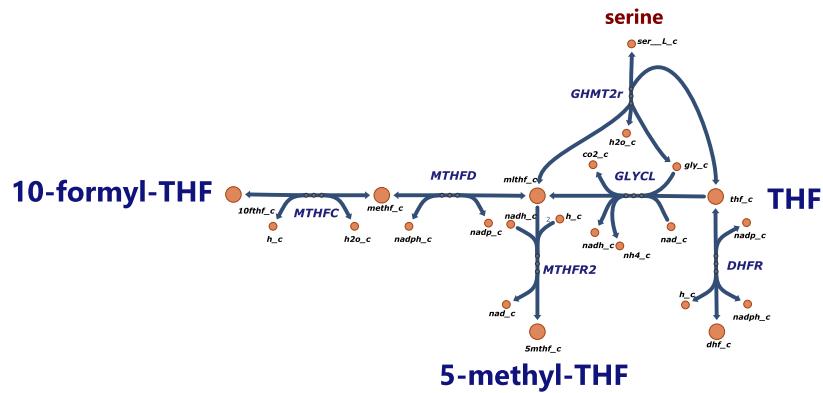
Supplementary Figure S2: Biosynthesis of pyrimidine and purine (deoxy)ribonucleotides (blue labels) from core and aminoacid metabolism precursors (red labels) in iCH360. PRPP: 5-phosphoribosyl-1-pyrophosphate.

Biosynthesis of fatty acids

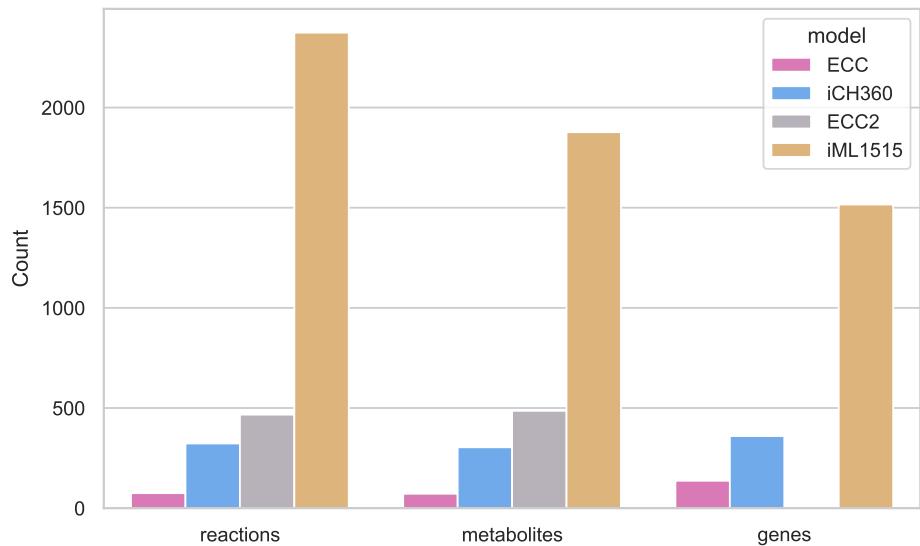


Supplementary Figure S3: Biosynthesis of saturated and unsaturated fatty acids from acetyl-CoA. The map for saturated fatty acids was taken from Escher [18].

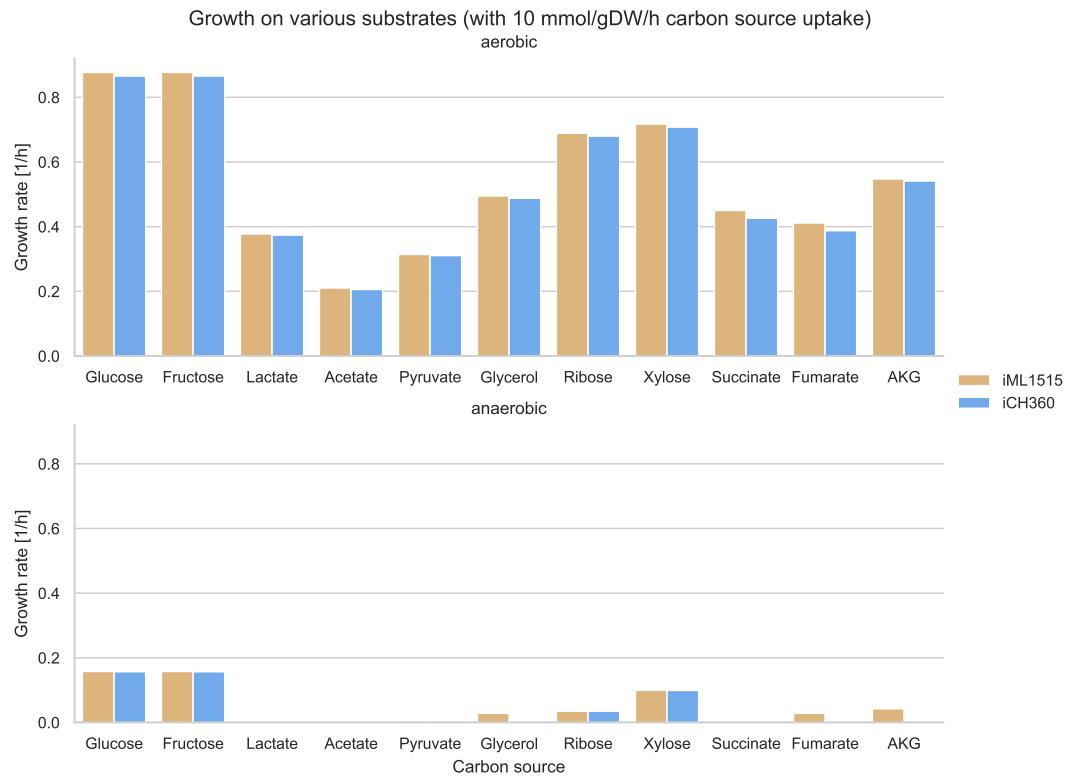
C1 Pool



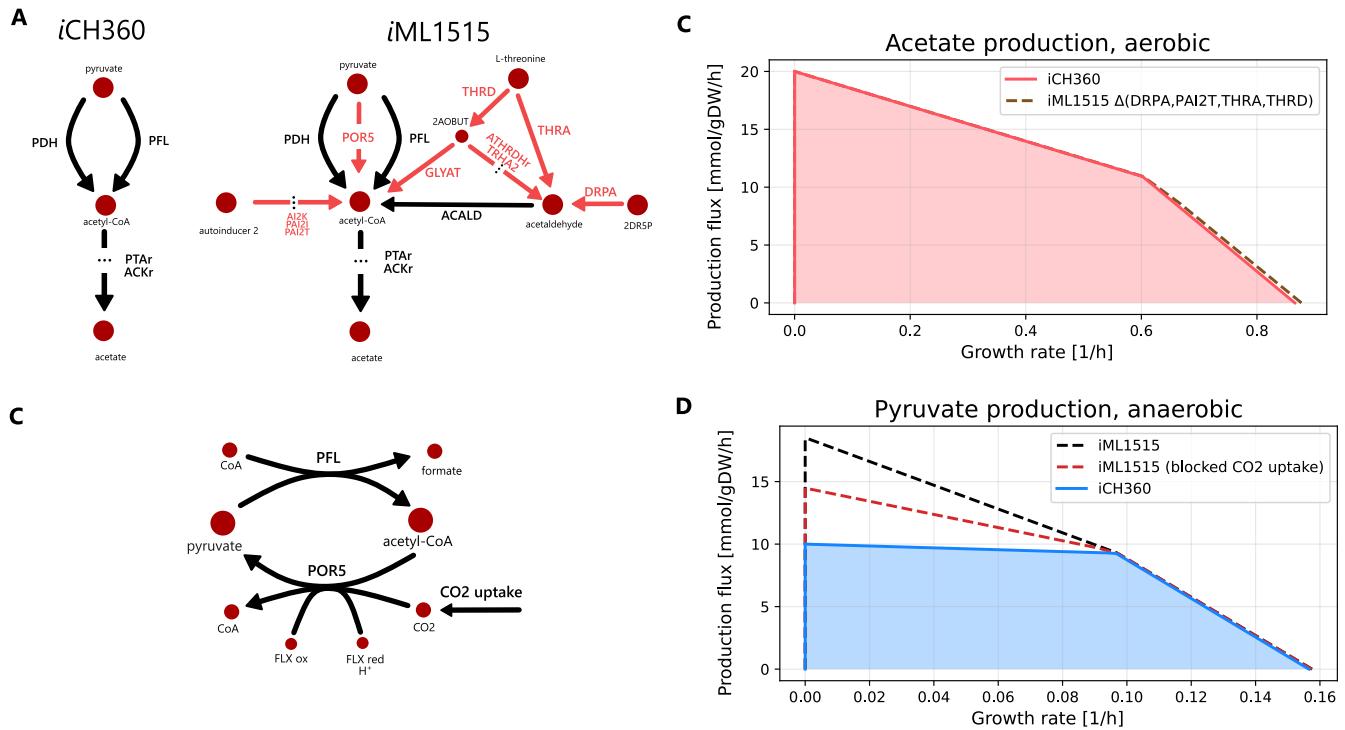
Supplementary Figure S4: Metabolism of one-carbon compounds in *iCH360*. THF: tetrahydrofolate.



Supplementary Figure S5: Comparison of model sizes between ECC, ECC2, *iCH360* and *iML1515*. To allow for a fair comparison, pseudo-reactions (e.g. exchange reactions) were excluded from the count. Note that gene annotations were not available in the SBML model of ECC2 accompanying its publication [14].

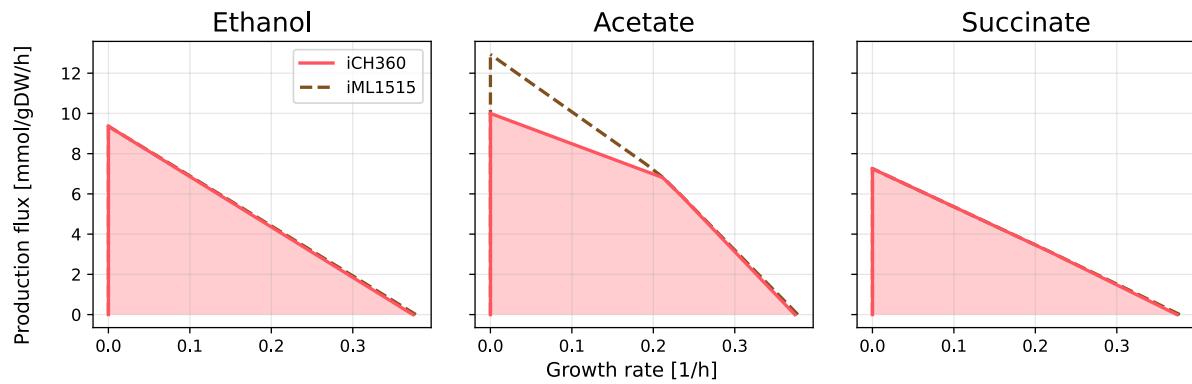


Supplementary Figure S6: Maximal biomass fluxes achieved by *iCH360* and its parent model *iML1515*, for aerobic and anaerobic growth across multiple carbon sources. In all cases, the substrate uptake flux was bounded to 10 mmol/gDW/h.

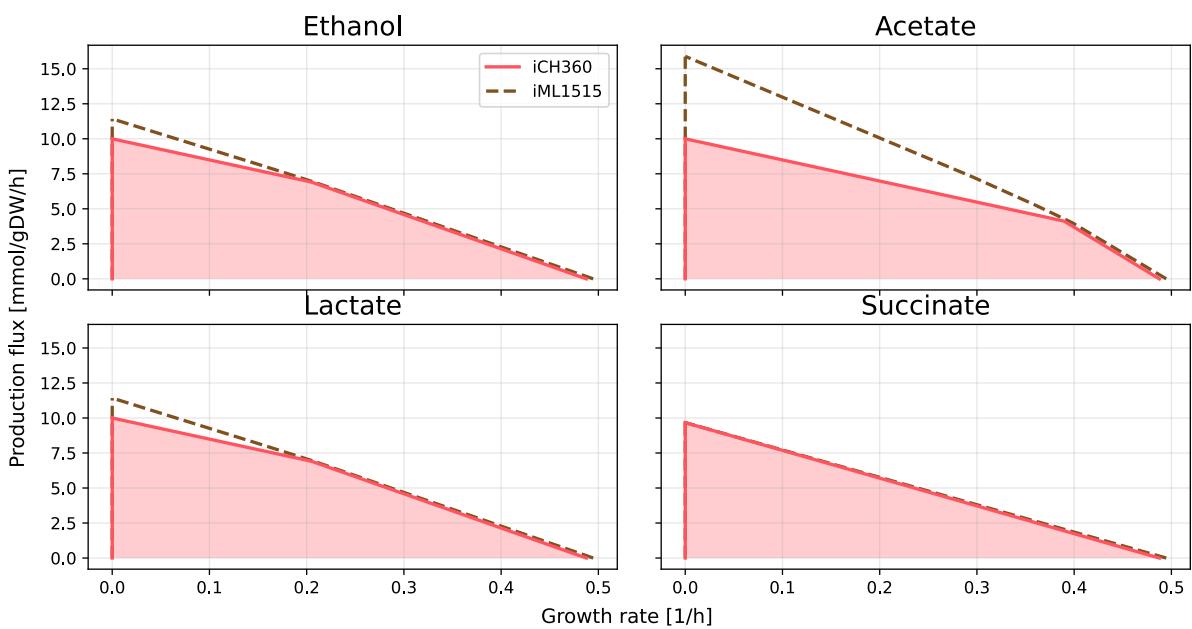


Supplementary Figure S7: Analysis of the differences in acetate production between *iCH360* and its genome-scale parent, *iML1515*. **A:** available metabolic routes for the production of acetate in both models. Left: *iCH360* can only produce acetyl-CoA, precursor for acetate, via the oxidation of pyruvate by either pyruvate dehydrogenase (PDH) or pyruvate-formate-lyase (PFL). Note that the latter reaction is known not to be active under aerobic conditions, but we did not block it for the purpose of this analysis. *iML1515* can additionally produce acetate via additional pathways not present in *iCH360* (in red). Note that only the main substrate and products for each reaction are shown for clarity. 2AOBUT: L-2-Amino-3-oxobutanoate; 2DR5P: 2-Deoxy-D-ribose 5-phosphate. **B:** Blocking these degradation routes by simultaneous knockout of four reactions (DRPA, PAI2T, THRA, THRD) results in the two model sharing a virtually identical production envelope under aerobic conditions. The production envelope shown here was computed for aerobic growth using glucose as a carbon source. **C:** Under anaerobic conditions, the differences between the two model are further exacerbated by the ability of the genome-scale network to achieve higher pyruvate production yield. The genome-scale model can produce higher amounts of pyruvate by uptaking external CO₂ and using it as an electron sink using the POR5 reaction. FLX ox/red: oxidised/reduced flavodoxin. **D:** Blocking CO₂ uptake reduces the differences in pyruvate production between the two models, but does not remove them completely, implying the existence of additional mechanisms used by *iML1515* to achieve higher pyruvate yields. The production envelopes shown here was computed for anaerobic growth using glucose as a carbon source.

Aerobic growth on Lactate

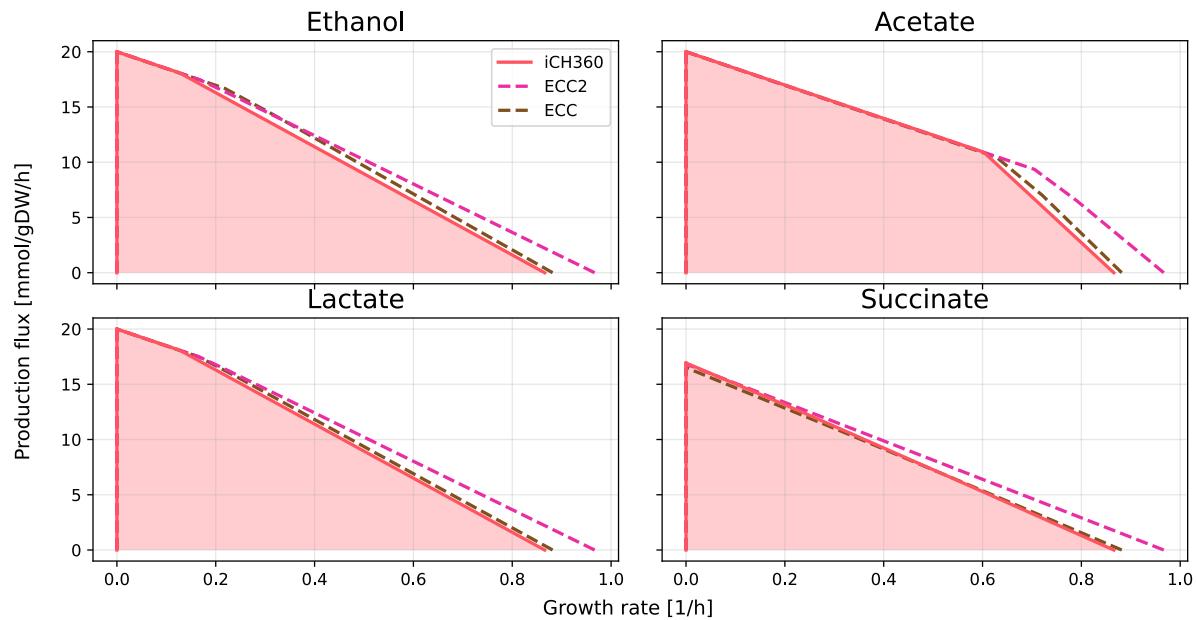


Aerobic growth on glycerol

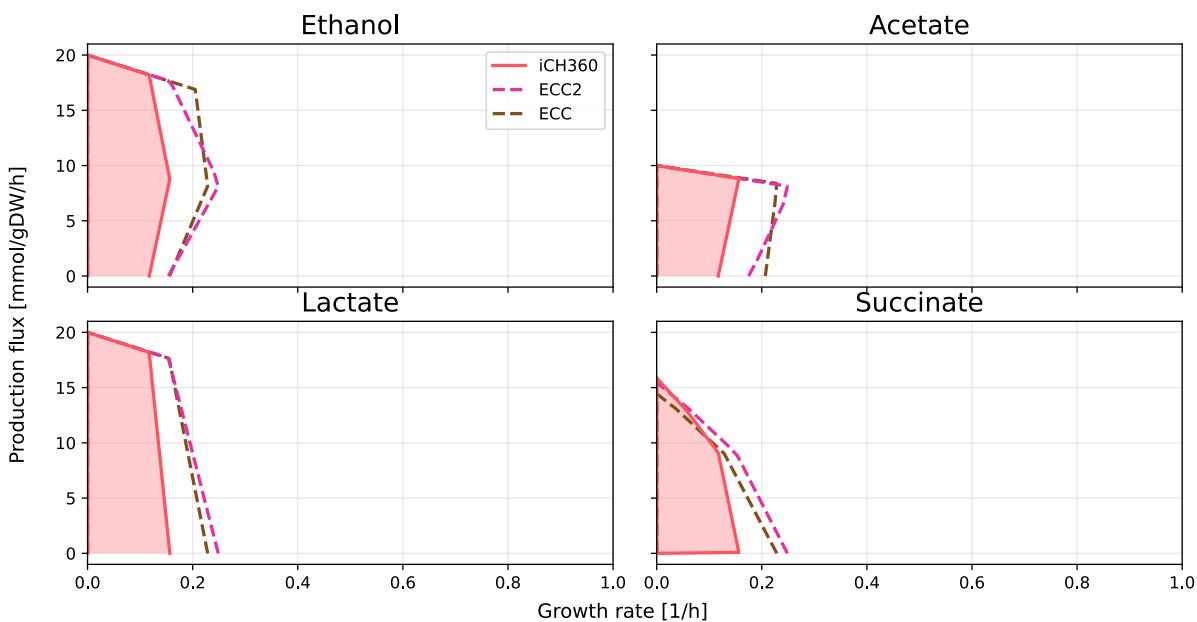


Supplementary Figure S8: Comparison of production envelopes between *iCH360* and its parent model *iML1515*. Top: production of ethanol, acetate and succinate during aerobic growth on lactate. Bottom: production of ethanol, acetate, lactate and succinate during aerobic growth on glycerol. Note that the dashed line representing the production envelope of *iML1515* is sometimes hidden behind the blue line corresponding to *iCH360*.

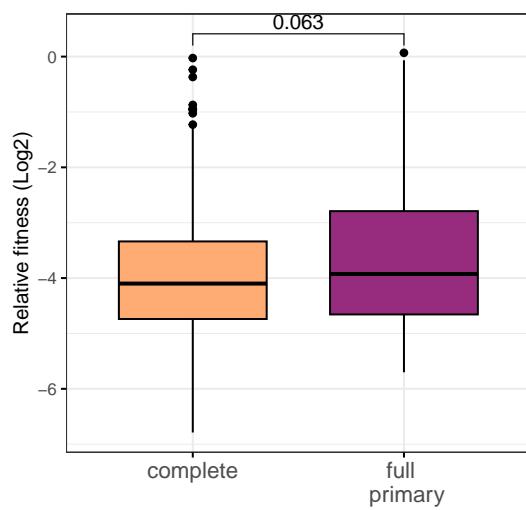
Aerobic growth on glucose



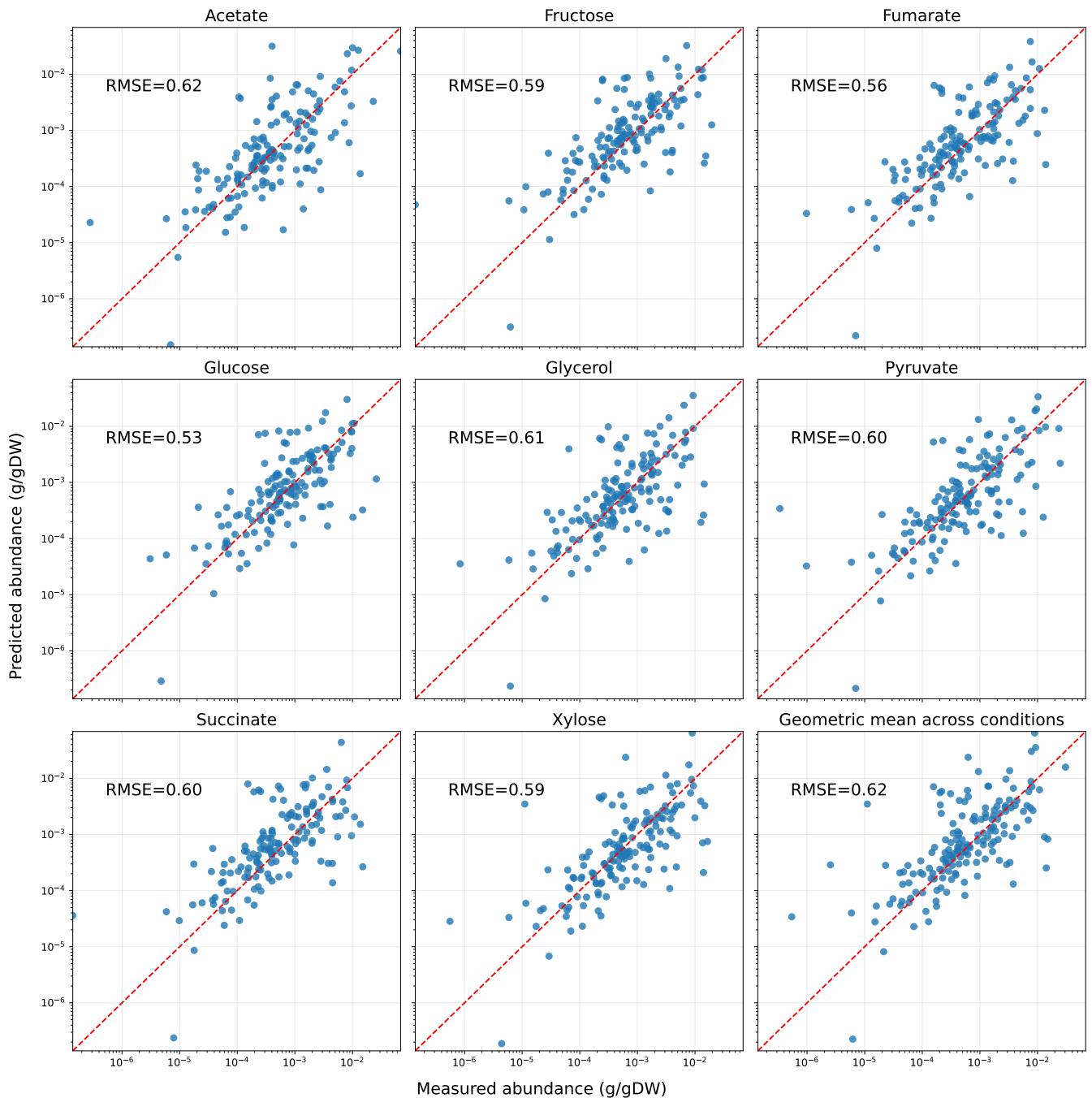
Anaerobic growth on glucose



Supplementary Figure S9: Comparison of production envelopes between iCH360 and other medium-scale models, namely *E. coli* Core (ECC) and *E. coli* Core 2 (ECC2) for growth on glucose as a sole carbon source. Top: production of ethanol, acetate and succinate under aerobic conditions. Bottom: production of ethanol, acetate, lactate and succinate under anaerobic conditions. Additional comparisons between the three models are available in the repository supporting this manuscript.

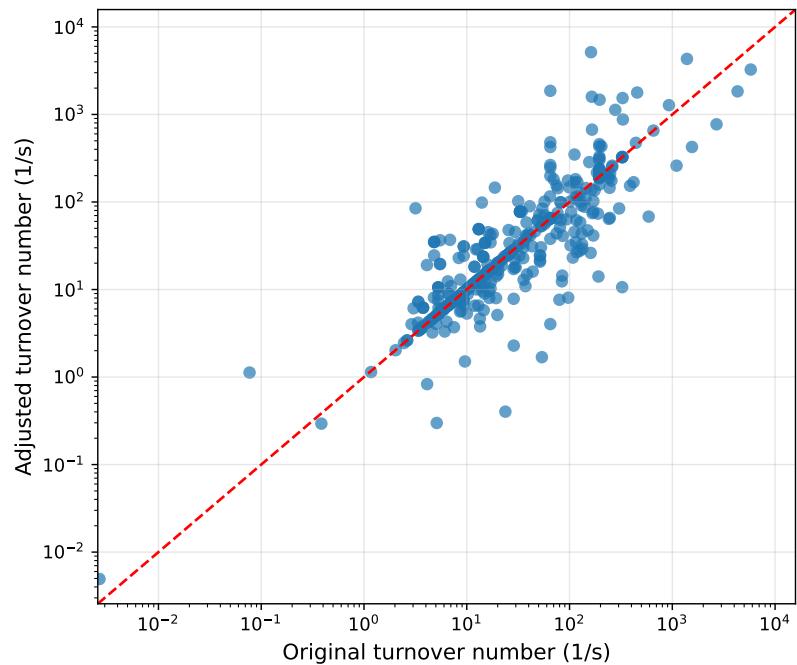


Supplementary Figure S10: Fitness losses associated with disruptions of catalytic edges. There is no significant difference between the fitness effects of disruptions classified as complete disruptions (disruption of all catalytic edges for a reaction) and full primary disruption (disruption of all primary catalytic associations, but with remaining secondary ones), according to a Wilcoxon rank-sum test, $p = 0.063$.

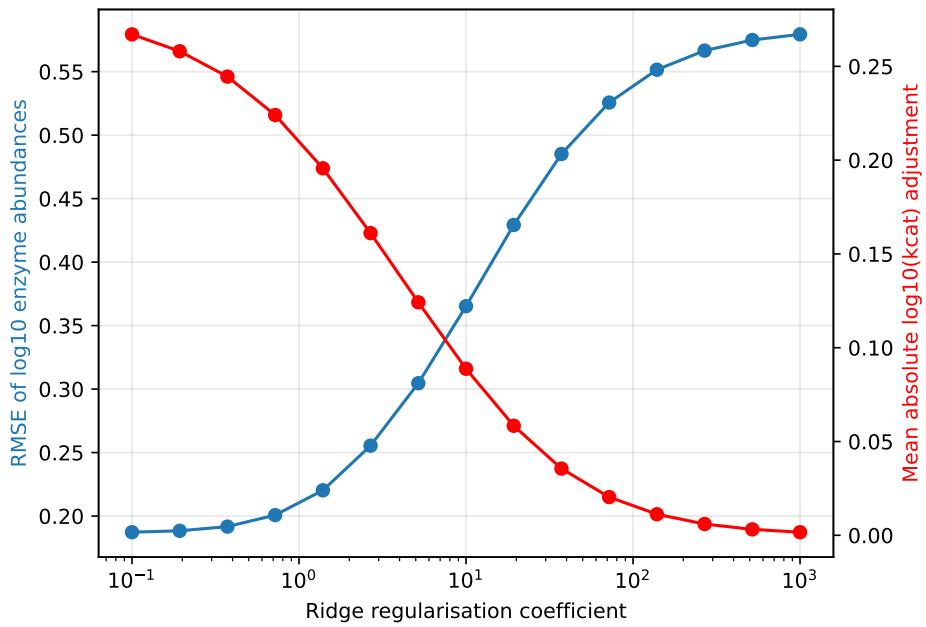


Supplementary Figure S11: Predicted proteome allocation across growth conditions, obtained with the turnover parameter set from Heckmann et al. (2020) [31]. The bottom-right panel shows the geometric means of measurements and predictions across conditions.

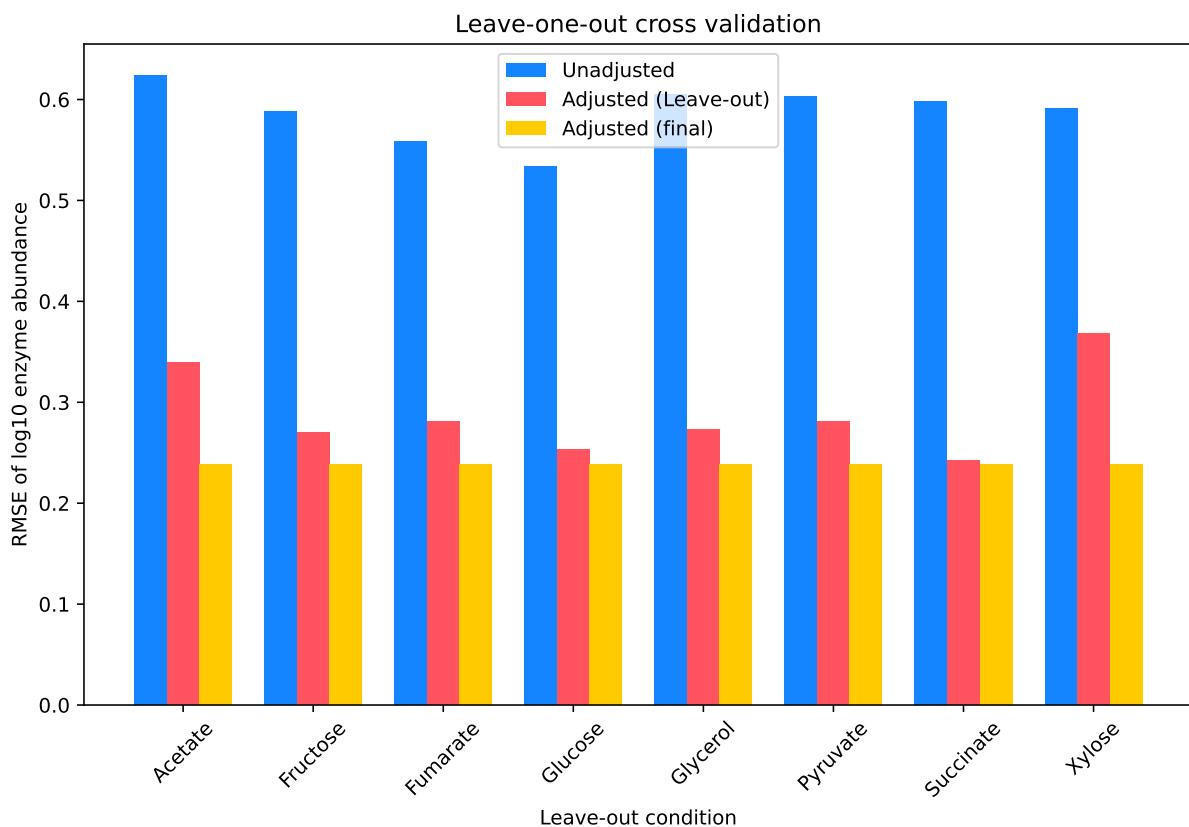
Adjustment of turnover numbers



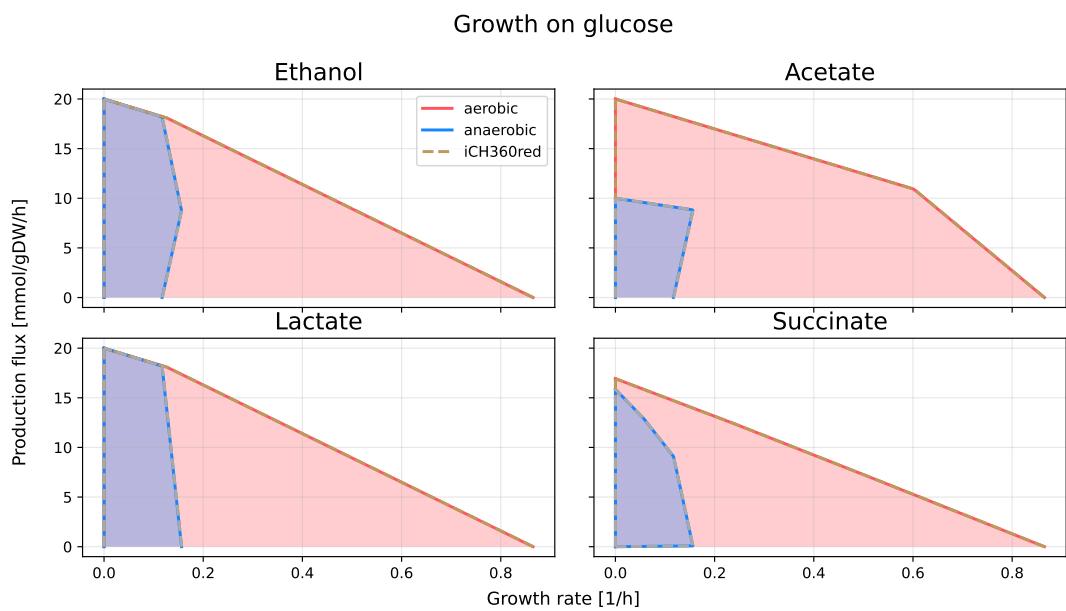
Supplementary Figure S12: Turnover parameters used in EC-*i*CH360 before and after the adjustment procedure (see section 2.4 and Methods in the main text, as well as Supplementary Information A.5 for details). The original parameter set was directly parsed from Heckmann et al. (2020) [31].



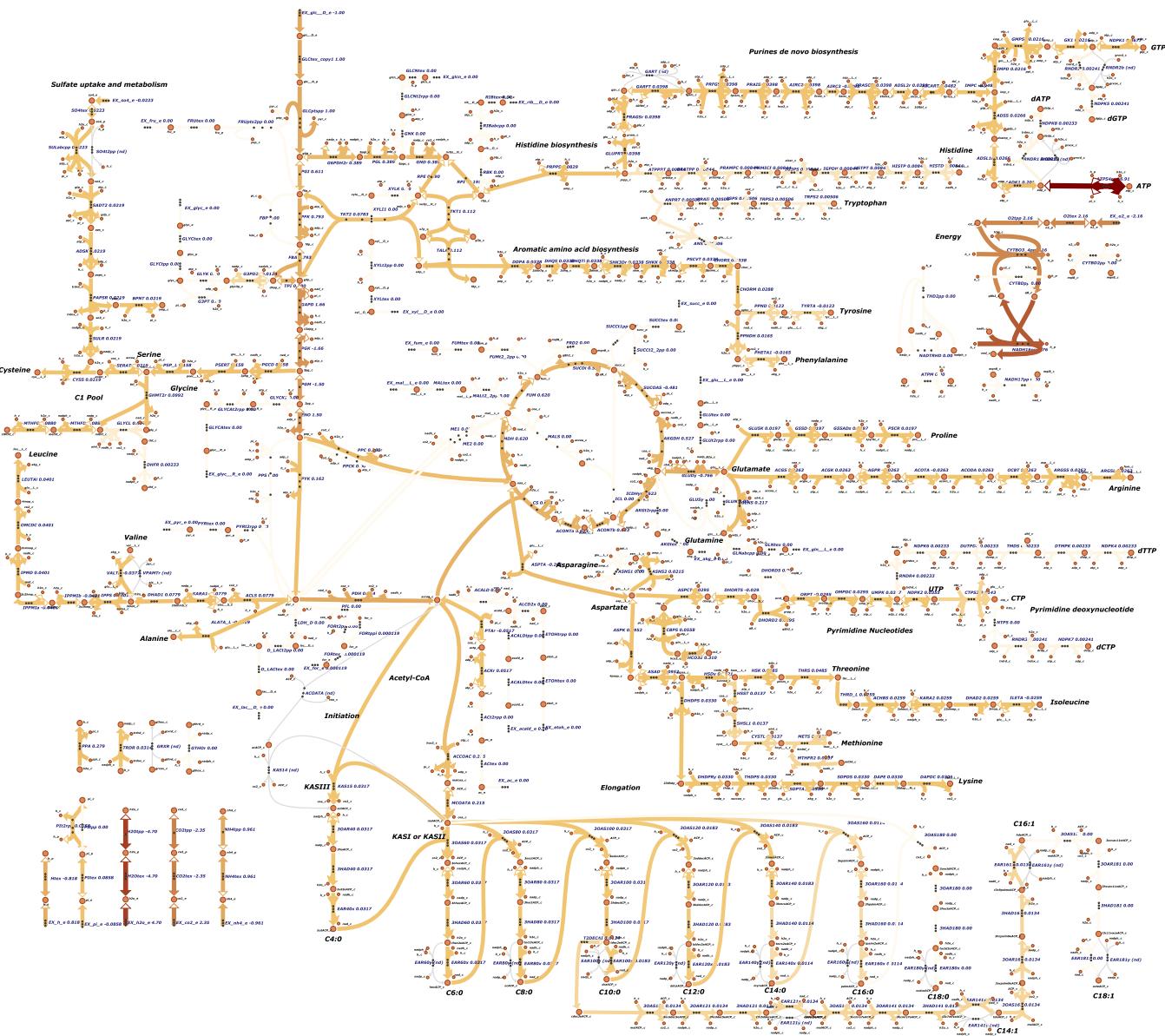
Supplementary Figure S13: Impact of the magnitude of the ridge regularisation coefficient (ρ in Eq. (34)) on the outcome of the adjustment procedure for turnover numbers. The blue curve represents the RMSE between measured and predicted log-enzyme abundances and decreases monotonically as less regularisation is applied to the problem. The red curve represents the mean absolute deviation between original and adjusted log-turnover numbers, which follows the opposite trend and converges to 0 as more regularisation is applied.



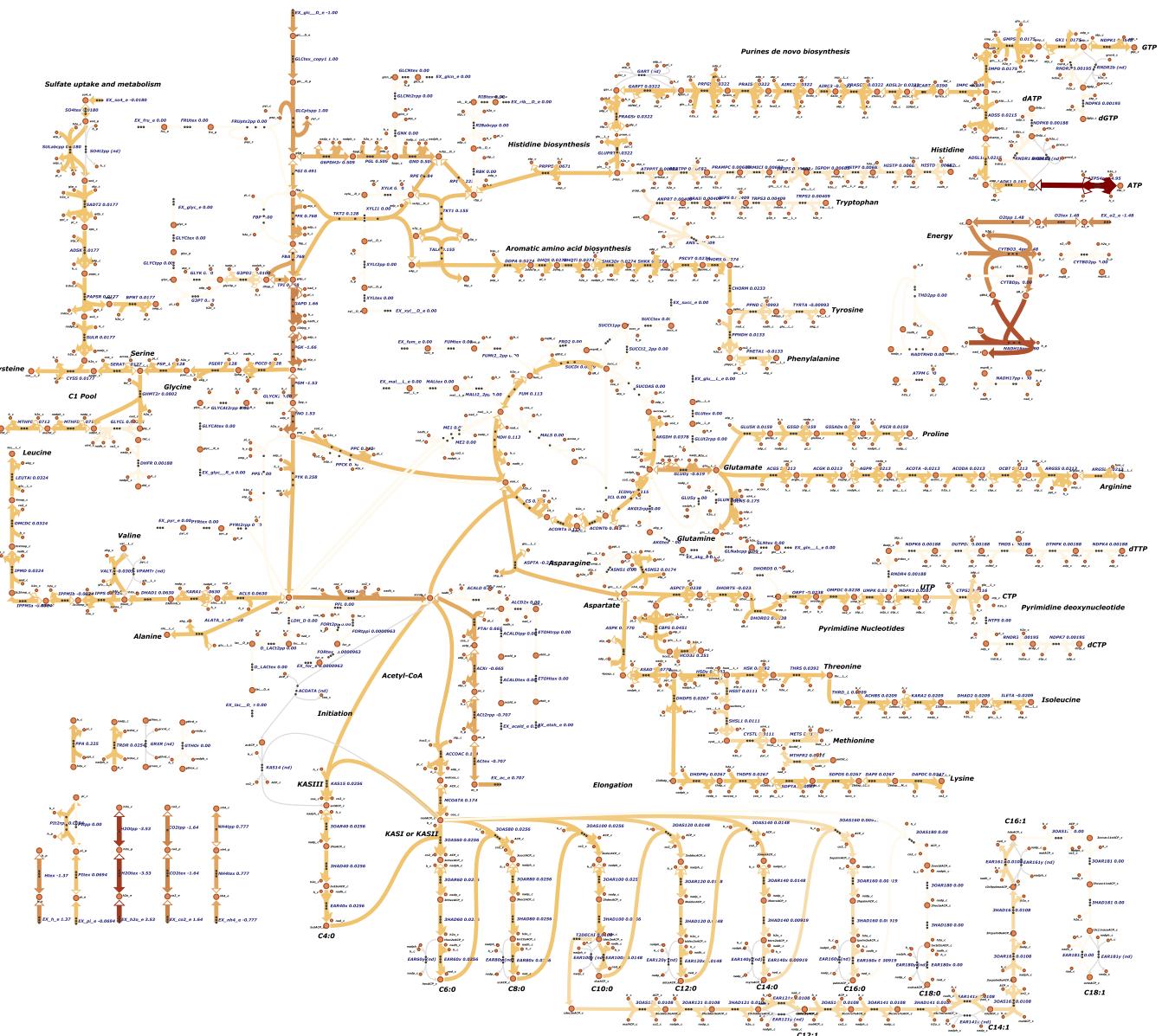
Supplementary Figure S14: Leave-one-out cross-validation for the turnover parameter adjustment procedure. For each condition in the dataset, the graph shows the RMSE (computed for \log_{10} -transformed values) between measurements and predictions of enzyme abundances in that condition. Blue bars show the RMSE computed using the initial, unadjusted parameter set from Heckmann et al. (2020) [31]. Red bars show the RMSE computed after parameter fitting, but excluding the condition for which the RMSE is evaluated from the training dataset. Finally, yellow bars show the RMSE computed using the final adjusted parameter set, obtained including all conditions in the training dataset.



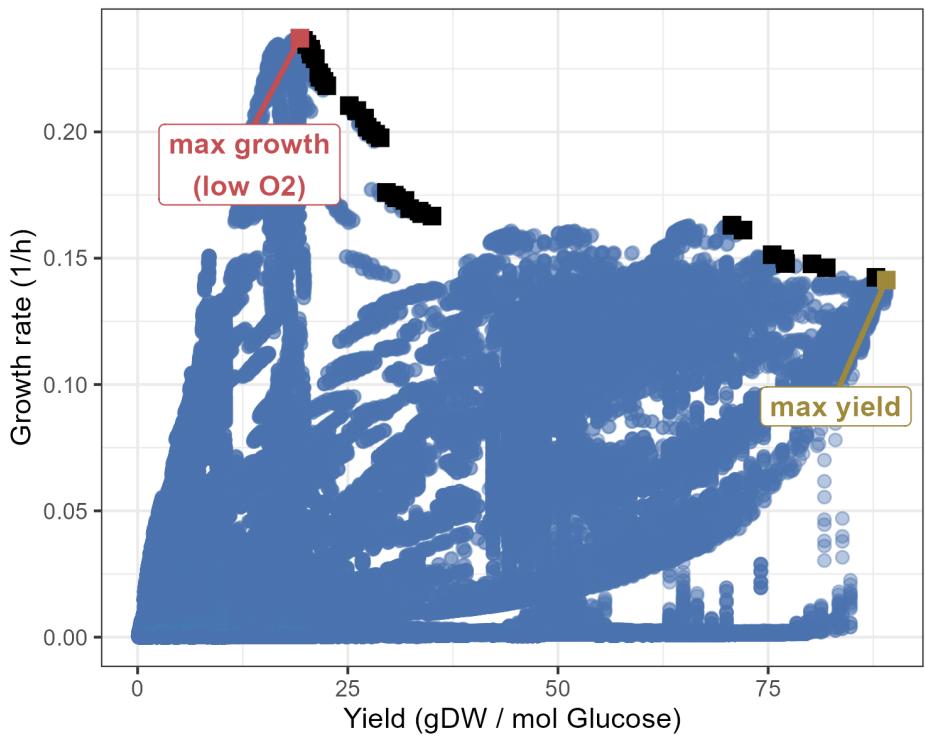
Supplementary Figure S15: Comparison of production envelopes for growth on glucose between *iCH360* and the reduced variant *iCH360_{red}* used for elementary flux mode enumeration and analysis. For the products and growth conditions shown, the two models have virtually identical solution spaces.



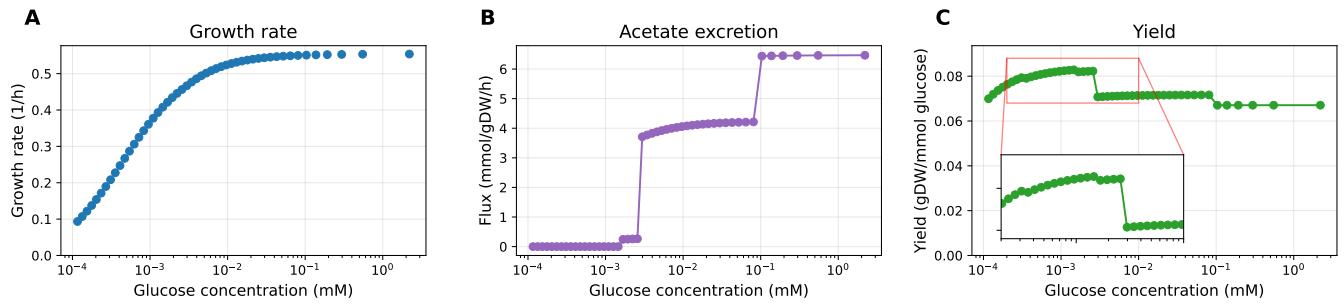
Supplementary Figure S16: Metabolic flux distribution corresponding to the maximum yield elementary flux mode (see main text, section 2.5). The mode is purely respiratory, with no excretion of typical fermentation by-products such as acetate, ethanol, or lactate. The graphics was produced in Escher [18].



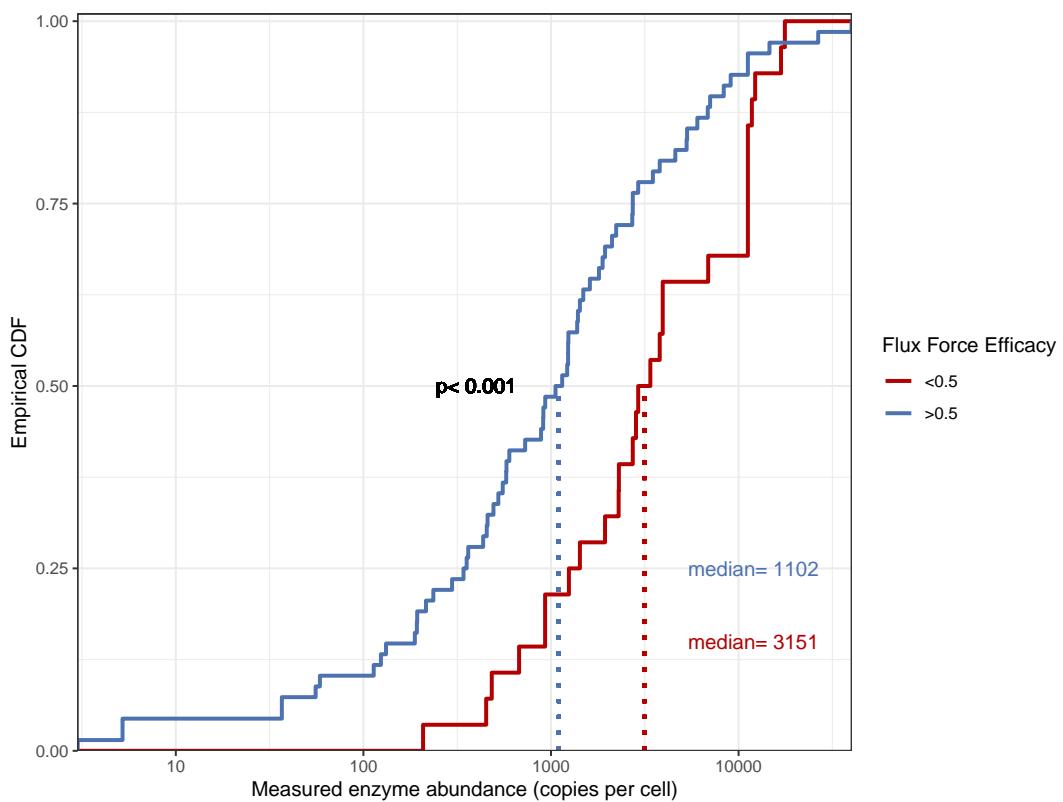
Supplementary Figure S17: Metabolic flux distribution corresponding to the maximum growth elementary flux mode (see main text, section 2.5). The mode shows a mixed respiratory/fermentative phenotype, with significant acetate excretion. The graphics was produced in Escher [18].



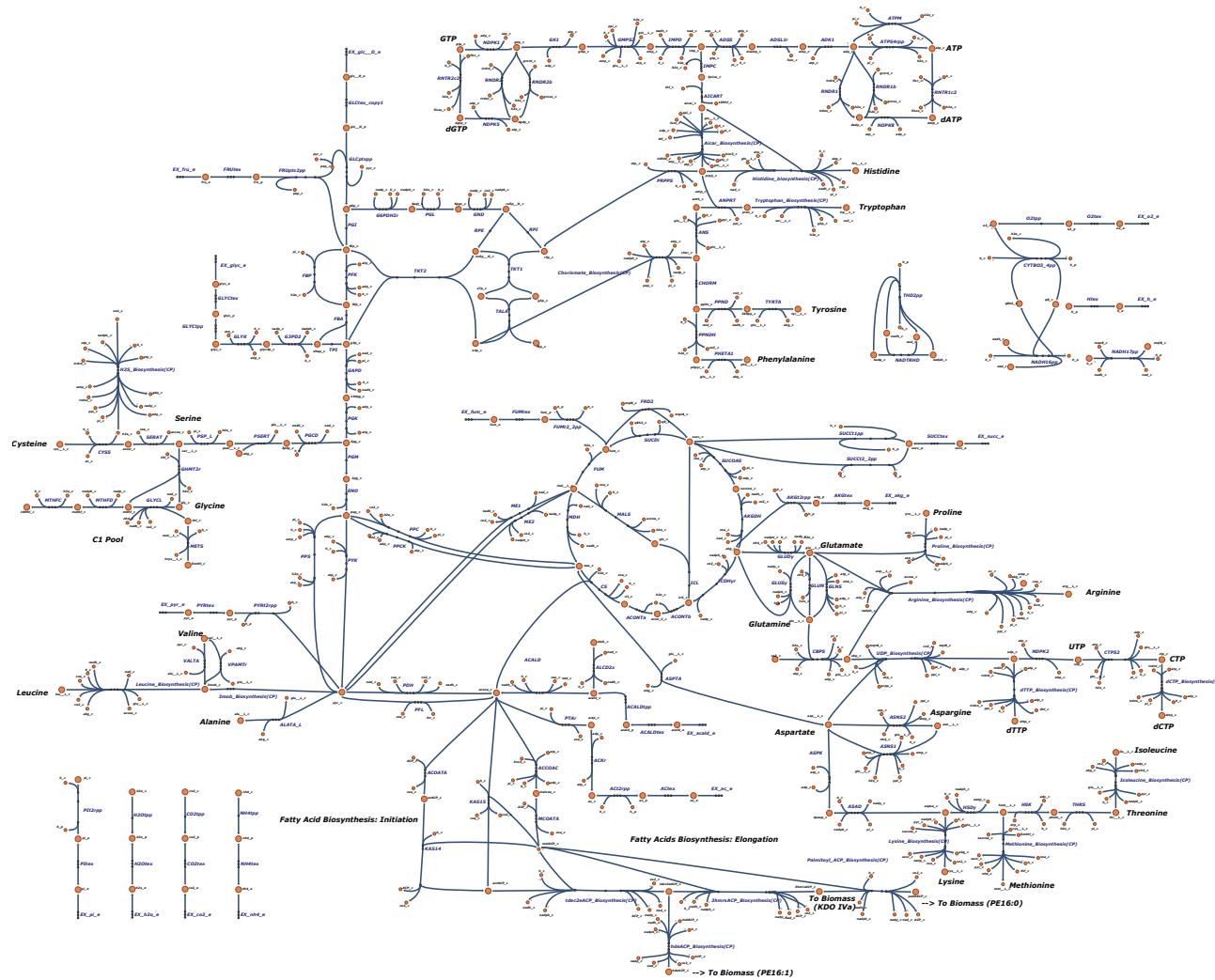
Supplementary Figure S18: EFM growth-yield trade-off front under simulated low oxygen conditions. The figure resembles Figure 5 in the main text, but was obtained by setting a 1000-fold higher cost for the cytochrome reactions. The higher enzyme cost is used to mimic the growth condition in which the oxygen-consuming reaction must operate at very low enzyme saturation, thus requiring higher enzyme-mass investment per unit flux. Black squares denote Pareto optimal EFMs. The Pareto front is much broader than the one obtained with the original enzyme cost for the oxygen-consuming reactions.



Supplementary Figure S19: satFBA results (corresponding to Figure 6 in the main text) obtained by enforcing a lower bound on the ATP maintenance flux equal to 6.86 mmol/gDW/h (taken directly from the parent model *iML1515*). In this case, the optimal solution is no longer an EFM. This is evident from the yield profile (C), which is not piecewise constant with respect to the external glucose concentration.



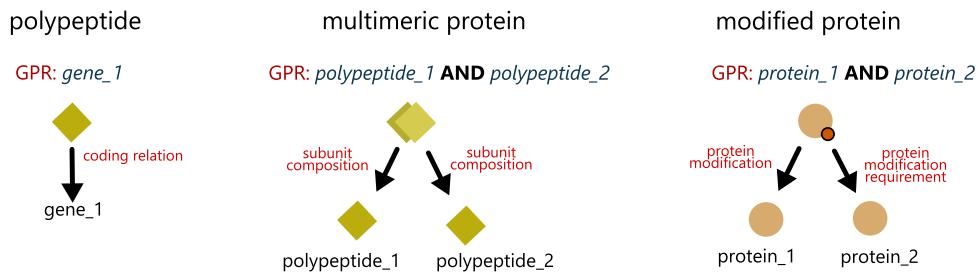
Supplementary Figure S20: Empirical cumulative distribution functions (CDF) of measured enzyme abundances for reactions with predicted flux force efficacies above or below 50% (blue and red curves, respectively). The two functions are significantly different ($p < 0.001$, two-sided Wilcox rank-sum test). The low-efficacy group shows an approximately 3-fold higher median enzyme abundance than the high efficacy group.



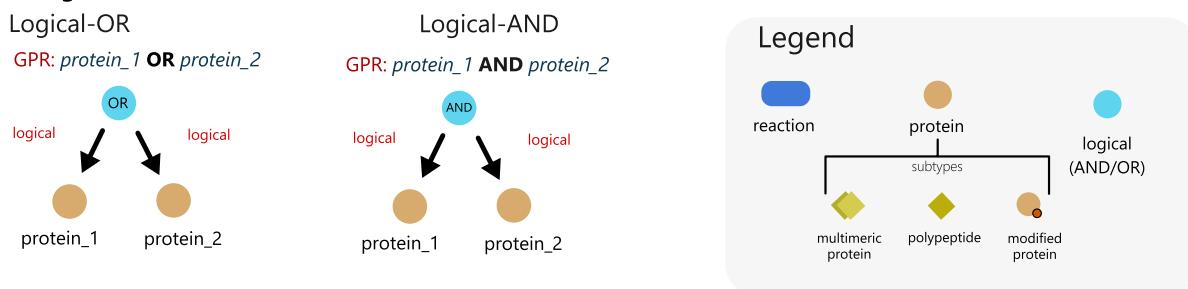
Supplementary Figure S21: Compressed metabolic map of iCH360, wherein linear pathways longer than two reactions were lumped into single effective reactions. The map was produced in Escher [18] and can be used to visualise fluxes, metabolite concentration, and gene expression data.

Computation of boolean rules for different node types in the knowledge graph

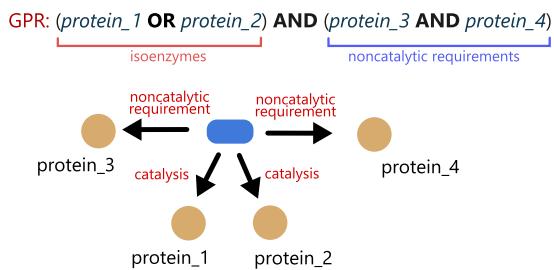
A. Protein nodes



B. Logical nodes



C. Reaction nodes



Supplementary Figure S22: Computation of boolean gene-protein-reaction (GPR) rules with the knowledge graph. For each node in the graph (the top node in each diagram), a boolean expression can be constructed to describe the state of a node (active or inactive) in terms of the state of its child nodes. The exact form of this boolean expression depends both on the node type (e.g. protein, reaction, or logical) and on the type of edges that connect it with its neighbours (e.g. catalysis, subunit composition, etc.). Using these rules, the state of a reaction can be expressed, ultimately, solely in terms of genes in the model, as per convention in standard metabolic models. These GPR rules define a map between a genotype (set of active genes) to a phenotype (set of active reactions, and can be used to simulate *in silico* the effect of given gene knockouts.

Appendix C Supplementary Tables

Supplementary Table 1: Some anecdotal examples of unrealistic predictions obtained when running FBA analyses on large-scale models (here, we considered the genome-scale parent of *iCH360*, *iML1515*). We stress that these behaviours are not the result of errors in the large model. Rather, they result from applying simple methods with few constraints, such as FBA, to a large model with many degrees of freedom. Note that the examples shown here were not obtained through a thorough or systematic investigation of the model prediction abilities.

| Example | Description |
|---|---|
| Production of fatty acids | In the (p)FBA solution computed on <i>iML1515</i> using glucose as a carbon source and growth as an objective, the canonical fatty acid production pathway in <i>E. coli</i> (see Supplementary Figure S3) is completely unused. Instead, the genome-scale model runs beta-oxidation in reverse to produce fatty acids. |
| Anaerobic pyruvate production | Under anaerobic conditions, <i>iML1515</i> can uptake external CO ₂ and use it as a sink for glycolytic electrons, reducing acetyl-CoA produced by PFL back into pyruvate (see Main text, section 2.2). This behaviour, which allows the model to channel additional carbon towards pyruvate, is thermodynamically unrealistic under ambient CO ₂ conditions. |
| Pyruvate auxotrophic strain knock-outs prediction | To construct a pyruvate auxotrophic strain, only a few knock-outs in the central metabolism are necessary (MaeA, MaeB, pcK [63]). However, <i>iML1515</i> uses reactions in amino acid degradations bypassing these knock-outs. |
| Acetyl-CoA auxotrophic strain knock-outs prediction | Similarly, knock-out of 4 genes (aceEF, pflB, poxB), disrupting PDH, PFL and POX reactions, respectively, results in an acetyl-CoA auxotrophic strain that is unable to grow on glucose as the sole carbon source [63]. These knockouts are bypassed <i>iML1515</i> , which can use a number of additional pathways to produce acetyl-CoA (see Main text, section 2.2). |

Supplementary Table 2: Description of node types in the knowledge graph supporting the stoichiometric model. The IDs used in the "Example" column refer to node identifiers in the graph.

| Node type | Node subtype | Description | Example |
|----------------------------|------------------|---|-----------------|
| reaction | | A mass balanced chemical or biochemical reaction | bigg:GAPD |
| protein | polypeptide | A single polypeptide coded by a gene | GAPDH-A-MONOMER |
| | multimer | A complex formed by stoichiometric binding of different polypeptides and/or other complexes | GAPDH-A-CPLX |
| | modified protein | A polypeptide or multimer which underwent post-translational modification | LIPOYL-GCVH |
| gene | | A gene | b1779 |
| metabolite | | An organic or inorganic molecule | L-ASPARTATE |
| logical-AND/ logical-OR | | Used as intermediates node to encode arbitrary logical rules in GPRs generated thorough the graph (e.g. one of two proteins being required by another node) | THIOREDOXINS |

Supplementary Table 3: Classification and properties of edge types (and biological meaning of their associated weights, when applicable) in the iCH360 graph data structure supporting the stoichiometric model.

| Edge type | Parent node type(s) | Child node type(s) | Subtype | Description | Example |
|----------------------------------|----------------------------|---------------------|-------------|---|---|
| catalysis | reaction | protein | primary | The default catalytic relationship between a reaction and an enzyme. | bigg:PFK → 6PFK-1-CPX |
| | | | secondary | A catalytic relationship between a reaction and an enzyme, where the enzyme had been shown in literature (based on <i>in vitro</i> or <i>in vivo</i> evidence) to account for only minor catalytic activity for the reaction when compared to another (primary) isoenzyme. Notes and references to the relevant literature for the secondary annotation are included as edge metadata. | bigg:PFK → 6PFK-2-CPX |
| | | | inactive | Indicates that the child protein is an enzyme for the parent reaction, but it's inactive in the K-12 strain due e.g. to a frameshift mutation. There are only two such edges in the model, and they both involve the same enzyme (Acetohydroxy-acid synthase II) encoded by the <i>ilvG</i> and <i>ilvM</i> genes | bigg:ACHBS → ACETOLACTSYNII-CPLX |
| non-catalytic requirement | reaction | protein | | Indicates that the child protein is required by the parent reaction, although not as a catalyst. Typical examples include proteins used as cofactors (e.g. glutaredoxins) in the reaction or featuring as prosthetic groups for a metabolite involved in the reaction (e.g. Acyl-Carrier-protein) | bigg:ACOATA → ACP-MONOMER |
| subunit composition | protein | protein | requirement | Indicates that the child node is a subunit of the parent node and is required for the correct functioning of the complex. The weight of the edge indicates the stoichiometry of the subunit in the complex. | FABA-CPLX ² → FABA-MONOMER |
| | | | accessory | Indicates that the child protein is an accessory subunit of the parent protein, meaning it can be part of the complex (potentially enhancing or modulating its function), but it's not strictly required for the complex to perform its physiological function. The weight of the edge indicates the stoichiometry of the subunit in the complex. | ATPSYN-CPLX ¹ → EG10106-MONOMER |
| protein modification | protein | protein | | Indicates that the parent protein is obtained by post-translational modification of the child protein. | PYRUVFORMLY-CPLX → PYRUVFORMLY-INACTIVE-CPLX |
| protein modification requirement | protein (modified-protein) | protein | | Indicates that the child protein is required to accomplish the post-translational modification leading to the parent protein. | PYRUVFORMLY-CPLX → PFLACTENZ-MONOMER |
| coding relation | protein | gene | | Indicates that the child gene codes for the parent polypeptide | RIBOKIN-MONOMER → b3752 |
| regulation | protein | metabolite, protein | | Indicates that the child metabolite or protein is a regulator for the enzyme. Information about the regulation mode (activation vs inhibition), the regulation mechanism (competitive vs allosteric) and the regulated reaction (if the enzyme catalyses multiple) is provided as edge metadata whenever available. If present, the weight of the edge denotes the activation/inhibition constant for the interaction as reported in EcoCyc, with units indicated as edge metadata. | SHIKIMATE ^{160.0 μM} → AROE-MONOMER |
| putative association | reaction | protein | | Indicates that a putative association between the reaction and the protein has been proposed in literature. | bigg:PFL → EG11910-MONOMER |
| logical | logical AND/OR | any | | Connects logical operator nodes to downstream nodes. Used to create arbitrary complex logic relations in the graph. | THIOREDOXINS → RED-THIOREDOXIN-MONOMER → RED-THIOREDOXIN2-MONOMER (In this example, the logical edges are used to indicate that the THIOREDOXINS node (of type <i>logical-OR</i>) is active when any of the two child nodes (representing two thioredoxins found in <i>E. coli</i>) are active.) |

Supplementary Table 4: Manual curation of the reaction pruning process used to construct *i*CH360_{red}. Each reaction set represents a set of alternative routes for the production of the same metabolite (but using, for example, different cofactors). For each set, the most physiologically relevant option, based on available literature, was preserved in the reduced model variant.

| Reaction set | Pruned in <i>i</i> CH360red | Notes |
|------------------------------|-----------------------------|--|
| EAR(n)x, EAR(n)y * | EAR(n)y | Enzyme FabI can work with both NADH/NADPH, but higher activity was found with NADH [64] |
| ACOATA, KAS14, KAS15 | ACOATA, KAS14 | Initiation of fatty acid biosynthesis can occur by either direct condensation of acetyl-CoA with malonyl-ACP (KAS15) or by transacylation of acetyl-CoA followed by condensation with malonyl-ACP (ACOATA + KAS15). Because the transacylase activity of FabH (ACOATA) has been to be significantly lower than its condensation activity (KAS15) [65], only the former pathway is maintained in <i>i</i> CH360red. |
| VALTA, VPAMTr | VPAMTr | These are both routes to production of valine. We keep VALTA (<i>ilvE</i>) as it is the last step in the canonical valine biosynthesis route. |
| RNDR1, RNDR2, RNDR1b, RNDR2b | RNDR1b, RNDR2b | The ribonucleoside diphosphate reductase can work both with the thioredoxin and glutaredoxin redox systems. <i>i</i> CH360 retains only the thioredoxin version. |
| SULabcpp, SO4t2pp | SO4t2pp | SULabcpp is an ATP-mediated active transport of sulfate in the cell via an ATP-binding-cassette (ABC) transporter [66], while SO42tpp (<i>cysZ</i>) is a proton symporter [67]. We maintain the former as its impairment was shown to lead to cysteine auxotrophy [66]. |

* $n \in (60, 80, 100, 120, 140, 160, 180, 121, 141, 161, 181)$

Supplementary Table 5: Numbers of elementary flux modes enumerated for the reduced model variant *iCH360red* under different growth conditions. Numbers in brackets represent the numbers of EFM after filtering. Filtered modes include only those supporting Biomass flux and, for aerobic conditions, those that a) have nonzero oxygen uptake and b) do not use either of three reactions (PFL, DHORD5, FRD2), which are known to be only physiologically active under anaerobic conditions

| | number of EFMs (filtered) | |
|----------|---------------------------|------------------|
| | Aerobic | Anaerobic |
| Glucose | 13468719 (1035696)) | 204028 (195670)) |
| Pyruvate | 1763631 (135266)) | 6949 (6480)) |
| Glycerol | 922217 (82112)) | NA |
| Acetate | 38099 (7596) | NA |
| Lactate | 1270315 (5897)) | 1497 (1424)) |