

Toward a Genome Scale Dynamic Model of Cell-Free Protein Synthesis in *Escherichia coli*

Nicholas Horvath, Michael Vilkhovoy, Joseph Wayman, Kara Calhoun¹, James Swartz¹ and Jeffrey D. Varner*

Robert Frederick Smith School of Chemical and Biomolecular Engineering
Cornell University, Ithaca NY 14853

¹School of Chemical Engineering
Stanford University, Stanford, CA 94305

Running Title: Dynamic modeling of cell-free protein synthesis

To be submitted: *Scientific Reports*

*Corresponding author:

Jeffrey D. Varner,

Professor, Robert Frederick Smith School of Chemical and Biomolecular Engineering,
244 Olin Hall, Cornell University, Ithaca NY, 14853

Email: jdv27@cornell.edu

Phone: (607) 255 - 4258

Fax: (607) 255 - 9166

Abstract

Cell-free protein expression systems have become widely used in systems and synthetic biology. In this study, we developed an ensemble of dynamic *E. coli* cell-free protein synthesis (CFPS) models. Model parameters were estimated from measurements of glucose, organic acids, energy species, amino acids, and the protein product, chloramphenicol acetyltransferase (CAT). The ensemble described all of the training data, especially the central carbon metabolism. The model predicted a carbon yield for CAT production that was equal to 23% of the maximum theoretical yield, calculated using sequence-specific flux balance analysis. This suggests that CAT production could be further optimized. The dynamic modeling approach predicted that substrate consumption of glucose and pyruvate and oxidative phosphorylation were most important to both CAT production and the system as a whole, while CAT production alone depended heavily on the CAT synthesis reaction. Conversely, CAT production was robust to allosteric control, as was most of the network, with the exception of the organic acids in central carbon metabolism. This study is the first to model dynamic protein production in *E. coli*, and should provide a foundation for genome-scale, dynamic modeling of cell-free *E. coli* protein synthesis.

Keywords: Biochemical engineering, systems biology, cell-free protein synthesis

1 Introduction

2 Cell-free systems offer many advantages for the study, manipulation and modeling of
3 metabolism compared to *in vivo* processes. Central amongst these, is direct access to
4 metabolites and the biosynthetic machinery without the interference of a cell wall, or com-
5 plications associated with cell growth. This allows us to interrogate the chemical environ-
6 ment while the biosynthetic machinery is operating, potentially at a fine time resolution.
7 Cell-free protein synthesis (CFPS) systems are arguably the most prominent examples
8 of cell-free systems used today [1]. However, CFPS is not new; CFPS in crude *E. coli*
9 extracts has been used since the 1960s to explore fundamentally important biological
10 mechanisms [2, 3]. Today, cell-free systems are used in a variety of applications ranging
11 from therapeutic protein production [4] to synthetic biology [5, 6]. However, if CFPS is to
12 become a mainstream technology for applications such as point of care manufacturing,
13 we must first understand the performance limits of these systems. One tool to address
14 this question is mathematical modeling.

15 Mathematical modeling has long contributed to our understanding of metabolism. Dec-
16 ades before the genomics revolution, mechanistically structured metabolic models arose
17 from the desire to predict microbial phenotypes resulting from changes in intracellular
18 or extracellular states [7]. The single cell *E. coli* models of Shuler and coworkers pio-
19 neered the construction of large-scale, dynamic metabolic models that incorporated multi-
20 ple, regulated catabolic and anabolic pathways constrained by experimentally determined
21 kinetic parameters [8]. Shuler and coworkers generated many single cell kinetic mod-
22 els, including single cell models of eukaryotes [9, 10], minimal cell architectures [11], as
23 well as DNA sequence based whole-cell models of *E. coli* [12]. In the post genomics
24 world, large-scale stoichiometric reconstructions of microbial metabolism popularized by
25 techniques such as flux balance analysis (FBA) have become a standard approach [13].
26 Since the first genome-scale stoichiometric model of *E. coli*, developed by Edwards and

Palsson [14], well over 100 organisms, including industrially important prokaryotes are now available [15–17]. Stoichiometric models rely on a pseudo-steady-state assumption to reduce unidentifiable genome-scale kinetic models to an underdetermined linear algebraic system, which can be solved efficiently even for large systems. Traditionally, stoichiometric models have also neglected explicit descriptions of metabolic regulation and control mechanisms, instead opting to describe the choice of pathways by prescribing an objective function on metabolism. Interestingly, similar to early cybernetic models, the most common metabolic objective function has been the optimization of biomass formation [18], although other metabolic objectives have also been estimated [19]. Recent advances in constraint-based modeling have overcome the early shortcomings of the platform, including capturing metabolic regulation and control [20]. Thus, modern constraint-based approaches have proven extremely useful in the discovery of metabolic engineering strategies and represent the state of the art in metabolic modeling [21, 22]. However, genome-scale kinetic models of industrial important organisms such as *E. coli* have yet to be constructed.

In this study, we developed an ensemble of kinetic cell-free protein synthesis (CFPS) models using dynamic metabolite measurements in an *E. coli* cell free extract. Model parameters were estimated from measurements of glucose, organic acids, energy species, amino acids, and the protein product, chloramphenicol acetyltransferase (CAT). Characteristic values for model parameters and initial conditions, estimated from literature, were used to constrain the parameter estimation problem. The ensemble of parameter sets described the training data with a median cost that was greater than two orders of magnitude smaller than random sets constructed using the literature parameter constraints. We then used the ensemble of kinetic models to analyze the CFPS reaction. First, sensitivity analysis of the dynamic model suggested that CAT production was most sensitive to CAT synthesis parameters, as well as reactions in oxidative phosphorylation and pyruvate con-

sumption. Sensitivity analysis also showed that the system as a whole was most sensitive to these same parts of the network and glucose consumption. CAT production and other metabolites, specifically organic acid intermediates such as pyruvate, were sensitive to the presence of allosteric control mechanisms. Next, to gauge the performance of the cell-free reaction, we compared the observed CAT carbon yield with the maximum theoretical CAT carbon yield calculated using sequence-specific flux balance analysis. The CAT yield estimated from the kinetic model was 23% of the maximum theoretical yield, but 36% of the theoretical yield when physiologically realistic constraints were used. Taken together, we have integrated traditional kinetics with a logical rule-based description of allosteric control to simulate a comprehensive CFPS dataset. This study provides a foundation for genome-scale, dynamic modeling of cell-free *E. coli* protein synthesis.

Results

The ensemble of kinetic CFPS models captured the time evolution of CAT biosynthesis (Fig. 1 - 3). The cell-free *E. coli* metabolic network was constructed by removing growth associated reactions from the MG1655 reconstruction [23], and by adding reactions describing chloramphenicol acetyltransferase (CAT) biosynthesis, a model protein for which we have a comprehensive training dataset [24]. The CFPS model equations were formulated using the hybrid cell-free modeling framework of Wayman et al. [25]. An ensemble of model parameters ($N > 10,000$) was estimated from measurements of glucose, CAT, organic acids (pyruvate, lactate, acetate, succinate, malate), energy species (A(x)P, G(x)P, C(x)P, U(x)P), and 18 of the 20 proteinogenic amino acids using a constrained Markov Chain Monte Carlo (MCMC) approach. The MCMC algorithm minimized the error between the training data and model simulations starting from an initial parameter set assembled from literature and inspection. Parameter sets were selected for the ensemble based upon their error, and the Pearson correlation coefficient between the candidate and existing sets in the ensemble. The parameter set with the lowest error value was defined as the best-fit set. Central carbon metabolism (Fig. 1, top), energy species (Fig. 2), and amino acids (Fig. 3) were captured by the ensemble and the best-fit set. The constrained MCMC approach estimated parameter sets with a median error greater than two-order of magnitude less than random parameter sets generated within the same parameter bounds (Fig. 4); thus, we have confidence in the predictive capability of the estimated parameters. The model captures the biphasic nature of CFPS shown by the experimental data: in the first hour glucose consumption powers the system, and CAT is produced at $\sim 10 \mu\text{M/h}$; subsequently, pyruvate and lactate reserves are consumed to power metabolism, and CAT is produced less efficiently at $\sim 5 \mu\text{M/h}$. Allosteric control was important to this biphasic operation, as without it, both phases are sped up, and after 1.5 hours CAT synthesis has ceased (Fig. 1, bottom). In addition, acetate accumulation also

stops after 1.5 hours since there is no more substrate to power the system. Interestingly, malate captures the experimental measurements during the glucose consumption phase, but increases sharply during the pyruvate consumption phase.

While the model ensemble described the experimental data, it was unclear whether the performance of CFPS was optimal. To address this question, we used sequence-specific flux balance analysis in combination with a T7 promoter model to estimate the theoretical maximum performance of the system. We first validated the ssFBA approach by comparing simulated and measured concentrations of CAT for the first hour under glucose consumption (Fig. 7A). The ssFBA ensemble was able to predict the CAT dataset despite having no adjustable parameters. Uncertainty in experimental factors such as RNA polymerase, ribosome concentrations, elongation rates, or the upper bounds for oxygen and glucose consumption rates did not alter the qualitative performance of the model. Thus, the metabolic network and molecular description of transcription and translation were consistent with experimental measurements. To gauge the performance of CFPS, we next calculated the CAT carbon yield for three classes of constraints: (i) theoretical maximum glucose, amino acid and oxygen upper bounds, and no transcriptional/translational constraints; (ii) theoretical maximum glucose, amino acid and oxygen upper bounds, and realistic transcriptional/translational constraints; and (iii) metabolite fluxes constrained by the data, and realistic transcriptional/translational constraints (Fig. 7B). The unconstrained theoretical maximum CAT carbon yield (i) was 0.363 ± 0.02 (Fig. 7B, left); On the other hand, for realistic constraints on transcription and translation (ii), the CAT carbon yield was 0.226 ± 0.03 (Fig. 7B, middle). Lastly, when using realistic metabolite and transcription and translation constraints (iii), the predicted carbon yield was 0.062 ± 0.02 . By comparison, the best-fit parameter set had a CAT carbon yield of 0.086 ± 0.004 , equivalent to 23% of the theoretical maximum (i) and 36% of the physiological case (ii). The experimental dataset had a CAT carbon yield of 0.0821, similar to both the kinetic model and

the experimentally constrained case (iii).

To better understand which parameters and parameter combinations influenced model performance we performed sensitivity analysis (Fig. 5). According to the sensitivity results, CAT production was most dependent on coenzymes and pyruvate. CAT production was most sensitive to the CAT synthesis reaction, oxidative phosphorylation, and the pyruvate-consuming alanine synthesis reaction (Fig. 5, top, section A). Taking into account these three reactions as well as the next 16 most important reactions (section B), we see a common theme of reactions that involve the coenzymes ATP, GTP, NADH, and NADPH as well as the metabolites pyruvate and glutamate. Glutamate is important as a reagent for the synthesis of all other amino acids which are required to assemble the CAT protein. Meanwhile, the coenzymes are needed to provide energy to power the synthesis of CAT, and pyruvate becomes important for energy supply after glucose is exhausted. In addition, pyruvate is a key metabolite required for the synthesis of several amino acids. The pairwise sensitivities (off-diagonal elements) were different from the corresponding first-order sensitivities (diagonal elements), and led to interesting outcomes. The combination of certain reactions had a greater effect on CAT production than that of the reactions by themselves. For example, glutamine synthesis and arginine degradation were both among the most important reactions to CAT production (they rank 5th and 10th, respectively). This was likely because they both affect the sensitive glutamine-glutamate balance; glutamine synthesis consumes glutamate, while arginine degradation produces it. However, when both were perturbed, their combined effect on the model was low, as the respective contributions to consumption and production of glutamate cancelled.

The system state was also most dependent on coenzymes and substrates; however, instead of pyruvate and glutamate, the substrates driving metabolism were glucose and pyruvate. The metabolism as a whole was most sensitive to glucose uptake via GTP and the forward reaction of lactate dehydrogenase, consuming pyruvate (Fig. 5, bottom, sec-

tion F). These two and the next 30 most important reactions (section G) largely involved
coenzymes, especially ATP and NADPH, as well as substrate-consuming reactions and
oxidative phosphorylation. The system state had even more pairwise sensitivities that
differed from the corresponding first-order sensitivities and stood out as significant. For
example, the first-order effect of alanine synthesis was large; it consumes both pyruvate
and glutamate, two key species in the network. However, there were enough reactions
that, when paired with alanine synthesis, had little effect on the model; malic enzyme is
one of these, as it produces the pyruvate that alanine synthesis consumes. Thus, the
total-order alanine synthesis sensitivity was low, placing it at the very bottom of section I.
Another interesting result was the intersection of sections F and G with section J. The 53
reactions in section J were turned off in the best-fit set ($V^{max} = 0$); therefore, the pertur-
bation of these reactions had no effect on the model. As a result, all pairwise sensitivities
with reactions in section J were pseudo first-order sensitivities for the other reactions.
Interestingly, many reactions in section F and several in section G showed their highest
sensitivities when paired with the "non-effects" of section J. Of these, three involved pyru-
vate, strengthening its role as a key metabolite; the others were glucose consumption via
GTP/CTP-specific hexokinases, fumarate reductase, and SO_4 utilization. This suggested
that these reactions' effects on the model were canceled out or lessened by most other
reactions, but were of course not affected by the reactions in section J. This was also
likely the reason that reactions in section J rank above those in section K, despite having
no effect on the model themselves. Taken together, sensitivity analysis identified blocks
of parameters that either individually, or in combination influenced model performance.

The sensitivity analysis showed oxidative phosphorylation was significant for CAT pro-
duction as well as the system state. To further investigate this, we knocked out key re-
actions in oxidative phosphorylation to examine its effect on glucose uptake and CAT
production (Fig. 6). A single knockout of the reaction *cyd* was detrimental to model per-

168 formance and CAT production, reducing the CAT carbon yield from 8.6% to 2.8%. In
169 addition, the glucose uptake rate was reduced compared to that of the control (no knock-
170 outs). A knockout of *nuo* showed a less drastic effect, reducing the CAT carbon yield to
171 6.8%; however, the glucose uptake rate remained similar to that of the control. Knocking
172 out *app* showed a CAT yield to 8.8%, but was not statistically different from that of the con-
173 trol. With all three reactions knocked out, CAT yield was 2.7%, not statistically different
174 from the *cyd* knockout. Thus, the model suggests that the majority of the energetic needs
175 of the system as well as the production of CAT are met by oxidative phosphorylation,
176 specifically by *cyd*.

177 To investigate the difference in yields between the unconstrained and constrained
178 cases, we compared the flux distributions from the ssFBA simulations (Fig. 8). The con-
179 strained cases (ii & iii) heavily utilized the first step in the pentose phosphate pathway
180 to generate NADPH; the carbon flux continued through the Entner–Doudoroff pathway
181 toward pyruvate, a key metabolite as shown by sensitivity analysis. For case ii, the ma-
182 jority of the flux proceeded toward acetate accumulation, whereas in case iii, the flux
183 accumulated as pyruvate, lactate, and acetate with some going through the TCA cycle.
184 In comparison, the unconstrained case (i) showed the majority of flux traveling through
185 glycolysis towards pyruvate, leading to an accumulation of lactate, acetate and malate.
186 In all cases the energy source was primarily oxidative phosphorylation, and to a lesser
187 extent the TCA cycle. However, the accumulation of acetate and lactate signifies that the
188 system is not operating at its highest efficiency. The system produces NADH through
189 lactate dehydrogenase as well as through pyridine nucleotide transhydrogenase (*pntAB*)
190 to power oxidative phosphorylation. Oxidative phosphorylation leads to a high redox ra-
191 tio contributing to the accumulation of acetate overflow and diverting flux away from the
192 TCA cycle. This suggests that there is potential to increase CAT production by reducing
193 the accumulation of acetate and lactate. To investigate this further, we simulated poten-

194 tial knockouts with constrained transcription/translation rates, but constrained the specific
195 glucose and amino acid uptake rates to the same values as simulated with no knock-
196 outs. Knocking out the *gnd* reaction decreased flux of acetate production but increased
197 flux through *pntAB*, which is responsible for regenerating NADPH. The simulation showed
198 carbon was diverted toward lactate; however, since CAT production is constrained by the
199 translation rate, we expected no increase in CAT production. A second simulation with a
200 knockout of *gnd* and phosphate acetyltransferase showed carbon being diverted toward
201 lactate and succinate; however, it required a higher flux through oxidative phosphorylation
202 and the TCA cycle to meet the energetic needs of the system. The decrease in acetate
203 production and the diversion of carbon flux shows a promising mechanism to increase
204 CAT yield.

Discussion

In this study we present an ensemble of *E. coli* cell-free protein synthesis (CFPS) models that accurately predict a comprehensive CFPS dataset of glucose, CAT, central carbon metabolites, energy species, and amino acid measurements. We used the hybrid cell-free modeling approach of Wayman and coworkers, which integrates traditional kinetic modeling with a logic-based description of allosteric regulation. CFPS is seen to be biphasic relying on glucose during the first hour and pyruvate and lactate afterward. Allosteric control was essential to the maintenance of the network and production of CAT, as without it, central carbon metabolism is exhausted within 1.5 hours leading to low CAT production. Having captured the experimental data, we investigated if CAT yield and CFPS performance could be further improved. We showed that the model produces CAT at 23% of the theoretical maximum in terms of carbon yield, and at 36% of a physiological case in which transcription and translation are constrained. The accumulation of waste byproducts, especially acetate, is responsible for this sub-optimal yield. Sensitivity analysis showed that certain substrates and energy species are instrumental to CAT production and overall metabolism. The system heavily relied on oxidative phosphorylation for the system's energetic needs as well as for CAT synthesis. A single knockout in oxidative phosphorylation reduced the CAT carbon yield ~3-fold, as well as disrupting the system state showing its crucial role in CFPS. In comparing flux distributions between low and high yield cases, carbon flux could be potentially diverted toward CAT by reducing acetate overflow and minimizing flux through the Entner-Doudoroff pathway. Taken together, these findings represent the first dynamic model of *E. coli* cell-free protein synthesis, and an important step toward a functional genome scale description.

We present an ensemble of models that quantitatively describes the system behavior of cell-free metabolism and production of CAT. Experimental observations of the metabolites and cometabolites validate the structure of the model and the estimation of kinetic

parameters. This is important in applying metabolic engineering principles to rationally design cell-free production processes and predict the redirection of carbon fluxes to product forming pathways. In analyzing the model parameters' effect on CAT production, CAT synthesis is the most important, followed by oxidative phosphorylation and the glutamate and pyruvate consuming reactions, as well as coenzyme reactions which are necessary to drive CAT synthesis. For example, the conversion of ATP to GTP shows significance since it is necessary for CAT synthesis. While Jewett and coworkers have shown that ATP may be at saturation in CFPS [26], GTP is also required for CAT synthesis and may be a limiting reactant. Thus, supplementation with additional GTP may improve the efficiency of CAT production. A similar theme is seen in the sensitivity of overall model state, where the most important reactions are glucose and pyruvate consuming reactions and coenzyme reactions which are vital to drive CFPS. This can be seen in the biphasic operation of CFPS, with the first phase operating on glucose and the second phase operating on pyruvate. During the first phase, there is an accumulation of byproducts from central carbon with the majority of flux going toward acetate and some toward pyruvate, lactate, and succinate; with the exception of acetate, these are all consumed in the second phase. This shows that CAT production can be sustained by pyruvate and glutamate in the absence of glucose, which provides alternative strategies to optimize CFPS performance. This is in accordance with literature, which showed pyruvate provided a relatively slow but continuous supply of ATP [27]. Taken together, this shows CFPS can be designed towards a specified application either requiring a slow stable energy source or faster production. This outstanding control on model performance was expected as these metabolites are responsible for driving CFPS and represent the first step in the model network. Nevertheless, there are further reactions with considerable impact on model performance. In examining oxidative phosphorylation activity, knockouts in the electron transport pathways disrupt metabolism across the network and show CAT carbon yield dropping from

8.6% to 2.7%; Jewett and coworkers also saw a decrease in CAT yield, ranging from 1.5-fold to 4-fold, when knocking out oxidative phosphorylation reactions[26]. Oxidative phosphorylation is vital, since it provides most of the energetic needs of CFPS. However, it is unknown how active oxidative phosphorylation is compared to that of *in vivo* systems, and both of our modeling approaches suggest its importance to CAT production and CFPS. Thus, oxidative phosphorylation is a potential area for improvement for CFPS performance and protein yield. Comparing the theoretical maximum carbon yield of CAT from ssFBA predictions to those of the kinetic model and experimental measurements suggests that there is potential for increasing CAT yield as well as CFPS performance. The model and experimental yields were 36% of the theoretical maximum and 23% of a physiologically constrained case. Knockouts of *gnd* and phosphate acetyltransferase show that carbon can be diverted away from acetate and potentially toward CAT or other proteins of interest expressed in CFPS. Another limitation to be addressed in CFPS is the transcription and translation description, since protein production is ultimately bounded by these kinetic rates. Li et al. have increased productivity of firefly luciferase by 5-fold in CFPS systems by adding and adjusting factors that affect transcription and translation such as elongation factors, ribosome recycling factor, release factors, chaperones, BSA, and tRNAs [28]. Underwood and coworkers have also shown that an increase in ribosome levels does not significantly increase protein yields or rates; however, adding elongation factors increased yields by 23% at 30 minutes [29].

A logical next step for this work would be sequence-specific dynamic modeling, as the kinetic modeling approach in this study used a single reaction to approximate CAT synthesis. Including specific transcription and translation steps for CAT would allow more accurate modeling of the complexity and the resource cost of protein synthesis. In addition, sensitivity analysis could be performed on these new parameters to determine the robustness of CAT synthesis to the processes of transcription and translation. Another

area for future work is to more thoroughly sample parameter space. Parameters were varied so as to best fit the dataset; however, the resulting ensemble may not represent every biological possibility. In a different region of parameter space, the system may behave differently but still fit the experimental data. This could include the flux distribution through the network, the variation of predictions across the ensemble, and the relative sensitivity values. Testing the model under a variety of conditions could strengthen or challenge the findings of this study. Further experimentation could also be used to gain a deeper understanding of model performance under a variety of conditions. Specifically, CAT production performed in the absence of amino acids could inform the system's ability to manufacture them, while experimentation in the absence of glucose or oxygen could shed light on how important they are to protein synthesis, and under which conditions. Finally, the approach should be extended to other protein products. CAT is only a test protein used for model identification; the modeling framework, and to some extent the parameter values, should be protein agnostic. An important extension of this study would be to apply its insights to other protein applications, where possible.

Materials and Methods

Formulation and solution of the model equations. We used ordinary differential equations (ODEs) to model the time evolution of metabolite (x_i) and scaled enzyme abundance (ϵ_i) in hypothetical cell-free metabolic networks:

$$\frac{dx_i}{dt} = \sum_{j=1}^{\mathcal{R}} \sigma_{ij} r_j(\mathbf{x}, \epsilon, \mathbf{k}) \quad i = 1, 2, \dots, \mathcal{M} \quad (1)$$

$$\frac{d\epsilon_i}{dt} = -\lambda_i \epsilon_i \quad i = 1, 2, \dots, \mathcal{E} \quad (2)$$

where \mathcal{R} denotes the number of reactions, \mathcal{M} denotes the number of metabolites and \mathcal{E} denotes the number of enzymes in the model. The quantity $r_j(\mathbf{x}, \epsilon, \mathbf{k})$ denotes the rate of reaction j . Typically, reaction j is a non-linear function of metabolite and enzyme abundance, as well as unknown kinetic parameters \mathbf{k} ($\mathcal{K} \times 1$). The quantity σ_{ij} denotes the stoichiometric coefficient for species i in reaction j . If $\sigma_{ij} > 0$, metabolite i is produced by reaction j . Conversely, if $\sigma_{ij} < 0$, metabolite i is consumed by reaction j , while $\sigma_{ij} = 0$ indicates metabolite i is not connected with reaction j . Lastly, λ_i denotes the scaled enzyme activity decay constant. The system material balances were subject to the initial conditions $\mathbf{x}(t_o) = \mathbf{x}_o$ and $\epsilon(t_o) = 1$ (initially we have 100% cell-free enzyme abundance).

The reaction rate was written as the product of a kinetic term (\bar{r}_j) and a control term (v_j), $r_j(\mathbf{x}, \mathbf{k}) = \bar{r}_j v_j$. We used multiple saturation kinetics to model the reaction term \bar{r}_j :

$$\bar{r}_j = V_j^{max} \epsilon_i \prod_{s \in m_j^-} \frac{x_s}{K_{js} + x_s} \quad (3)$$

where V_j^{max} denotes the maximum rate for reaction j , ϵ_i denotes the scaled enzyme activity which catalyzes reaction j , K_{js} denotes the saturation constant for species s in reaction j and m_j^- denotes the set of *reactants* for reaction j . On the other hand, the control term $0 \leq v_j \leq 1$ depended upon the combination of factors which influenced

rate process j . For each rate, we used a rule-based approach to select from competing control factors. If rate j was influenced by $1, \dots, m$ factors, we modeled this relationship as $v_j = \mathcal{I}_j(f_{1j}(\cdot), \dots, f_{mj}(\cdot))$ where $0 \leq f_{ij}(\cdot) \leq 1$ denotes a transfer function quantifying the influence of factor i on rate j . The function $\mathcal{I}_j(\cdot)$ is an integration rule which maps the output of regulatory transfer functions into a control variable. We used hill-like transfer functions and $\mathcal{I}_j \in \{min, max\}$ in this study [25].

We included 17 allosteric regulation terms, taken from literature, in the CFPS model. PEP was modeled as an inhibitor for phosphofructokinase [30, 31], PEP carboxykinase [30], PEP synthetase [30, 32], isocitrate dehydrogenase [30, 33], and isocitrate lyase/malate synthase [30, 33, 34], and as an activator for fructose-biphosphatase [30, 35–37]. AKG was modeled as an inhibitor for citrate synthase [30, 38, 39] and isocitrate lyase/malate synthase [30, 34]. 3PG was modeled as an inhibitor for isocitrate lyase/malate synthase [30, 34]. FDP was modeled as an activator for pyruvate kinase [30, 40] and PEP carboxylase [30, 41]. Pyruvate was modeled as an inhibitor for pyruvate dehydrogenase [30, 42, 43] and as an activator for lactate dehydrogenase [44]. Acetyl CoA was modeled as an inhibitor for malate dehydrogenase [30].

Estimation of kinetic model parameters. We estimated an ensemble of diverse parameter sets using a constrained Markov Chain Monte Carlo (MCMC) random walk strategy. Starting from a single best fit parameter set estimated by inspection and literature, we calculated the cost function, equal to the sum-squared-error between experimental data and model predictions:

$$\text{cost} = \sum_{i=1}^{\mathcal{D}} \left[\frac{w_i}{\mathcal{Y}_i^2} \sum_{j=1}^{\mathcal{T}_i} \left(y_{ij} - x_i|_{t(j)} \right)^2 \right] \quad (4)$$

where \mathcal{D} denotes the number of datasets ($\mathcal{D} = 37$), w_i denotes the weight of the i^{th} dataset, \mathcal{T}_i denotes the number of timepoints in the i^{th} dataset, $t(j)$ denotes the j^{th} time-

point, y_{ij} denotes the measurement value of the i^{th} dataset at the j^{th} timepoint, and $x_i|_{t(j)}$ denotes the simulated value of the metabolite corresponding to the i^{th} dataset, interpolated to the j^{th} timepoint. Lastly, the cost calculation was scaled by the maximum experimental value in the i^{th} dataset, $\mathcal{Y}_i = \max_j (y_{ij})$. We then perturbed each model parameter between an upper and lower bound that varied by parameter type:

$$k_i^{new} = \min(\max(k_i \cdot \exp(a \cdot r_i), l_i), u_i) \quad i = 1, 2, \dots, \mathcal{P} \quad (5)$$

where \mathcal{P} denotes the number of parameters ($\mathcal{P} = 815$), which includes 163 maximum reaction rates (V^{max}), 163 enzyme activity decay constants, 455 saturation constants (K_{js}), and 34 control parameters, k_i^{new} denotes the new value of the i^{th} parameter, k_i denotes the current value of the i^{th} parameter, a denotes a distribution variance, r_i denotes a random sample from the normal distribution, l_i denotes the lower bound for that parameter type, and u_i denotes the upper bound for that parameter type. Maximum reaction rates were bounded between 0 and 500,000 mM/h [45]. Assuming a total enzyme concentration of 50 μ M, this corresponds to catalytic rate bounds of 0 and 2778 s^{-1} . These bounds resulted in a median catalytic rate of 0.016 s^{-1} across the ensemble. Enzyme activity decay constants were bounded between 0 and 1 h^{-1} , corresponding to half lives of 42 minutes and infinity; median = 25 min. Saturation constants were bounded between 0.001 and 10 mM; median = 0.16 mM. Control parameters (gains and orders) were left unbounded; gain median = 0.076, order median = 0.69. For each newly generated parameter set, we re-solved the balance equations and calculated the cost function. All sets with a lower cost (and some with higher cost) were accepted into the ensemble. After generating greater than 10,000 sets, we selected $N = 100$ sets with minimal set to set correlation to avoid over-sampling any region of parameter space.

Sensitivity analysis of the kinetic CFPS model. We determined the reactions most important to protein production by computing the local sensitivity of CAT concentration (denoted as CAT) to each individual maximum reaction rate, and each pair of maximum reaction rates in the network. The sensitivity index was formulated as:

$$\mathcal{S}_{ij}^{\text{CAT}} = \|\text{CAT}(p_i, p_j, t) - \text{CAT}(\alpha \cdot p_i, \alpha \cdot p_j, t)\|_2 \quad i, j = 1, 2, \dots, \mathcal{P} \quad (6)$$

where $\mathcal{S}_{ij}^{\text{CAT}}$ denotes the sensitivity of CAT production to the i^{th} and j^{th} parameters, $\text{CAT}(p_i, p_j, t)$ denotes CAT concentration as a function of time and the i^{th} and j^{th} parameters, α denotes the perturbation factor, and \mathcal{P} denotes the number of maximum reaction rates ($\mathcal{P} = 163$). In calculating the pairwise sensitivities, each parameter was perturbed by 1%; first-order sensitivities ($i = j$) were subject to two 1% perturbations. Parameters and parameter combinations were stratified into five degrees of importance, from least to most sensitive.

Likewise, we determined which reactions were most important to global system performance by computing the sensitivity of all species for which data exists (denoted as X) to each maximum reaction rate in the network. In this case, each sensitivity index was formulated as:

$$\mathcal{S}_{ij}^{\text{X}} = \|\text{X}(p_i, p_j, t) - \text{X}(\alpha \cdot p_i, \alpha \cdot p_j, t)\|_2 \quad i, j = 1, 2, \dots, \mathcal{P} \quad (7)$$

where $\mathcal{S}_{ij}^{\text{X}}$ denotes the sensitivity of the system state to the i^{th} and j^{th} parameters, and $\text{X}(p_i, p_j, t)$ denotes the system state, an array consisting of the concentration of every species for which data exists as a function of time and the i^{th} and j^{th} parameters. The parameter sensitivities were stratified into five degrees of importance, from least to most sensitive, as above.

Sequence specific calculation of carbon yield. We estimated the theoretical maximum carbon yield of CAT using sequence-specific flux balance analysis (ssFBA) [46]. The CAT carbon yield (Y_C^{CAT}) was calculated as the ratio of carbon produced as CAT divided by the carbon consumed as reactants (glucose and amino acids):

$$Y_C^{CAT} = \frac{\Delta CAT \cdot C_{CAT}}{\sum_{i=1}^{\mathcal{R}} \max(\Delta m_i, 0) \cdot C_{m_i}} \quad (8)$$

where ΔCAT denotes the abundance of CAT produced, C_{CAT} denotes carbon number of CAT, \mathcal{R} denotes the number of reactants, Δm_i denotes the amount of the i^{th} reactant consumed (never allowed to be negative), and C_{m_i} denotes the carbon number of the i^{th} reactant. Arginine or glutamate were not considered in the yield calculations, as no experimental measurements were available for these amino acids. Yield of the best-fit parameter set and the experimental data were calculated by setting ΔCAT equal to the final minus the initial CAT concentration, and setting Δm_i equal to the initial minus the final reactant concentration.

The sequence specific flux balance analysis problem was formulated as:

$$\begin{aligned} & \max_{\mathbf{w}} (w_{obj} = \boldsymbol{\theta}^T \mathbf{w}) \\ & \text{Subject to : } \mathbf{S}\mathbf{w} = \mathbf{0} \\ & \alpha_i \leq w_i \leq \beta_i \quad i = 1, 2, \dots, \mathcal{R} \end{aligned} \quad (9)$$

where \mathbf{S} denotes the stoichiometric matrix, \mathbf{w} denotes the unknown flux vector, $\boldsymbol{\theta}$ denotes the objective selection vector and α_i and β_i denote the lower and upper bounds on flux w_i , respectively. The objective was to maximize the CAT translation rate. The glucose uptake rate was bounded in the range [0,40 mM/h] according to experimental data; while the amino acid uptake rates were bounded by [0,30 mM/h], but did not reach the maximum

flux. The transcription and translation rates were modeled using template reactions

The transcription rate (w_{TX}) was fixed in the ssFBA calculation as:

$$w_{TX} = \left[R_1 \left(\frac{K_{T7}}{1 + K_{T7}} \right) \left(\frac{v_{RNAP}}{l_{mRNA}} \right) \right] \left(\frac{G}{K_{TX} + G} \right) \quad (10)$$

where R_1 denotes the concentration of RNA polymerase, v_{TX} denotes RNA polymerase elongation rate (nt/hr), l_G denotes the gene length in nucleotides, G denotes the gene concentration and K_{TX} denotes the plasmid saturation coefficient. The gene and protein sequence for CAT was determined from literature. The last term in the w_{TX} expression describes T7 promoter activity, where K_{T7} quantifies T7 RNA polymerase binding [REFHERE]. The translation rate was formulated as:

$$w_{TL} = K_P \text{ Ribo} \left(\frac{v_{Ribo}}{l_{Protein}} \right) mRNA_{SS} \quad (11)$$

where K_P denotes the polysome amplification constant, $Ribo$ denotes the ribosome concentration, v_{Ribo} denotes the elongation rate of the ribosome in amino acids per hour, $l_{Protein}$ denotes the number of amino acids in the protein of interest, and $mRNA_{SS}$ denotes the mRNA concentration at steady state, equal to the transcription rate divided by the degradation rate of mRNA.

An ensemble of 100 sets of flux distributions was calculated for three different cases: unconstrained, constrained by transcription/translation rates, and constrained by transcription/translation rates and experimental measurements. For the unconstrained case, all rates were left unbounded, except the specific glucose uptake rate. An ensemble of flux distributions was then calculated by randomly sampling the maximum specific glucose uptake rate from within a range of 30 to 40 mM/h, determined from experimental data. For the case constrained by transcription/translation rates, an ensemble was generated by randomly sampling RNAP polymerase levels, ribosome levels, and elongation rates in

a physiological range determined from literature. RNA polymerase levels were sampled between 60 and 80 nM, ribosome levels between 7 and 16 μ M, the RNA polymerase elongation rate between 20 and 30 nt/sec, and the ribosome elongation rate between 1.5 and 3 aa/sec [29, 47]. For the case constrained by transcription/translation rates and experimental measurements, the lower and upper bounds on the fluxes for the data-informed metabolites were sampled within the range given by the experimental noise. This included the data for glucose, organic acids, energy species, and amino acids; CAT was not constrained by experimental data, but by the transcription/translation rates as stated above. The flux balance analysis problem was solved using the GNU Linear Programming Kit (v4.52) [48].

Competing interests

The authors declare that they have no competing interests.

Author's contributions

J.V and A.Y directed the study. R.T, H.J and J.C conducted the cell culture measurements. J.V and W.D developed the reduced order HL-60 models and the parameter ensemble. W.D analyzed the model ensemble, and generated figures for the manuscript. The manuscript was prepared and edited for publication by W.D, A.Y and J.V.

Acknowledgements

We gratefully acknowledge the suggestions from the anonymous reviewers to improve this manuscript.

Funding

This study was supported by a National Science Foundation Graduate Research Fellowship (DGE-1333468) to N.H and by an award from the US Army and Systems Biology of Trauma Induced Coagulopathy (W911NF-10-1-0376) to J.V for the support of M.V.

References

1. Jewett MC, Calhoun KA, Voloshin A, Wu JJ, Swartz JR (2008) An integrated cell-free metabolic platform for protein production and synthetic biology. *Mol Syst Biol* 4: 220.
2. Matthaei JH, Nirenberg MW (1961) Characteristics and stabilization of dnaase-sensitive protein synthesis in e. coli extracts. *Proc Natl Acad Sci U S A* 47: 1580-8.
3. Nirenberg MW, Matthaei JH (1961) The dependence of cell-free protein synthesis in e. coli upon naturally occurring or synthetic polyribonucleotides. *Proc Natl Acad Sci U S A* 47: 1588-602.
4. Lu Y, Welsh JP, Swartz JR (2014) Production and stabilization of the trimeric influenza hemagglutinin stem domain for potentially broadly protective influenza vaccines. *Proc Natl Acad Sci U S A* 111: 125-30.
5. Hodgman CE, Jewett MC (2012) Cell-free synthetic biology: thinking outside the cell. *Metab Eng* 14: 261-9.
6. Pardee K, Slomovic S, Nguyen PQ, Lee JW, Donghia N, et al. (2016) Portable, on-demand biomolecular manufacturing. *Cell* 167: 248-59.e12.
7. Fredrickson AG (1976) Formulation of structured growth models. *Biotechnol Bioeng* 18: 1481-6.
8. Domach MM, Leung SK, Cahn RE, Cocks GG, Shuler ML (1984) Computer model for glucose-limited growth of a single cell of escherichia coli b/r-a. *Biotechnol Bioeng* 26: 203-16.
9. Steinmeyer D, Shuler M (1989) Structured model for *Saccharomyces cerevisiae*. *Chem Eng Sci* 44: 2017-30.
10. Wu P, Ray NG, Shuler ML (1992) A single-cell model for cho cells. *Ann N Y Acad Sci* 665: 152-87.
11. Castellanos M, Wilson DB, Shuler ML (2004) A modular minimal cell model: purine

and pyrimidine transport and metabolism. *Proc Natl Acad Sci U S A* 101: 6681-6.

12. Atlas JC, Nikolaev EV, Browning ST, Shuler ML (2008) Incorporating genome-wide dna sequence information into a dynamic whole-cell model of *escherichia coli*: application to dna replication. *IET Syst Biol* 2: 369-82.
13. Lewis NE, Nagarajan H, Palsson BØ (2012) Constraining the metabolic genotype-phenotype relationship using a phylogeny of in silico methods. *Nat Rev Microbiol* 10: 291-305.
14. Edwards JS, Palsson BØ (2000) The *escherichia coli* mg1655 in silico metabolic genotype: its definition, characteristics, and capabilities. *Proc Natl Acad Sci U S A* 97: 5528-33.
15. Feist AM, Herrgård MJ, Thiele I, Reed JL, Palsson BØ (2009) Reconstruction of biochemical networks in microorganisms. *Nat Rev Microbiol* 7: 129-43.
16. Feist AM, Henry CS, Reed JL, Krummenacker M, Joyce AR, et al. (2007) A genome-scale metabolic reconstruction for *escherichia coli* k-12 mg1655 that accounts for 1260 orfs and thermodynamic information. *Mol Syst Biol* 3: 121.
17. Oh YK, Palsson BØ, Park SM, Schilling CH, Mahadevan R (2007) Genome-scale reconstruction of metabolic network in *bacillus subtilis* based on high-throughput phenotyping and gene essentiality data. *J Biol Chem* 282: 28791-9.
18. Ibarra RU, Edwards JS, Palsson BØ (2002) *Escherichia coli* k-12 undergoes adaptive evolution to achieve in silico predicted optimal growth. *Nature* 420: 186-9.
19. Schuetz R, Kuepfer L, Sauer U (2007) Systematic evaluation of objective functions for predicting intracellular fluxes in *escherichia coli*. *Mol Syst Biol* 3: 119.
20. Hyduke DR, Lewis NE, Palsson BØ (2013) Analysis of omics data with genome-scale models of metabolism. *Mol Biosyst* 9: 167-74.
21. McCloskey D, Palsson BØ, Feist AM (2013) Basic and applied uses of genome-scale metabolic network reconstructions of *escherichia coli*. *Mol Syst Biol* 9: 661.

- 496 22. Zomorodi AR, Suthers PF, Ranganathan S, Maranas CD (2012) Mathematical opti-
497 mization applications in metabolic networks. *Metab Eng* 14: 672-86.
- 498 23. Feist AM, Henry CS, Reed JL, Krummenacker M, Joyce AR, et al. (2007) A genome-
499 scale metabolic reconstruction for escherichia coli k-12 mg1655 that accounts for
500 1260 orfs and thermodynamic information. *Molecular Systems Biology* 3.
- 501 24. Calhoun KA, Swartz JR (2005) An economical method for cell-free protein synthesis
502 using glucose and nucleoside monophosphates. *Biotechnology Progress* 21: 1146-
503 53.
- 504 25. Wayman JA, Sagar A, Varner JD (2015) Dynamic modeling of cell-free biochemical
505 networks using effective kinetic models. *Processes* 3: 138.
- 506 26. Jewett MC, Calhoun KA, Voloshin A, Wu JJ, Swartz JR (2008) An integrated cell-free
507 metabolic platform for protein production and synthetic biology. *Molecular Systems*
508 *Biology* 4.
- 509 27. Swartz J (2001) A pure approach to constructive biology. *Nature Biotechnology* 19:
510 732-3.
- 511 28. Li J, Gu L, Aach J, Church GM (2014) Improved cell-free rna and protein synthesis
512 system. *PLoS ONE* 9: 1-11.
- 513 29. Underwood KA, Swartz JR, Puglisi JD (2005) Quantitative polysome analysis iden-
514 tifies limitations in bacterial cell-free protein synthesis. *Biotechnology and Bioengi-*
515 *neering* 91: 425-35.
- 516 30. Kotte O, Zaugg JB, Heinemann M (2010) Bacterial adaptation through distributed
517 sensing of metabolic fluxes. *Mol Syst Biol* 6: 355.
- 518 31. Cabrera R, Baez M, Pereira HM, Caniuguir A, Garratt RC, et al. (2011) The crys-
519 tal complex of phosphofructokinase-2 of Escherichia coli with fructose-6-phosphate:
520 kinetic and structural analysis of the allosteric ATP inhibition. *J Biol Chem* 286: 5774-
521 83.

32. Chulavatnatol M, Atkinson DE (1973) Phosphoenolpyruvate synthetase from *Escherichia coli*. Effects of adenylate energy charge and modifier concentrations. *J Biol Chem* 248: 2712-5.
33. Ogawa T, Murakami K, Mori H, Ishii N, Tomita M, et al. (2007) Role of phosphoenolpyruvate in the NADP-isocitrate dehydrogenase and isocitrate lyase reaction in *Escherichia coli*. *J Bacteriol* 189: 1176-8.
34. MacKintosh C, Nimmo HG (1988) Purification and regulatory properties of isocitrate lyase from *Escherichia coli* ML308. *Biochem J* 250: 25-31.
35. Donahue JL, Bownas JL, Niehaus WG, Larson TJ (2000) Purification and characterization of glpX-encoded fructose 1, 6-bisphosphatase, a new enzyme of the glycerol 3-phosphate regulon of *Escherichia coli*. *J Bacteriol* 182: 5624-7.
36. Hines JK, Fromm HJ, Honzatko RB (2006) Novel allosteric activation site in *Escherichia coli* fructose-1,6-bisphosphatase. *J Biol Chem* 281: 18386-93.
37. Hines JK, Fromm HJ, Honzatko RB (2007) Structures of activated fructose-1,6-bisphosphatase from *Escherichia coli*. Coordinate regulation of bacterial metabolism and the conservation of the R-state. *J Biol Chem* 282: 11696-704.
38. Pereira DS, Donald LJ, Hosfield DJ, Duckworth HW (1994) Active site mutants of *Escherichia coli* citrate synthase. Effects of mutations on catalytic and allosteric properties. *J Biol Chem* 269: 412-7.
39. Robinson MS, Easom RA, Danson MJ, Weitzman PD (1983) Citrate synthase of *Escherichia coli*. Characterisation of the enzyme from a plasmid-cloned gene and amplification of the intracellular levels. *FEBS Lett* 154: 51-4.
40. Zhu T, Bailey MF, Angley LM, Cooper TF, Dobson RC (2010) The quaternary structure of pyruvate kinase type 1 from *Escherichia coli* at low nanomolar concentrations. *Biochimie* 92: 116-20.
41. Wohl RC, Markus G (1972) Phosphoenolpyruvate carboxylase of *Escherichia coli*.

Purification and some properties. J Biol Chem 247: 5785-92.

42. Kale S, Arjunan P, Furey W, Jordan F (2007) A dynamic loop at the active center of the Escherichia coli pyruvate dehydrogenase complex E1 component modulates substrate utilization and chemical communication with the E2 component. J Biol Chem 282: 28106-16.
43. Arjunan P, Nemeria N, Brunskill A, Chandrasekhar K, Sax M, et al. (2002) Structure of the pyruvate dehydrogenase multienzyme complex E1 component from Escherichia coli at 1.85 Å resolution. Biochemistry 41: 5213-21.
44. Okino S, Suda M, Fujikura K, Inui M, Yukawa H (2008) Production of D-lactic acid by Corynebacterium glutamicum under oxygen deprivation. Appl Microbiol Biotechnol 78: 449-54.
45. Keseler IM, Mackie A, Peralta-Gil M, Santos-Zavaleta A, Gama-Castro S, et al. (2013) Ecocyc: fusing model organism databases with systems biology. Nucleic Acids Res 41: 605-12.
46. Allen TE, Palsson BØ (2003) Sequence-based analysis of metabolic demands for protein synthesis in prokaryotes. J Theor Biol 220: 1-18.
47. Garamella J, Marshall R, Rustad M, Noireaux V (2016) The all e. coli tx-tl toolbox 2.0: A platform for cell-free synthetic biology. ACS Synth Biol 5: 344-55.
48. (2016). GNU Linear Programming Kit, Version 4.52. URL <http://www.gnu.org/software/glpk/glpk.html>.

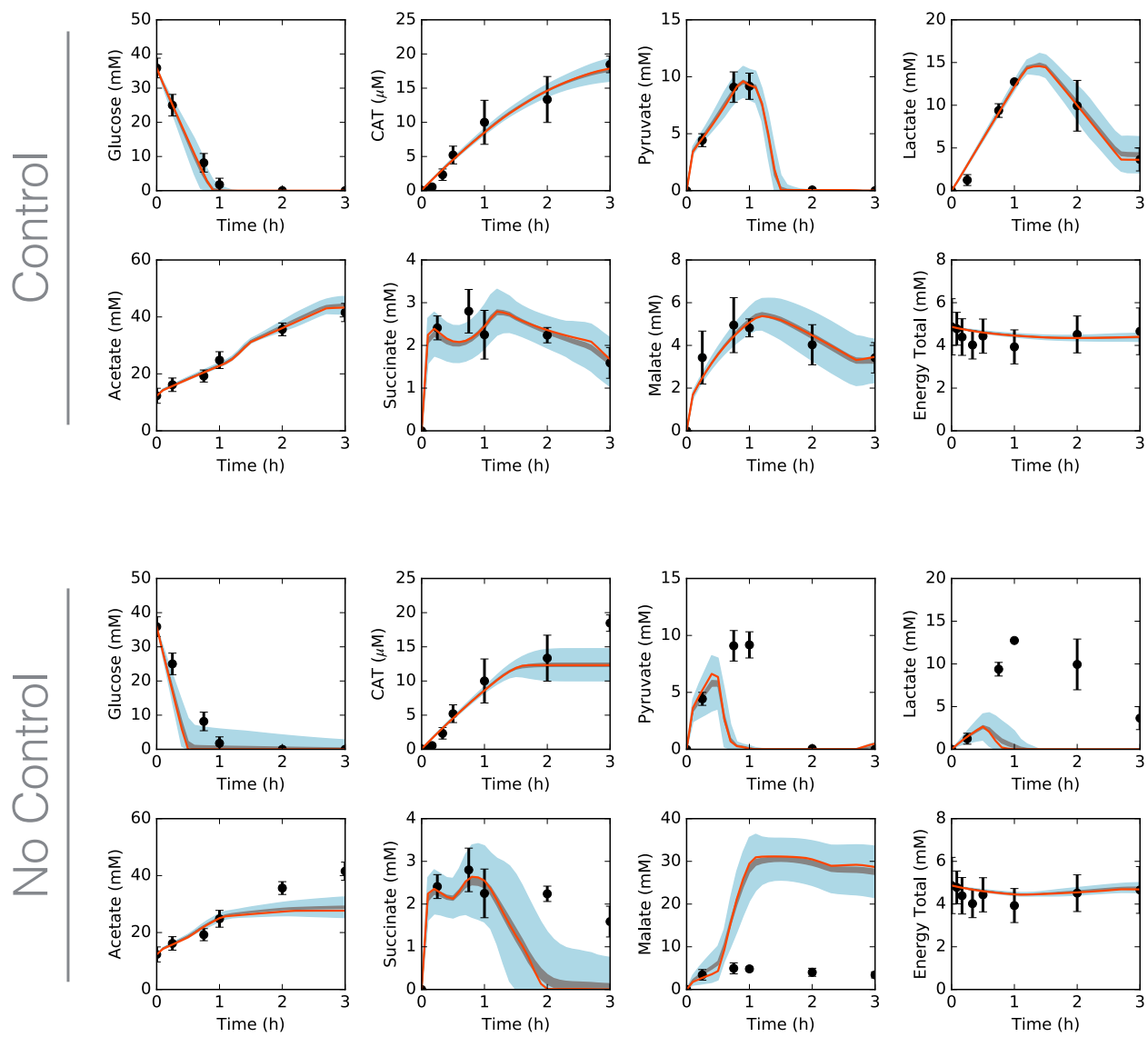


Fig. 1: Central carbon metabolism in the presence (top) and absence (bottom) of allosteric control, including glucose (substrate), CAT (product), and intermediates, as well as total concentration of energy species. Best-fit parameter set (orange line) versus experimental data (points). 95% confidence interval (blue shaded region) and 95% confidence interval of the mean (gray shaded region) over the ensemble of 100 sets.

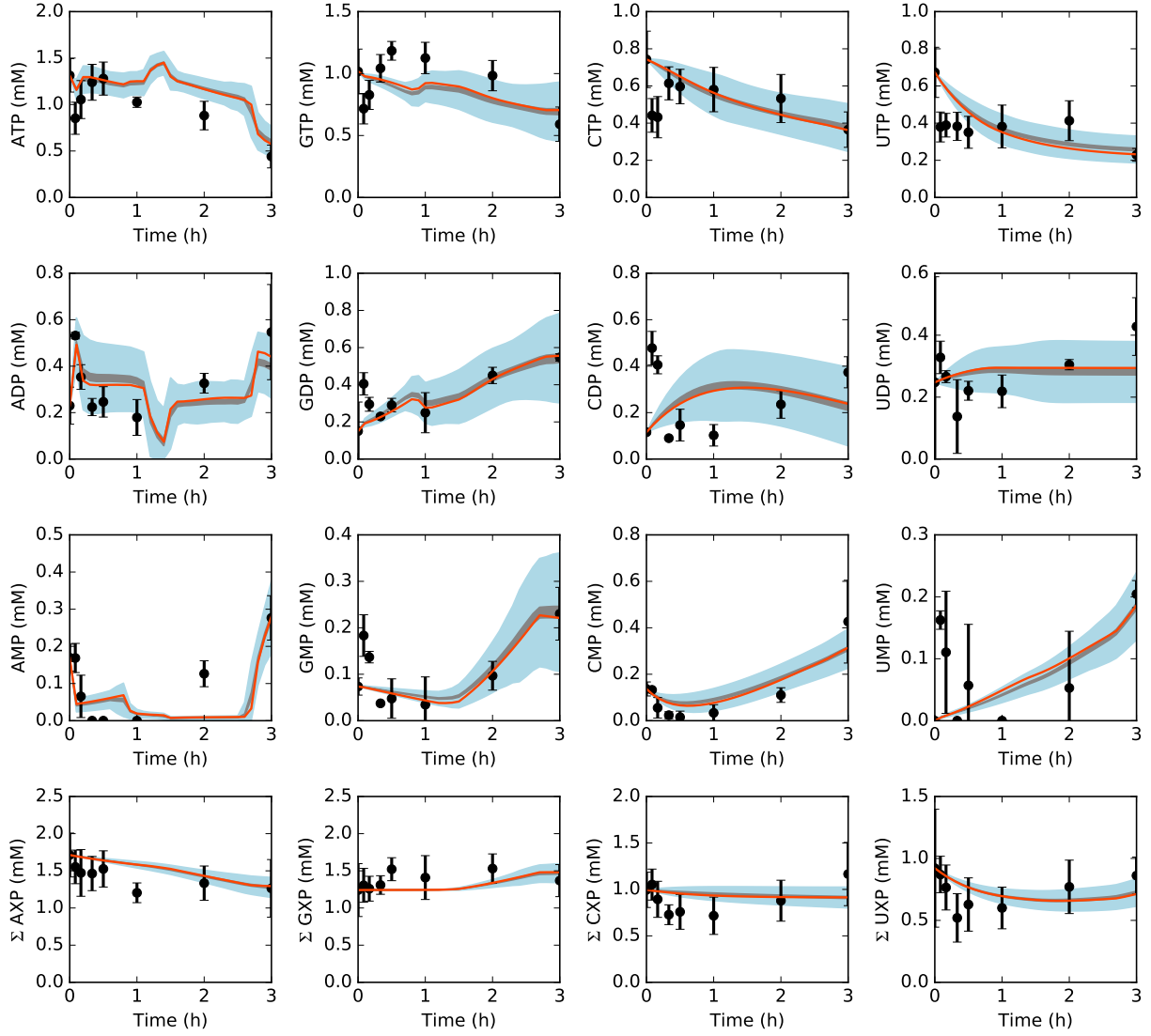


Fig. 2: Energy species and energy totals by base in the presence of allosteric control. Best-fit parameter set (orange line) versus experimental data (points). 95% confidence interval (blue shaded region) and 95% confidence interval of the mean (gray shaded region) over the ensemble of 100 sets.

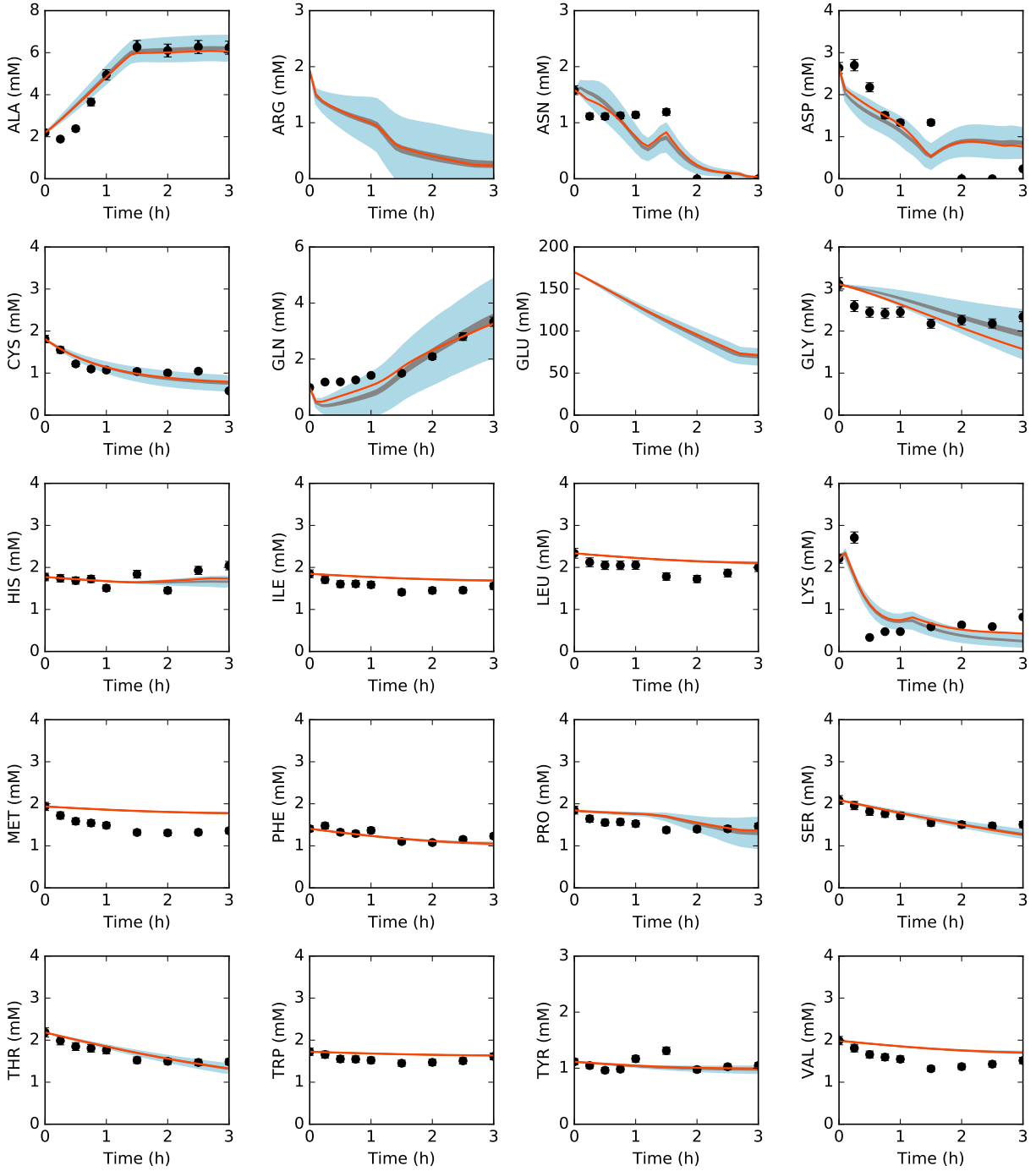


Fig. 3: Amino acids in the presence of allosteric control. Best-fit parameter set (orange line) versus experimental data (points). 95% confidence interval (blue shaded region) and 95% confidence interval of the mean (gray shaded region) over the ensemble of 100 sets.

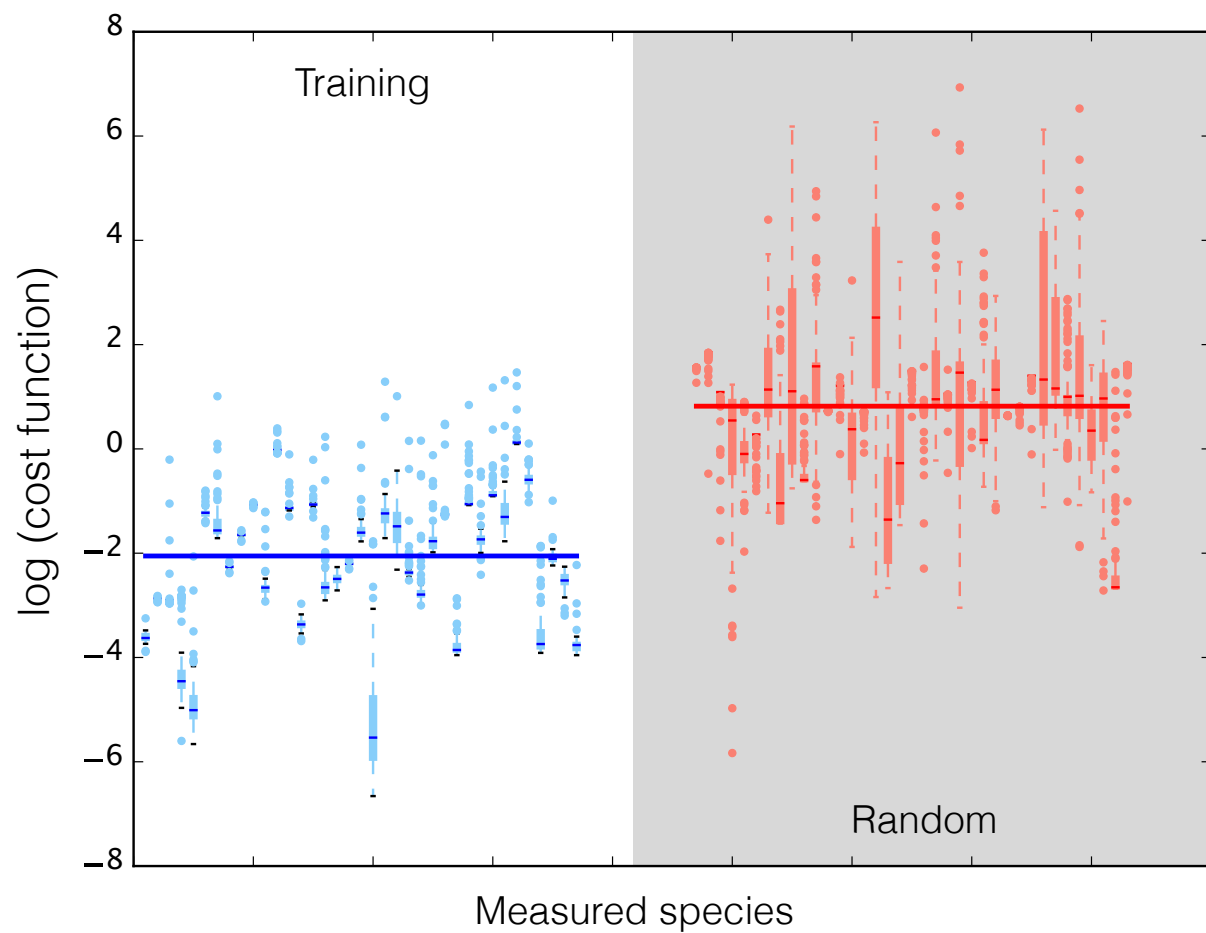


Fig. 4: Log of cost function across 37 datasets for data-trained ensemble (blue) and randomly generated ensemble (red, gray background). Median (bars), interquartile range (boxes), range excluding outliers (dashed lines), and outliers (circles) for each dataset. Median across all datasets (large bar overlaid).

Fig. 5: Normalized first-order and pairwise sensitivities of CAT production (top) and system state (bottom) to maximum reaction rates.

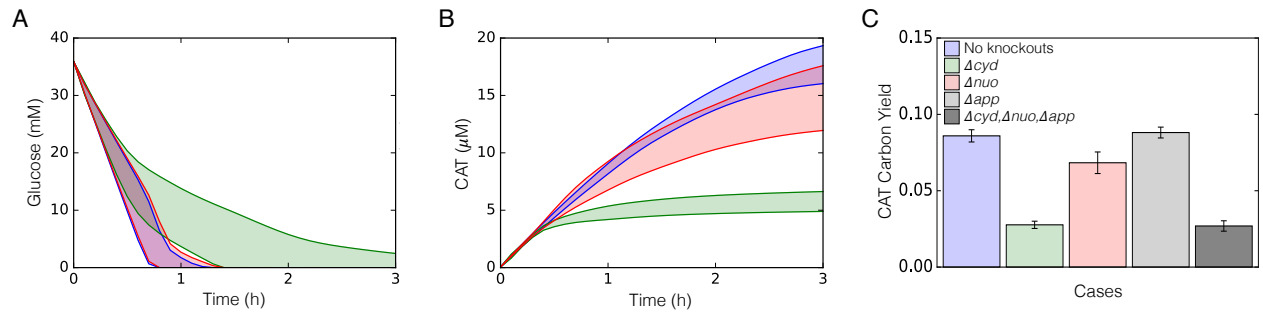


Fig. 6: The effect of oxidative phosphorylation on glucose uptake, CAT production and CAT carbon yield. A. 95% confidence interval of an ensemble for glucose concentration versus time for no knockouts (blue shaded region), *cyd* knockout (green shaded region), and *nuo* knockout (red shaded region). B. 95% confidence interval of an ensemble for CAT concentration versus time for no knockouts (blue shaded region), *cyd* knockout (green shaded region), and *nuo* knockout (red shaded region). C. CAT carbon yield for 5 different cases of oxidative phosphorylation: no knockouts (blue), *cyd* knockout (green), *nuo* knockout (red), *app* knockout (light grey), and a combination of *cyd*, *nuo*, *app* knockouts (dark grey).

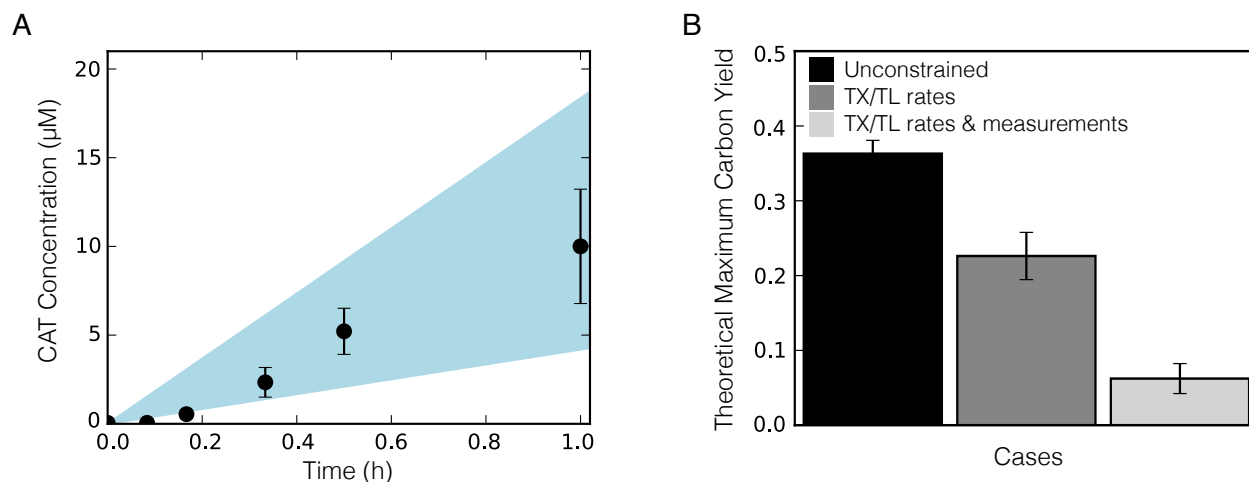


Fig. 7: Sequence-specific flux balance analysis of CAT production and yield. A. 95% confidence interval of the ensemble (light blue region) for CAT concentration versus time. B. Theoretical maximum carbon yield of CAT calculated by ssFBA for three different cases: unconstrained except for glucose uptake (black), constrained by transcription/translation (TX/TL) rates (grey), and constrained by transcription/translation (TX/TL) rates and experimental measurements where available (light grey). Error bars represent standard deviation of the ensemble.

