

A Machine Learning Approach to Predict Metabolic Pathway Dynamics from Time Series Multiomics Data

Zak Costello^{1,2,3} and Hector Garcia Martin^{1,2,3,4,*}

¹Biological Systems and Engineering Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA

²DOE Agile Biofoundry, Emeryville, CA, USA

³DOE Joint BioEnergy Institute, Emeryville, CA, USA

⁴BCAM, Basque Center for Applied Mathematics, Bilbao, Spain

*hgmartin@lbl.gov

ABSTRACT

New synthetic biology capabilities hold the promise of dramatically improving our ability to engineer biological systems. However, a fundamental hurdle in realizing this potential is our inability to accurately predict biological behavior after modifying the corresponding genotype. Kinetic models have traditionally been used to predict pathway dynamics in bioengineered systems, but they take significant time to develop, and rely heavily on domain expertise. Here, we show that the combination of machine learning and abundant multiomics data (proteomics and metabolomics) can be used to effectively predict pathway dynamics in an automated fashion. The new method outperforms a classical kinetic model, and produces qualitative and quantitative predictions that can be used to productively guide bioengineering efforts. This method systematically leverages arbitrary amounts of new data to improve predictions, and does not assume any particular interactions, but rather implicitly chooses the most predictive ones.

Introduction

Biology has been transformed in the second half of the 20th century from a descriptive science to a design science. This transformation has been produced by a combination of the discovery of DNA as the repository of genetic information¹, and of recombinant DNA as an effective way to alter this instruction set². The subsequent advent of genetic engineering and synthetic biology as effective tools to engineer biological cells has produced numerous beneficial applications ranging from the production of renewable biofuels and other bioproducts³⁻⁶ to applications in human health⁷⁻⁹, creating the expectation of an industrialized biology affecting almost every facet of human activity¹⁰.

However, effective design of biological systems is precluded by our inability to predict their behavior. We can engineer changes faster than ever, enabled by DNA synthesis productivity that improves as fast as Moore's law¹¹, and new tools like CRISPR-enabled genetic editing, which have revolutionized our ability to modify the DNA in vivo¹². In general, we can make the DNA changes we intend (in model systems), but the end result on cell behavior is usually unpredictable¹³. At the same time, there is an exponentially increasing amount of functional genomics data available to the experimenter in order to phenotype the resulting bioengineered organism: transcriptomics data volume has a doubling rate of 7 months¹⁴, and high-throughput workflows for proteomics¹⁵ and metabolomics¹⁶ are becoming increasingly available. Furthermore, the miniaturization of these techniques and the progressive automation of laboratory work through microfluidics chips promises a future where data analysis will be the bottleneck in biological research¹⁷. Unfortunately, the availability of all this data does not translate into better predictive capabilities for biological systems: converting these data into actionable insights to achieve a given goal (e.g. higher bioproduct yields) is far from trivial or routine.

Mathematical modeling provides a systematic manner to leverage these data to predict the behavior of engineered systems. Hence, increasingly, computational biology is focusing on large-scale modeling of dynamical systems predicting phenotype from genotype^{18,19}. However, computational biology is still nascent and not able to provide the high accuracy predictions that we are accustomed to seeing in other engineering fields²⁰. Arguably, the most widely used and successful modeling technique in metabolic engineering involves analysis of internal metabolic fluxes (i.e. reaction rates) through stoichiometric models of metabolism. Metabolic flux values are constrained by stoichiometry, thermodynamic and evolutionary assumptions^{21,22}, or experimental data (e.g. ¹³C labeling experimental data²³⁻²⁵), and used to suggest genetic interventions that bring cell metabolism closer to the desired phenotype. While this approach has provided significant successes²⁶⁻²⁹, it has also shown its limitations³⁰ due to its simplicity. Stoichiometric models are limited for bioengineering purposes because they ignore enzyme

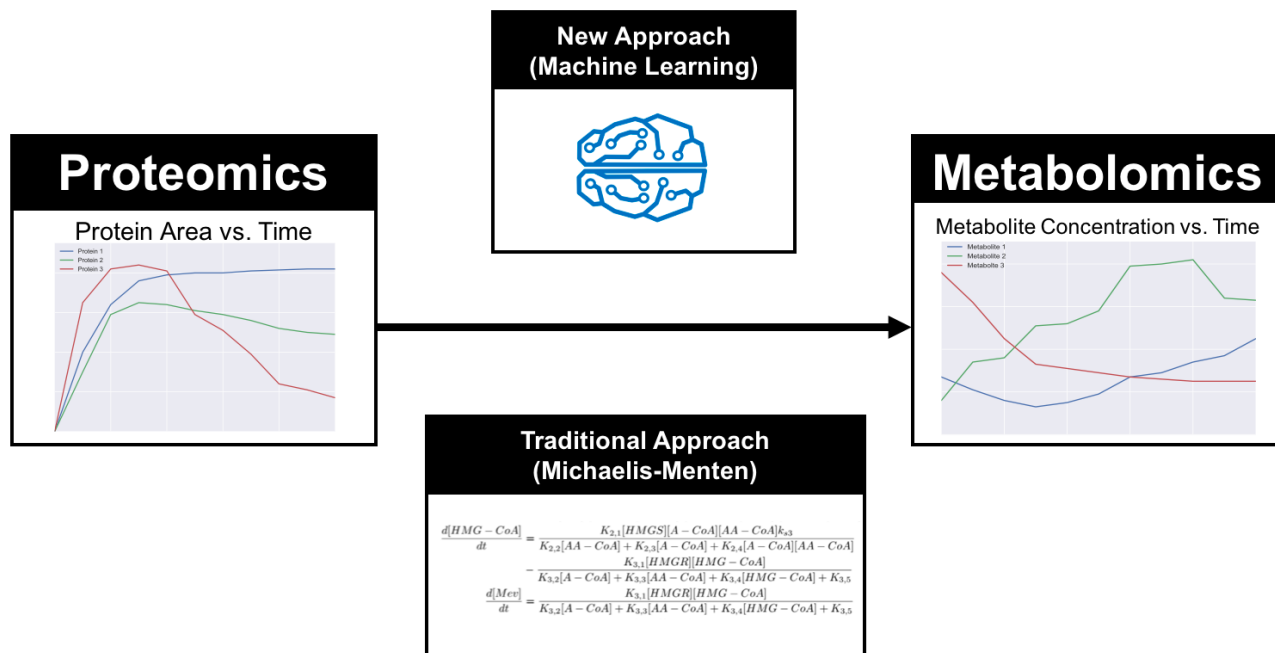


Figure 1. An alternative to traditional kinetic modeling by using machine learning. Our goal is to use time series proteomics data to predict time series metabolomics data (Fig. 2). The traditional approach involves using ordinary differential equations where the change in metabolites over time is given by Michaelis-Menten kinetics (Fig. 3 and 4). The alternative approach proposed here uses time series of proteomics and metabolomics data to feed machine learning algorithms in order to predict pathway dynamics (equation 1 and Fig. S1). While the machine learning approach necessitates more data, it can be automatically applied to any pathway or host, leverage systematically new data sets to improve accuracy and capture dynamic relationships which are unknown by the literature or have a different dynamic form than Michaelis-Menten kinetics.

39 kinetics and cannot accurately capture dynamic metabolic responses, nor offer a straightforward way to leverage ever more
40 abundant proteomics and metabolomics data for increased accuracy.

41 Kinetic models explicitly take into account enzyme kinetics and are able to predict metabolite concentrations as a function
42 of time from protein concentrations³¹. This type of prediction is useful to metabolic engineers in order to design pathways that
43 have the desired titers, rates and yields. Kinetic models rest on an explicit functional relationship connecting the rate of change
44 of a metabolite and the proteins and metabolites involved in the reaction (see Fig. 1): Michaelis-Menten kinetics^{32,33} is the
45 most common choice, but the fact is that the true mechanistic kinetic rate law for each specific reaction is unknown for most
46 enzymes³⁴ (alternatives include generalized mass action³⁵, lin-log kinetics^{36,37} or power-law models³⁸). However, there is a
47 lack of reliable data for the enzyme activity and substrate affinity parameters used in these models: in-vitro characterization may
48 not be extrapolatable to in vivo conditions, and the effect of activators and inhibitors are typically unknown. Approaches such as
49 ensemble modeling³⁹⁻⁴³ tackle the parsity of known kinetic constants by producing an ensemble of models displaying different
50 combinations of randomly chosen kinetic parameters and selecting only those models that match known experimental data, or
51 by optimizing the selection of these parameters through genetic algorithms^{44,45}. In a similar fashion, ORACLE^{46,47} produces
52 populations of models which are consistent with reaction stoichiometry, thermodynamics and available concentration and
53 fluxomic data. By design, these approaches are able to match measured final production levels and flux data, and the predictions
54 have been shown to improve as the model approaches genome-scale coverage⁴⁵. However, a significant problem remains in that
55 essential mechanisms are only sparsely known: allosteric regulation, for example, is known to be critical in order to determine
56 fluxes^{48,49}, and yet a comprehensive map of this regulatory mechanism is unavailable. Post-translational modifications of
57 proteins are also known to markedly affect catalytic activity⁵⁰, but are still largely unmapped. Pathway channeling, too,
58 significantly affects catalytic rates but the degree to which this phenomenon occurs in metabolism has only begun to be
59 explored⁵¹⁻⁵³. These and other gaps in our knowledge of mechanisms will require significant time and effort to be cleared.
60 Given the urgent need of predictive capabilities by the emerging biotech industry, it may be useful to consider a different

61 approach while this knowledge is gathered.

62 Here we propose an alternative to traditional kinetic modeling involving a machine learning approach (Figs. 1 and 2),
63 in which the function that determines the rate of change for each metabolite from protein and metabolite concentrations is
64 directly learned from training data (equation 1 and Fig. S1), without presuming any specific relationship. Machine learning
65 has shown remarkable success in well bounded problems where a mechanistic model is impossible or difficult to develop: e.g.
66 artificial vision for driverless cars⁵⁴, automated playing of the Go game⁵⁵, automated language translation⁵⁶ or private trait
67 prediction from digital records of human behavior⁵⁷ with direct impact on national elections⁵⁸. In biology, these methods have
68 recently been successfully applied to challenging problems such as predicting DNA and RNA protein binding sequences⁵⁹, skin
69 cancer diagnosis⁶⁰, single nucleotide polymorphism (SNP) and small indel variant calling⁶¹, and tumor detection in breast
70 histopathology⁶².

71 This alternative, machine-learning based, approach provides a faster development of predictive pathway dynamics models
72 since all required knowledge (regulation, host effects... etc) is inferred from experimental data, instead of arduously gathered
73 and introduced by domain experts (see supplementary material for an example). In this way, the method provides a general
74 approach, valid even if the host is poorly understood and there is little information on the heterologous pathway, and provides a
75 systematic way to increase prediction accuracy as more data is added. This method obtains better predictions than the traditional
76 Michaelis Menten approach for the limonene and isopentenol producing pathways studied here (Fig. 3) using only two times
77 series (strains), and is shown to significantly improve its prediction performance as more time series are added. The new
78 method is accurate enough to drive bioengineering efforts: we show it is able to predict the relative production ranking for
79 several designs, given enough data. This approach is a specific solution to the more general type of problem of determining
80 dynamics from observed data (system identification)^{63–65}, a problem generally recognized as hard. We believe this approach is
81 scalable to genome-scale models, or generally applicable to other types of data (e.g. transcriptomics) or dynamic systems (e.g.
82 microbiome dynamics).

83 **Mathematical Problem Formulation**

84 Here, we describe the problem and its solution in succinct mathematical terms. Let us assume we are given q sets of time
85 series metabolite $\tilde{\mathbf{m}}^i[t] \in \mathbb{R}^n$ (Fig. S2) and protein $\tilde{\mathbf{p}}^i[t] \in \mathbb{R}^\ell$ observations at times $\mathbf{T} = [t_1, t_2, \dots, t_s] \in \mathbb{R}_+^s$. The superscript
86 $i \in \{1, \dots, q\}$ indicates the time series index (strain), and $\tilde{\mathbf{m}}[t] = [\tilde{m}_1[t], \dots, \tilde{m}_n[t]]^T$ and $\tilde{\mathbf{p}}[t] = [\tilde{p}_1[t], \dots, \tilde{p}_\ell[t]]^T$ are vectors of
87 measurements at time t containing concentrations for the n metabolites and ℓ proteins considered in the model. We require the
88 number of observation time points to be dense enough to capture the dynamic behavior of the system.

89 We also assume that the underlying continuous dynamics of the system which generates these time series observations can
90 be described by coupled nonlinear ordinary differential equations of the general type used for kinetic modeling:

$$\dot{\mathbf{m}}(t) = f(\mathbf{m}(t), \mathbf{p}(t)), \quad (1)$$

91 where \mathbf{m} and \mathbf{p} are vectors that denote the metabolite and protein concentrations, as explained above. The function
92 $f: \mathbb{R}^{n+\ell} \rightarrow \mathbb{R}^n$ encloses all the information on the system dynamics. Deriving these dynamics from the time series data will be
93 formulated as a supervised learning problem where the function f is learned through machine learning methods which predict
94 the relationship between metabolomics and proteomics concentrations (input features, see Fig. S1) and the metabolite time
95 derivative $\dot{\mathbf{m}}(t)$ (output). In order to provide the training data set for this problem, the metabolite time derivative $\dot{\mathbf{m}}$ is obtained
96 from the times series data $\tilde{\mathbf{m}}(t)$, as shown in Fig. S2.

97 In order to parametrize the machine learning algorithm, the following optimization problem is solved (through scikit-learn,
98 see materials and methods):

Problem 1 (Supervised Learning of Metabolic Dynamics). *Find a function f which satisfies:*

$$\arg \min_f \sum_{i=1}^q \sum_{t \in \mathbf{T}} \|f(\tilde{\mathbf{m}}^i[t], \tilde{\mathbf{p}}^i[t]) - \dot{\tilde{\mathbf{m}}}^i(t)\|^2 \quad (2)$$

100 Solving Problem 1 is equivalent to finding the metabolic dynamics which best describe the time series data provided. Once the
101 dynamics are learned we can then predict the behavior of the metabolic pathway by solving an initial value problem (equations 3
102 and 4).

103 **Results and Discussion**

104 We used the supervised learning method described above (Figs 1 and 2, Eqns 1, 2, 3 and 4) to predict pathway dynamics (i.e.
105 metabolite concentrations as a function of time) from protein concentration data for two pathways of relevance to metabolic

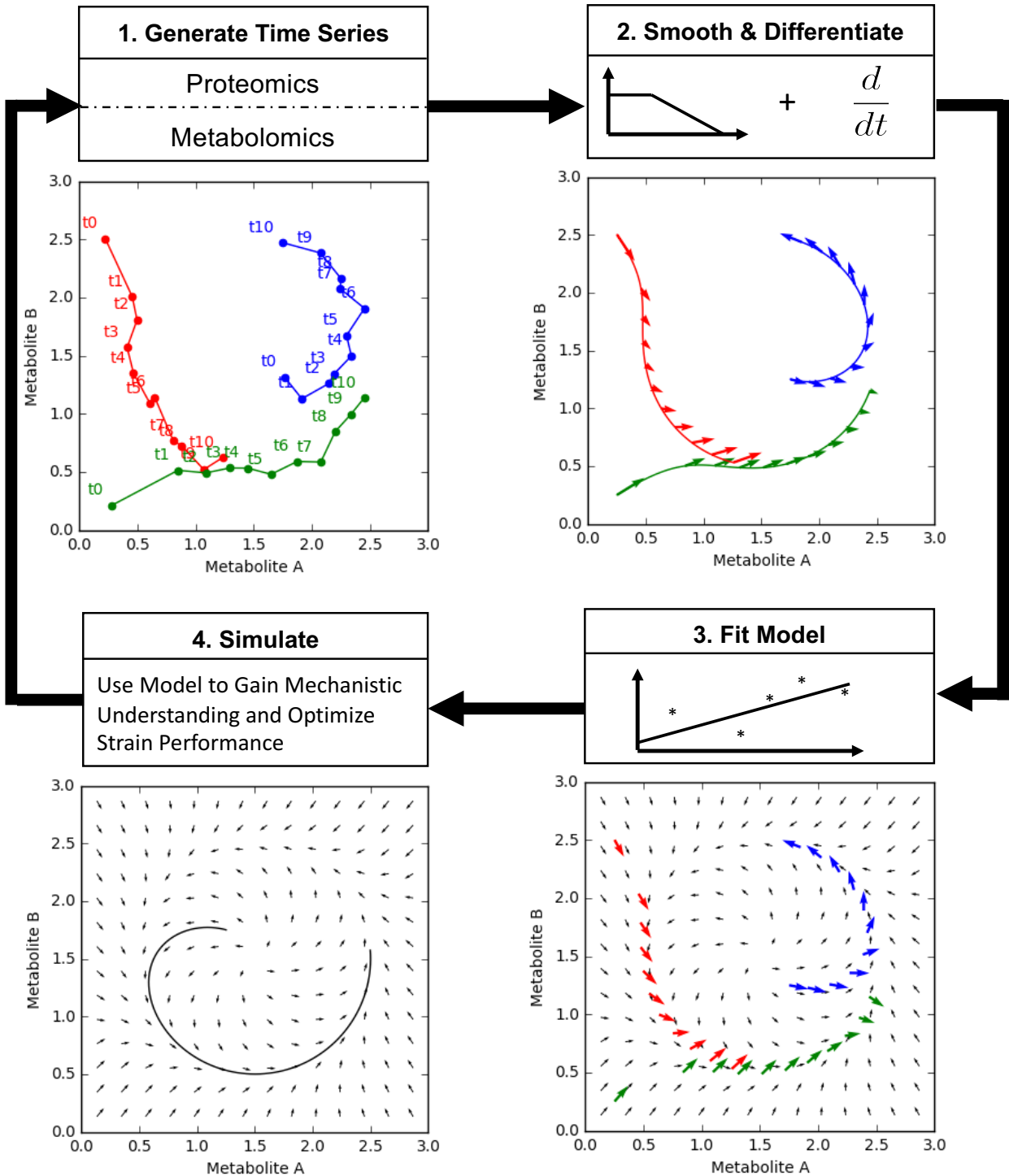


Figure 2. Cycle for learning metabolic system dynamics from time series proteomics and metabolomics data. (1) Experimentally, time series proteomics and metabolomics data are acquired for several strains of interest (represented by different colors). These data are represented in a metabolomics phase space, with an axis corresponding to each measured metabolite. (2) The time series data traces are smoothed and differentiated (Fig. S2). The derivatives provide the training data to derive the relationship between metabolomics and proteomics data and the metabolite change (Fig. S1, equation 1). (3) The state derivative pairs are fed into a supervised machine learning algorithm. The machine learning algorithm learns and generalizes the system dynamics from the examples provided by each strain. (4) The model can then be used to simulate virtual strains and explore the metabolic space looking for mechanistic insight or commercially valuable designs. This process can then be repeated using the model to create new strains, which will further improve the accuracy of the dynamic model in the next round.

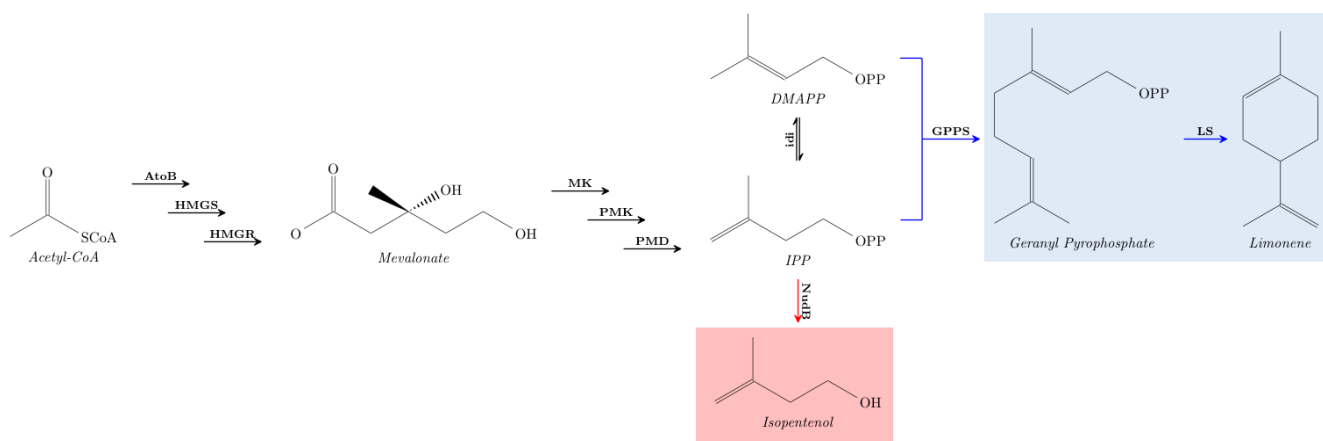


Figure 3. The new method was tested on the limonene & isopentenol metabolic pathways. The limonene (blue) and isopentenol (red) producing pathways are variants of the mevalonate pathway. Time series proteomics and metabolomics data are used to learn the dynamics of both the isopentenol and limonene producing strains. Additionally, a kinetic model is created and compared to the machine learning approach for the more complex limonene production pathway (Fig. 4). This pathway model is also used to generate simulated data to further evaluate the scaling properties of the proposed machine learning approach. See the original data set⁶⁶ for enzyme and metabolite acronyms.

106 engineering and synthetic biology: a limonene producing pathway and an isopentenol producing pathway (Fig. 3). For each
 107 pathway we used experimental times series data obtained from the low and high biofuel producing strains as training data sets,
 108 in order to predict the dynamics for the medium producing strains⁶⁶. Because of the paucity of dense multiomics time series
 109 data sets, we used simulated data sets (Fig. 4) to study the algorithm's performance as more training data sets (strains) were
 110 added.

111 Qualitative Predictions of Limonene & Isopentenol Pathway Dynamics can be Obtained with Two Time 112 Series Observations

113 Surprisingly, just two time-series (strains) were enough to train the algorithm to produce acceptable predictions for most
 114 metabolites. While the predictions of derivatives from proteomics and metabolomics were quite accurate (aggregate Pearson
 115 R value of 0.973), any small error in these predictions compounds quickly when solving the initial value problem given by
 116 equations 3 and 4. The reason is that predictions for a given time point depend on the accuracy of all previous time points. In
 117 spite of these hurdles, the method produced respectable qualitative and quantitative predictions of metabolite concentrations
 118 for a strain it has never seen before (Figs. 5 and 6). For some metabolites (33%), the predictions were quantitatively close
 119 to the measured profile: Acetyl-CoA (83.4 % error, Fig. 5a) and Isopentenol (43.7 % error, Fig. 5f) for the isopentenol
 120 producing pathway; Acetyl-CoA (128.2 % error, Fig. 6a), HMG-CoA (83.9 % error, Fig. 6b) and Limonene (82.9 % error,
 121 Fig. 6f) for the limonene producing pathway. For most metabolites (42%), the predictions were off by a scale factor, but
 122 they were able to qualitatively reproduce the metabolite behavior. For example, for Mevalonate in the isopentenol producing
 123 pathway (Fig. 5c) and mevalonate in the limonene producing pathway (Fig. 6c) the predictions reproduce the initial increase
 124 of metabolite concentration followed by a saturation. For IPP/DMAPP (Fig. 5e) or Mevalonate Phosphate (Fig. 5d) in the
 125 isopentenol pathway, the prediction reproduces qualitatively the concentration increase, followed by a peak and a decrease.
 126 The prediction of even just this type of qualitative behavior is useful to metabolic engineers in order to obtain an intuitive
 127 understanding of the pathway dynamics and design better versions of it. By simulating several scenarios the metabolic engineer
 128 can extract qualitative knowledge (e.g. metabolite x seems toxic, or protein y seems regulated by metabolite x) that can lead to
 129 testable hypotheses. Finally, in a minority of cases (25%), the predictions are wrong both quantitatively and qualitatively: e.g.
 130 HMG-CoA for the isopentenol producing pathway (Fig. 5b), Mevalonate Phosphate (Fig. 6d) and IPP/DMAPP (Fig. 6e) for the
 131 limonene producing pathway. Interestingly, the predictions for both final products (limonene and isopentenol) fell in the group
 132 of quantitatively accurate predictions. This is important because, for the purpose of guiding metabolic engineering, it is the
 133 final product predictions that are relevant.

134 The machine learning approach outperforms a handcrafted kinetic model of the limonene pathway (Fig. 6). A realistic
 135 kinetic model of this pathway was built and fit to the data, leaving all kinetic constants as free parameters (Figs. 3 and 4). The
 136 kinetic model notably fails to capture the qualitative dynamics for Acetyl-CoA, HMG-CoA, Mevalonate and IPP/DMAPP
 137 (Figs 6a,b,c and e). More quantitatively, the machine learning model produces an average 130% error (RMSE = 8.42) versus an

$$\begin{aligned}
\frac{d[A - CoA]}{dt} &= \frac{K_{1,1}[AtoB][A - CoA]}{K_{1,2} + K_{1,3}[A - CoA]} - \frac{K_{2,1}[HMGS][A - CoA][AA - CoA]k_{s3}}{K_{2,2}[AA - CoA] + K_{2,3}[A - CoA] + K_{2,4}[A - CoA][AA - CoA]} \\
\frac{d[AA - CoA]}{dt} &= \frac{K_{1,1}[AtoB][A - CoA]}{K_{1,2}K_{1,3}[A - CoA]} - \frac{K_{2,1}[HMGS][A - CoA][AA - CoA]k_{s3}}{K_{2,2}[AA - CoA] + K_{2,3}[A - CoA] + K_{2,4}[A - CoA][AA - CoA]} \\
\frac{d[HMGS - CoA]}{dt} &= \frac{K_{2,1}[HMGS][A - CoA][AA - CoA]k_{s3}}{K_{2,2}[AA - CoA] + K_{2,3}[A - CoA] + K_{2,4}[A - CoA][AA - CoA]} \\
&\quad - \frac{K_{3,1}[HMGR][HMGS - CoA]}{K_{3,2}[A - CoA] + K_{3,3}[AA - CoA] + K_{3,4}[HMGS - CoA] + K_{3,5}} \\
\frac{d[Mev]}{dt} &= \frac{K_{3,1}[HMGR][HMGS - CoA]}{K_{3,2}[A - CoA] + K_{3,3}[AA - CoA] + K_{3,4}[HMGS - CoA] + K_{3,5}} \\
&\quad - \frac{K_{4,1}[MK][Mev]}{K_{4,2}[GPP] + K_{4,3}[MevP] + K_{4,4}[Mev] + K_{4,5}} \\
\frac{d[MevP]}{dt} &= \frac{K_{4,1}[MK][Mev]}{K_{4,2}[GPP] + K_{4,3}[MevP] + K_{4,4}[Mev] + K_{4,5}} - \frac{K_{5,1}[PMK][MevP]}{K_{5,1} + [MevP]} \\
\frac{d[MevPP]}{dt} &= \frac{K_{5,1}[PMK][MevP]}{K_{5,1} + [MevP]} - \frac{K_{6,1}[PMD][MevPP]}{K_{6,2}[MevP] + K_{6,3}[Mev] + K_{6,4}[MevPP] + K_{6,5}} \\
\frac{d[IPP]}{dt} &= \frac{K_{6,1}[PMD][MevPP]}{K_{6,2}[MevP] + K_{6,3}[Mev] + K_{6,4}[MevPP] + K_{6,5}} - \frac{K_{7,1}[IDI][IPP]}{K_{7,2} + [IPP]} \\
&\quad - \frac{K_{8,1}[GPPS][IPP][DMAPP]}{K_{8,2} + K_{8,3}[IPP]K_{8,4}[DMAPP] + [IPP][DMAPP]} \\
\frac{d[DMAPP]}{dt} &= \frac{K_{7,1}[IDI][IPP]}{K_{7,2} + [IPP]} - \frac{K_{8,1}[GPPS][IPP][DMAPP]}{K_{8,2} + K_{8,3}[IPP] + K_{8,4}[DMAPP] + [IPP][DMAPP]} \\
\frac{d[GPP]}{dt} &= \frac{K_{8,1}[GPPS][IPP][DMAPP]}{K_{8,2} + K_{8,3}[IPP]K_{8,4}[DMAPP] + [IPP][DMAPP]} - \frac{K_{9,1}[LS][GPP]}{K_{9,2} + [GPP]} \\
\frac{d[Limonene]}{dt} &= \frac{K_{9,1}[LS][GPP]}{K_{9,2} + [GPP]}
\end{aligned}$$

Figure 4. Limonene pathway kinetic Michaelis Menten model. This kinetic model was compiled from sources in the BRENDA database along with guidance from Weaver *et al*⁶⁷. This system is composed of 10 nonlinear ordinary differential equations which describe the concentration for each metabolite in the pathway (see supplementary material for details). The dynamics of this model are rich and complex enough to pose a significant challenge to be predicted through machine learning. This model is used in this work to: 1) compare its predictions with machine learning predictions, and 2) generate simulated data sets to check scaling dependencies with the amount of time series used for training of machine learning algorithms. The method presented in this paper focuses on substituting these Michaelis Menten expressions by machine learning algorithms (see Fig. S1). Kinetic constants were left as free parameters when fitting experimental data in Fig. 6.

138 average 144% (RMSE = 10.04) for the kinetic model. Hence, even a machine learning model informed by the time series data
139 of just two strains is able to outperform the handcrafted kinetic model which required domain expertise and significant time
140 investment to construct. The machine learning approach, however, is more easily generalizable and it can be instantly reapplied
141 for a new pathway, host or product by feeding it the corresponding data. Once the predictions were made for the limonene
142 pathway, results for the isopentenol pathway can be obtained easily just by changing the time series data input. In contrast, in
143 order to make predictions for the isopentenol pathway a new kinetic model would have to be crafted. Kinetic models become
144 more difficult to construct as the size of the reaction network increases and as the knowledge of the relevant network decreases.
145 Additionally, all kinetic relationships must be known or inferred, whereas unknown relationships can be uncovered from data
146 using a machine learning approach. The machine learning approach only requires a sufficient amount of data to disentangle
147 these relationships. Determining how much data is a “sufficient amount” is the goal of the next section.

148 Interestingly, the model was able to perform well even though the training sets corresponded to pathways which differed in
149 more than just protein levels. This is important because the model is designed to take protein concentrations as input (Fig. 1) in
150 order to predict pathway dynamics, assuming the rest of pathway characteristics to remain the same. This use case covers a
151 wide range of metabolic engineering needs where e.g. promoters and ribosome binding sites (RBSs) are modified in order to

Prediction of Isopentenol Strain Dynamics

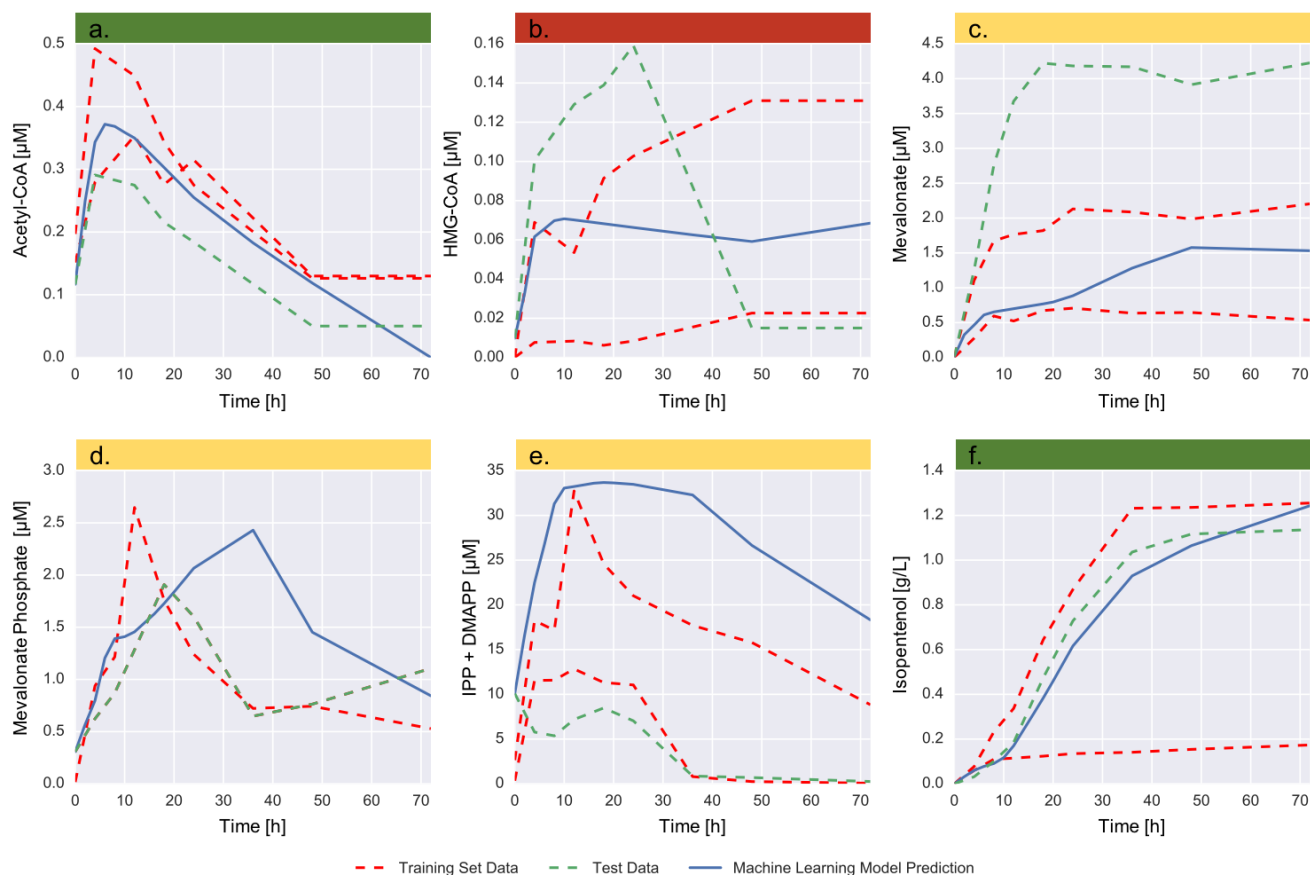


Figure 5. The machine learning method produces acceptable predictions of metabolite time series from proteomics data for the isopentenol producing *Escherichia coli* strain. The measured metabolomics and proteomics data⁶⁶ for the highest and lowest producing strains (training set data, red line) are used to train a model and learn the underlying dynamics (Fig. 2). The model is then tested by predicting the metabolite profiles (blue line) for a strain the model has never seen (medium producing strain, test data in green). A perfect prediction (blue line) would perfectly track the test data set (green line). Interestingly, reasonable qualitative agreement is achieved even with only two time series (strains) as training data. From a purely quantitative perspective, the average error is high: the total RMSE for the strain predictions is 40.34, which can be translated to 149.2% average error. However, for a couple of metabolites (green color band) the predictions quantitatively reproduce the measured data: Acetyl-CoA and Isopentenol (the final product, and most relevant for guiding bioengineering). For some metabolites (Mevalonate, Mevalonate Phosphate and IPP/DMAPP, yellow band), the model qualitatively reproduces the metabolite patterns, missing the scale factor. Only for HMG-CoA does the model fail to predict the metabolite concentration over time both quantitatively and qualitatively (red band).

152 affect the resulting protein concentrations. However, other typically used metabolic engineering strategies include changing a
 153 given enzyme in order to access faster or slower catalytic rates (i.e. k_{cat}). Even though this case was not explicitly contemplated,
 154 the model was able to provide good predictions (i.e. I3 was using a HMGR analogue from *Staphylococcus aureus* and I2
 155 uses a codon optimized HMGR, see strain description). We hypothesize that k_{cat} changes can be renormalized into (and be
 156 equivalent to) protein abundance changes. In order to fully address this type of engineering practice, this method may be
 157 expanded to include enzyme characteristics as input (besides the proteomics data): k_{cat} and K_M constants or even full kinetic
 158 characterization curves.

159 Increasing the Number of Strains Improves the Accuracy of Dynamic Predictions

160 We used simulated data to show that predictions improved markedly as more data sets are used for training. Simulated data sets
 161 have the advantage of providing unlimited samples to thoroughly test scaling behavior, and allow us to explore a wider variety
 162 of types of dynamics than experimentally accessible. Moreover, the dense multiomics time series data sets needed as training

Prediction of Limonene Strain Dynamics

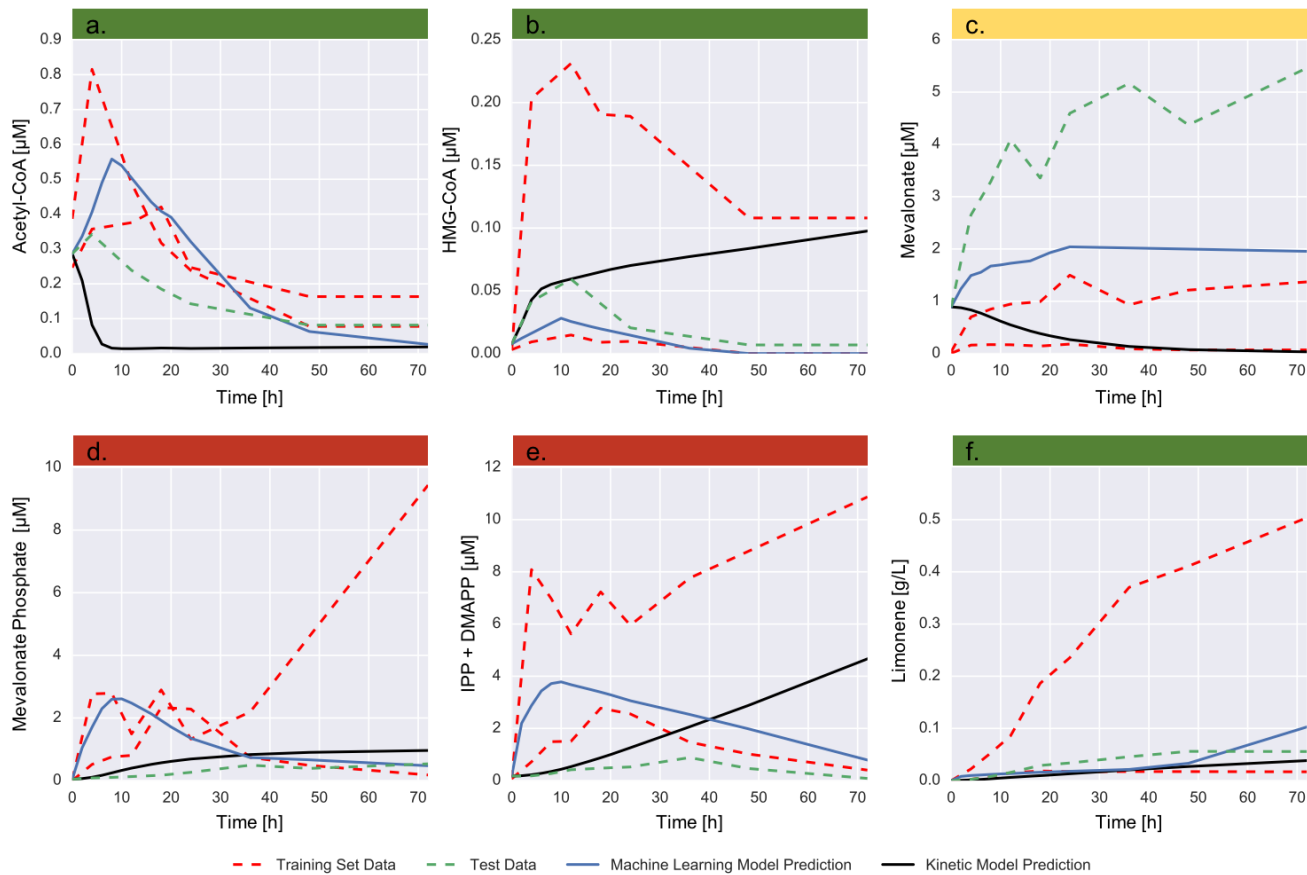


Figure 6. The machine learning method outperforms the handcrafted kinetic model for the limonene producing *E. coli* strain. The only metabolite for which the kinetic model (black line) provides a better fit than the machine learning method (blue line) is Mevalonate Phosphate, although both methods appear to track limonene (final product) production fairly well. The machine learning approach provides acceptable quantitative fits for Acetyl-CoA, HMG-CoA and Limonene (green band), a qualitative description of metabolite behaviour missing the scale factor for Mevalonate (yellow band), and fails quantitatively and qualitatively for Mevalonate Phosphate and IPP/DMAPP (red band). As in Fig. 5, the experimentally measured profiles correspond to high, low and medium producers of limonene. The training sets are the low and high producers (in red) and the model is used to predict the concentrations for the medium producing strain (in green). Kinetic constants for the handcrafted kinetic model in Fig.4 were left as free parameters when fitting the experimental data.

163 data are rare because they are very time consuming and expensive to produce. Since machine learning predictions generally
 164 improve as more data is used to train them, we expected our method to improve with the availability of more time series for
 165 training. We expected this improvement to be significant since initially only two time series (strains) were used for training,
 166 out of the three available for each product⁶⁶ (the other one was needed for testing). Hence, we used simulated data obtained
 167 from using the kinetic model developed for the limonene pathway (Figs. 3 and 4), in order to study: 1) how much predictions
 168 improve as more time series data sets are added and 2) how many time series are needed to guide pathway design effectively
 169 (next section). A pool of 10,000 sets of time series data with different protein profiles was created that shared the same kinetic
 170 constants. We fed the machine learning algorithm groups of 2, 10 and 100 times series randomly sampled from this pool in
 171 order to study how quickly the algorithm was able to recover the original simulated dynamics. In order to gauge the variability
 172 of the predictions (i.e., how predictions change as different training sets are used) as a function of training group size (2,10 or
 173 100), we repeated the predictions 10 times for each training group size.

174 The prediction error (RMSE, equation 6) decreased monotonically as a function of the number of time series (strains) used
 175 to train the algorithm in a nonlinear fashion (Fig. 7). Also, the standard deviation of the predictions significantly decreased with
 176 the number of training of data sets (Fig. 8). The standard deviation is an indication of the variability of pathway dynamics

177 predictions due to stochastic effects of the optimization algorithms (e.g. different seeds) and lack of extrapolability from a
178 reduced set of initial protein concentrations. Hence, a predictive model trained with 10 or 100 data sets produces much more
179 robust predictions than a model trained with 2 data sets. In fact, the high standard deviations observed for models trained on
180 only 2 data sets explain the prediction variability observed in the previous section due to stochastic effects. Interestingly, there
181 is a limited drop in error and standard deviation from 10 to 100 strains, with the decrease from 2 to 10 being the largest (Fig. 7).
182 This indicates that it is more productive to do 10 rounds of engineering collecting 10 time series data set than a single round
183 collecting 100 time series: in this way, 10 time series produce accurate enough predictions to pinpoint the desirable part of
184 proteomics of phase space, new strains can be engineered around that space so that new multiomics time series can be obtained
185 around the desirable phase space and optimize for prediction accuracy around that area of phase space. Doing this 10 times is
186 more accurate than a single prediction based on 100 time series that may not be close to the ultimately desirable proteomics
187 phase space. Furthermore, it indicates that the results from the previous section would have been much more reliable if only 8
188 time series more had been available for training.

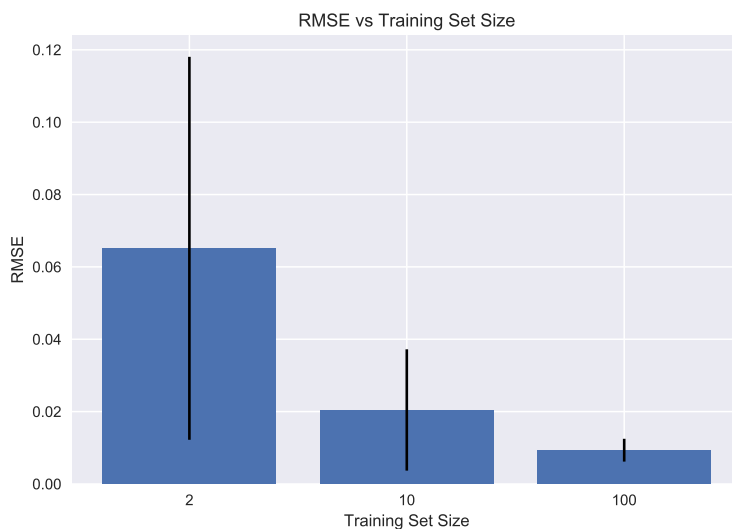


Figure 7. Prediction errors decrease markedly with increasing training set size. As the number of available proteomics and metabolomics times series data sets (strains) for training increases, the prediction error (RMSE, equation 6) decreases conspicuously. Moreover, the standard deviation of the predictions error (vertical bar) decreases notably as well. The change from 2 to 10 strains is more pronounced than the change from 10 to 100. This fact indicates that it is more productive to do 10 rounds of metabolic engineering collecting 10 time series data sets, than a single round collecting 100 time series.

189 **Model Predictions are Accurate Enough to Guide pathway Design and Produce Biological Insights**

190 The machine learning predictions do not need to be 100% quantitatively correct to accurately predict the relative ranking of
191 production for different strains. Being able to reliably predict which of several possible pathway designs will produce the
192 highest amount of product is very valuable in guiding bioengineering efforts and accelerating them in order to improve Titer,
193 Rate and Yield (TRY). These process characteristics are fundamental determinants of economic relevance⁶⁸.

194 The machine learning algorithm was able to reliably predict the relative production ranking for groups of three randomly
195 chosen strains (highest, lowest and medium producer, mimicking the available experimental data) chosen from the pool of
196 10,000 time series data sets mentioned above (Fig. 9, left panel). The success rate depended critically on the number of data
197 sets available for training: starting at 22% for only 2 strains up to 92% for 100 training sets. For 10 strains the success rate is ~
198 80%, which is reliable enough to practically guide metabolic engineer efforts to improve TRY. For models trained using 100
199 time series, the prediction errors were minimal (Fig. 9, right panel).

200 Biological insights can be generated by using the ML model to produce data in substitution of bench experiments. For
201 example, similarly to Principal Component Analysis of Proteomics (PCAP⁶⁹), we can use the ML simulations to determine
202 which proteins to over/underexpress, and for which base strain, in order to improve production (Fig. 10). Proteins LS, AtoB,
203 PMD and Idi are the most important drivers of production in the case of limonene: changing protein expression along the

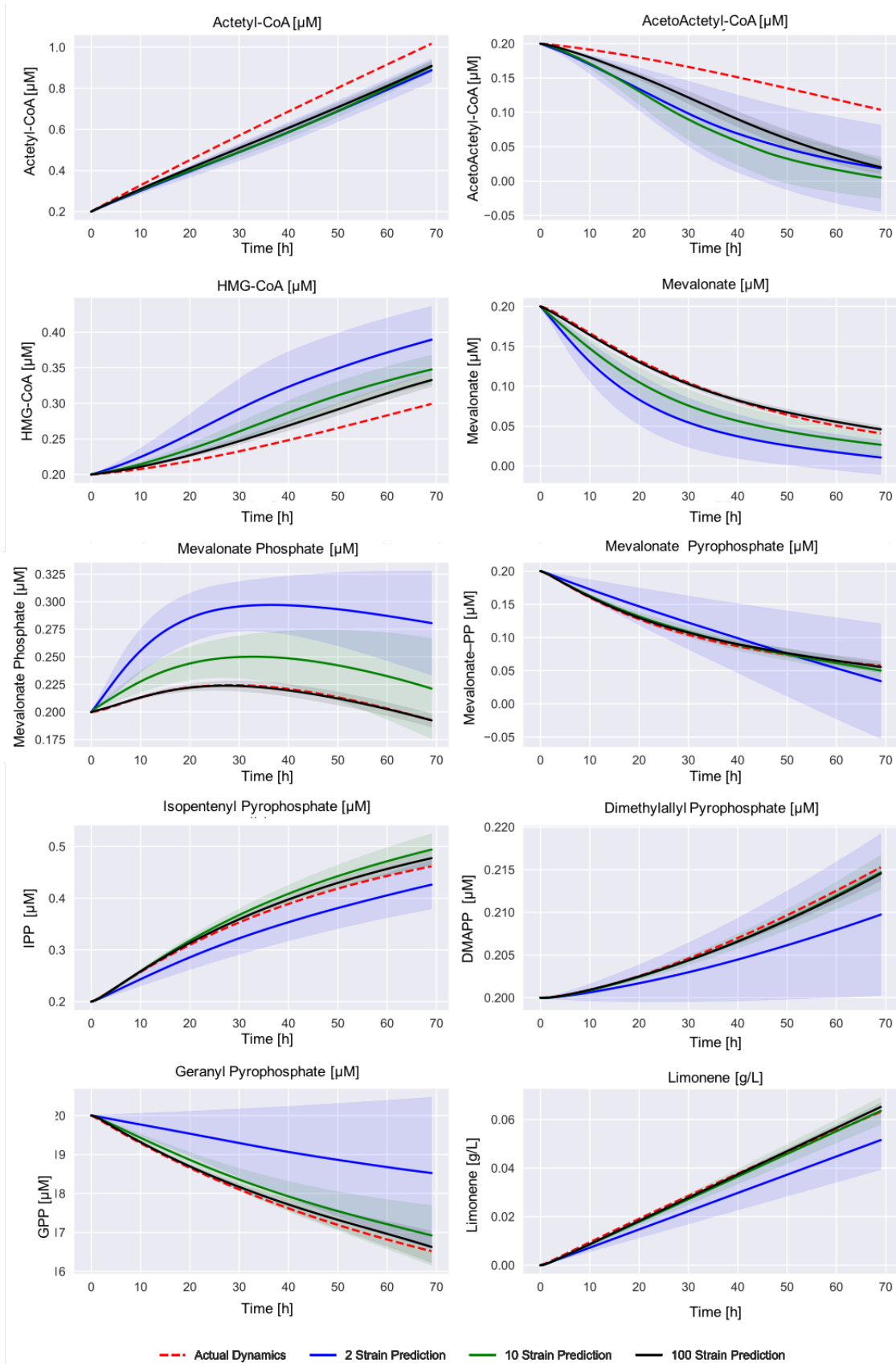


Figure 8. Predictions improve with more training data sets. The machine learning algorithm was used to predict kinetic models for varying sizes of training sets (2, 10 and 100 virtual strains in blue, red and black). 10 unique training sets were used for each size to show prediction variability (transparency) for each training set size. All models converge towards the actual dynamics with the 100 strain models in closest agreement. Standard deviations (shown by the transparency) also decrease markedly as the size of the training set increases.

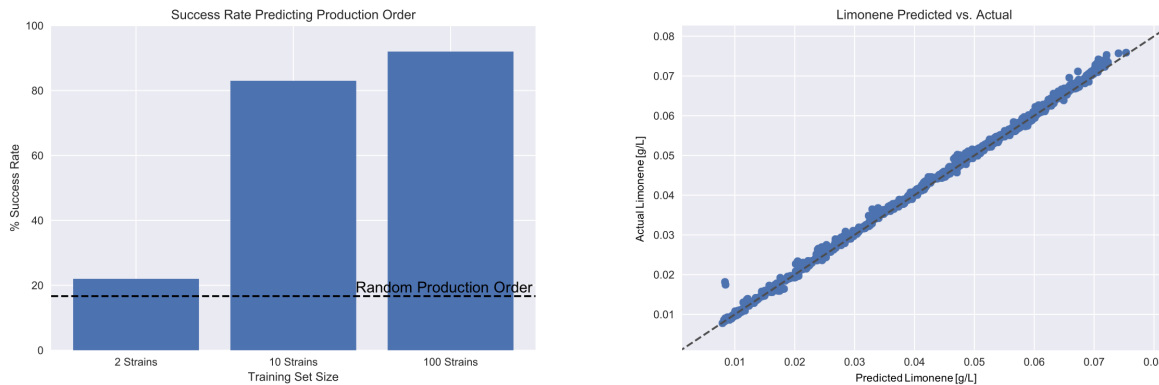


Figure 9. Success rate predicting production ranks increases with training set size. The left panel shows the success rate in predicting the relative production order (i.e. which strain produces most, which one produces least and which one is a medium producer) for groups of 3 time series (strains) randomly chosen from a pool of 10,000 strains, as a function of training data set size (strains). For 100 data sets, the failure rate to predict the top producer is < 10%. For 10 data sets the success rate is ~ 80%, which is reliable enough to guide engineering efforts. The horizontal line provides the rate of success (1/6) if order is chosen randomly. The right panel shows that prediction of limonene production is extremely accurate for the case of a training data set comprised of 100 time series (strains). These data shows that the machine learning model predictions are accurate enough to guide pathway design if enough training data is available.

204 principal component associated with them increases limonene creation (Fig. 10, left panel). Furthermore, this approach provides
 205 expected behavior for all metabolites in the pathway, providing hypotheses that can be tested experimentally (Fig. 10, right
 206 panel).

207 Data Constraints are Significant but Surmountable

208 Since the ML approach is purely data-based, data quantity and quality concerns are of paramount importance. Data quantity
 209 concerns involve both the availability of enough time series as well as time points sampled in each time series.

210 The training set used here⁶⁶ is one of the largest data sets characterizing a metabolically engineered pathway at regular
 211 time intervals through proteomics and metabolomics. There are no larger data sets that include: time series, several types of
 212 omics data, more than 7 time points, and several strains. For example: the *E. coli* multiomics database⁷⁰ has proteomics and
 213 metabolomics data for several strains, but no time series; Ma *et al*⁷¹ report proteomics and metabolomics data but only one time
 214 series with fewer time points (5 instead of 7); Yang *et al*⁷² only provide one time series and only one time point for proteomics;
 215 Doerfler *et al*⁷³ and Dyar *et al*⁷⁴ only provide time series metabolomics data; Patel *et al*⁷⁵ does not combine metabolomics
 216 and proteomics and data download was disabled at the time of testing; the DOE kbase⁷⁶ focuses on genomics and does not have
 217 any time series proteomics or metabolomics publicly available; and the Experiment Data Depot⁷⁷ does not have any studies
 218 surpassing this one in terms of data points and strains.

219 In order to get enough pairs of derivatives and proteomics and metabolomics data to train ML algorithms (Fig. S1), we have
 220 used data augmentation (filtering and interpolation, Figs 2 and S2), expanding the initial seven time points to two hundred
 221 by just assuming continuity in the multiomics data (a reasonable assumption in our experience). It would be desirable to
 222 have more time points available, so as to not to depend on these data augmentation techniques. However, data sets including
 223 more time points are nonexistent for physical, biological and economical reasons. Every time a sample is taken for -omics
 224 analysis, the volume in the culture flask diminishes and, if the total sampled volume is comparable to the total volume, it may
 225 significantly affect the strain physiology. Since taking excessive samples may affect measurements and these coupled omics
 226 analysis are expensive and require specialized personal, it is not surprising that the maximum amount of time points we have
 227 seen is ~ 7. Another reason more time points have not been typically collected is that experts in multiomics data collection
 228 consider this sampling rate to fully capture the physiology of strains based on previous experience^{78,79}. The fact that we are
 229 able to produce reasonable predictions for a third time series that the algorithm has never seen before (test strain) validates this
 230 and the multiomics data continuity assumption.

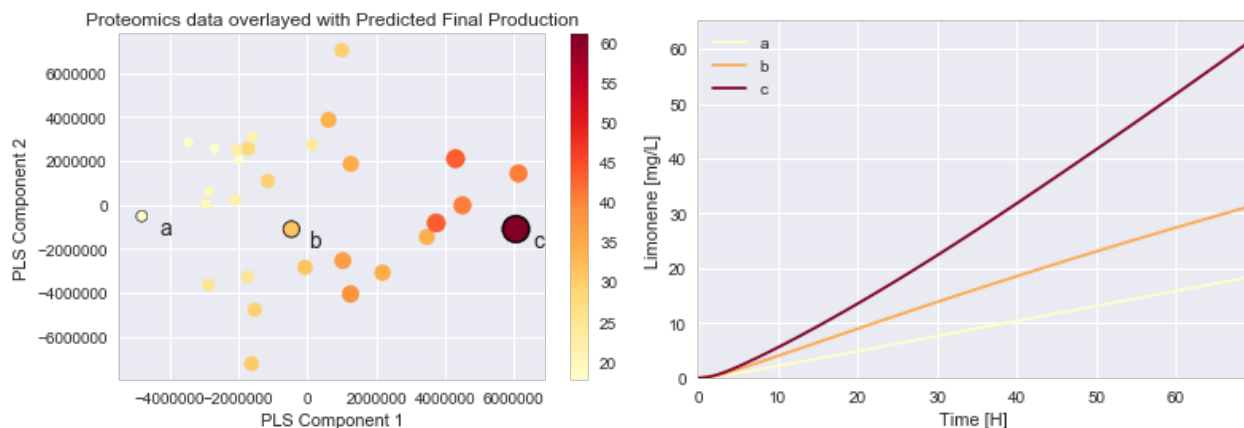


Figure 10. The ML approach can be used to produce biological insights. The left panel shows the final position in the proteomics phase space (similarly to the PCAP⁶⁹ approach) for fifty strains generated by the ML algorithm by learning from the Michaelis-Menten kinetic model (Fig.4) used as ground truth. Final limonene production is given by circle size and color. The PLS algorithm finds directions in the proteomics phase space that best align with increasing limonene production (component 1). Traveling in proteomics phase space along that direction (which involves overexpression of LS and underexpression of AtoB, PMD, and Idi, see Table S2) creates strains with higher limonene production. The ML approach not only produces biological insights to increase production but also predicts the expected concentration as a function of time for limonene and all other metabolites, generating hypotheses that can be experimentally tested (right panel).

Future Work

231

232 The application of machine learning to synthetic biology will hopefully open up new avenues of research as well as accelerate
 233 the adoption of modeling in bioengineering and beyond. This work is a first step demonstrating that a purely data-driven
 234 approach can fruitfully predict biological dynamics. There are plenty of possible ways to improve it.

235 An obvious first step involves adding other supervised learning techniques to improve predictions. The current approach
 236 uses TPOT to combine, through genetic algorithms, 11 different machine learning regressors and 18 different preprocessing
 237 (feature selection) algorithms. New supervised learning techniques can be added to this approach by adding them to the
 238 scikit-learn library⁸⁰. TPOT will automatically test them and use them if they provide more accurate predictions than the
 239 techniques used here. Among the most popular algorithms for ML are Deep Learning techniques based on neural networks.
 240 However, the small size of the available data sets for this study limited the use of machine learning techniques to classical
 241 methods. Modern Deep Learning (DL) techniques typically require orders of magnitude more data than was used in this study
 242 (~ 1000 strains as a starting point). While this amount of data is currently cost prohibitive, it is a worthy goal to move towards
 243 DL: these methods have demonstrated super human performance across a variety disciplines. These include, for example,
 244 image labeling tasks, in which humans have evolved proficiency. In domains where humans are less capable, such as the
 245 dynamical system characterization considered here, superhuman performance should be substantially easier to achieve. The
 246 payoff would involve radically improving engineering outcomes by making the predictability of complex biological systems
 247 proportional to the quantity of input data.

248 An often posed question is whether mechanistic insights can be inferred from ML approaches. While this is not trivial,
 249 there are a couple of possibilities for this inference: 1-) for any particular ML model that produces good fits, the most relevant
 250 features (i.e. protein X has the highest weight in determining Y molecule concentration) provides a prioritized list of putative
 251 mechanistically linked parts that can be further investigated. 2-) the ML model can be used as a surrogate for high-throughput
 252 experiments to derive mechanistic biological insights (Fig. 10). Another example of this last approach would involve studying
 253 toxicity by adding cell biomass (through optical density, OD) to the measurements and simulate for a variety of scenarios
 254 (protein inputs) the correlation between OD and all metabolites: a negative correlation would signal putative toxic metabolites.

255 It is instructive, however, to pause and reflect on the drive to find mechanisms. Mechanisms offer a causally related set of
 256 processes and parts that produce the observed phenomena. Understanding these processes, parts and causal relations produces a
 257 knowledge that can indeed be transferred to predict the behavior of different systems (pathways, strains, products.. etc) where
 258 the same mechanism is involved. However, biology has been particularly inefficient in making predictions of complex systems
 259 from known and tested mechanisms. If our final goal is to predict new biological systems, it may be more successful to look
 260 into ML techniques such as transfer learning⁸¹. These techniques tackle directly the challenge of predicting systems based on

261 data originated in related systems without the need to delve into mechanisms. Having said this, there is not doubt that the most
262 desirable outcome is a model that is both predictive and mechanistic, but if we are to do without one of these characteristics, the
263 mechanistic knowledge may be the least immediately useful for current bioengineering.

264 Infusing prior knowledge into the ML approach is a related possible future research avenue. Currently, our method does not
265 constrain the vector fields that are learned using any biological intuition. There are often biological facts known about these
266 dynamical systems that could be use to improve the performance of our method. Specifically, genome scale stoichiometric
267 constraints could provide guarantees that the resulting system dynamics conserve mass and conform to our prior knowledge
268 about the organism.

269 Since the procedure outlined here requires little prior biological knowledge, it is enticing to imagine extending this method
270 for use with different data inputs or other types of applications. An obvious extension is to use transcriptomics data as input.
271 Given the current exponential increase in sequencing capabilities, transcriptomics data is more amenable to high throughput
272 production than proteomics and metabolomics data. Our biological intuition says that transcriptomics data should be less
273 informative than proteomics, but it is surely interesting to explore whether that can be countered with more time series (and how
274 many). It would also be of interest to use the ML method to predict proteomics in addition to metabolomics time series. Another
275 logical proposition is to expand this method to encompass genome-scale multiomics data. We surmise that the extra predictive
276 capabilities of the machine learning with respect to the Michaelis-Menten approach proceed, in part, from indirectly accounting
277 for host metabolism effects through proxies (e.g. metabolites or proteins that are affected indirectly by host metabolism). Hence,
278 we expect more comprehensive metabolomics and proteomics (as well as transcriptomics) data sets to increase the method
279 predictive accuracy. A more intriguing and bold endeavor would be to apply this method to predict microbial community
280 dynamics using metaproteomics and metabolite concentration data as inputs. There is nothing in this approach that constrains it
281 to intracellular pathway prediction and microbiome research and industry have a definite need for increased predictive power⁸².
282 Finally, the incoming availability of dense multiomics data sets for human metabolism provides an alluring target^{83,84}.

283 Conclusion

284 We have demonstrated that it is possible to use a pure machine learning approach to qualitatively predict pathway dynamics.
285 This approach, using only two time series (strains) as training data, was able to outperform in predictive power a classical
286 Michaelis-Menten kinetic model. Unlike traditional kinetic modeling, we do not need to assume any particular interaction (e.g.
287 allosteric regulation), but we give full freedom to the system to implicitly choose the ones that best predict the experimental
288 data. Furthermore, we were able to produce predictions that, although not fully quantitatively accurate, are precise enough to
289 drive design decisions given enough data: production rankings can be predicted. The ability to predict the pathway dynamics is
290 of significant interest to metabolic engineers and synthetic biologists, since it allows for building an intuitive understanding of
291 the pathway that can produce testable hypotheses (yield increase, compound toxicity). This method is also an example of the
292 benefit of targeting the prediction of derivatives using machine learning in order to predict dynamic processes.

293 We have also shown that the machine learning approach improves markedly by using more time-series (strains) as training
294 sets, and used simulated data to estimate the number of time series required to guide engineering. Although the training set
295 used here⁶⁶ is one of the largest data sets characterizing a metabolically engineered pathway at regular time intervals through
296 proteomics and metabolomics, it is barely sufficient to train machine learning algorithms. Another limitation of this work
297 involves only being able to test the method with two pathways, which are the only ones for which dense time series multiomic
298 data sets are available. These limitations justify future efforts directed at methodic collection of large time series data sets
299 as enabled by multiomic pipelines⁸⁵⁻⁸⁸, since this method provides a systematic method to productively leverage those data.
300 Moreover, coupled with recent developments providing real-time metabolomics capabilities⁸⁹, this method opens the alluring
301 possibility of real-time prediction and control of biological pathways.

302 These results open the door to a data-centric approach to predicting metabolism that can greatly benefit the biotech and
303 synbio industries, much necessitated of predictive power in order to enable reliable production^{13,90}. This approach is agnostic as
304 to the pathway, host or product used, and can be systematically applied, as we have shown. Unlike previous approaches⁶⁶, it can
305 systematically leverage proteomics and metabolomics data in a fashion that increases accuracy as more data is available. Besides
306 being of immediate practical utility for bioengineering, this approach can be used as a first step in improving mechanistic
307 kinetic models by pinpointing the most relevant machine learning features for accurate predictions, that can then be followed
308 up by further experiments in order to obtain a mechanistic understanding of the reasons for their predictive power.

309 This work shows that, given sufficient data, the dynamics of complex coupled nonlinear systems relevant to metabolic
310 engineering can be systematically learned.

311 **Materials and Methods**

312 **Learning System Dynamics from Time Series Data**

313 The core of this method consists in using machine learning methods to predict the functional relationship between the metabolite
314 derivative and proteomics and metabolomics data, substituting the Michaelis Menten relationship (Eqn. 1, Figs. S1 and 4).
315 The first step involves creating a training set comprising sets of metabolomics and proteomics data and their corresponding
316 derivatives (Fig. S1). This entails computing the derivatives of the metabolite concentration time series data. Because the time
317 series data is subject to measurement noise, the derivatives must be carefully estimated. The second step involves finding the
318 best performing regression technique, among the many possibilities available⁸⁰. Finally, once the best prediction algorithm is
319 found and cross-validated, we can use it to predict metabolite concentrations given initial time points. The complete code to
320 implement these steps is provided in github (see below).

321 **Construction of the Training Data Set**

322 In order to train a machine learning model, a suitable training set must be created. We expect the trained machine learning
323 model to take in metabolite and protein concentrations at a particular point in time and return the derivative of the metabolite
324 concentrations at the same time point (Fig. S1). The observations provide us with the inputs to the model, $\tilde{\mathbf{m}}^i[t]$ and $\tilde{\mathbf{p}}^i[t]$. In
325 order to have examples of correct outputs for supervised learning we have to estimate the derivatives of the metabolite time
326 series data, $\dot{\mathbf{m}}^i(t)$ (Fig. S2).

327 Naively computing the derivative of a noisy signal will amplify the noise and make the result unusable. Derivatives of noisy
328 signals, like those obtained from experiment, require extra effort to estimate. In order to estimate the time derivatives on time
329 series of real data obtained from Brunk *et al*⁶⁶ accurately, we apply a Savitzky–Golay⁹¹ filter to the noisy time series data
330 to find a smooth estimate of the data (Fig. S2). This smooth function estimate can then be used to compute a more accurate
331 estimate of the derivative. We compute the derivative estimate of the signal using a central difference scheme from the filtered
332 experimental data. Specifically, the Savitzky-Golay filter is used with a filter window of 7 and a polynomial order of 2. The
333 derivative estimate, $\dot{\mathbf{m}}^i(t)$, is computed for all time points in T and time series i . This results in a training example associated
334 with each time point in every time series.

335 This work assumes that all relevant metabolites are measured and that the system has no unmeasured memory states. In
336 other words, the present set of metabolite and protein measurements completely determines the metabolite derivatives at the
337 next time instant. If this assumption does not hold practically, a limited time history of proteins and metabolites can be used to
338 predict the derivative at the next time instant. We observe that, for the specific pathways used in this paper, this assumption
339 produces good predictions.

340 **Model Selection**

341 The model selection process used a meta-learning package in python called TPOT⁹². Once the training data set is established,
342 a machine learning model must be selected learn the relationship between input and outputs (Fig. S1). TPOT uses genetic
343 algorithms to find a model with the best cross validated performance on the training set. Cross validation techniques are used to
344 score an initial set of models. The best performing models are mated to form a new population of models to test. This process
345 is repeated for a fixed number of generations and the best performing model is returned to the user.⁹³ If desired, the search
346 space for model selection can be specified before execution of the TPOT regressor search. This might be done to prune models
347 that require long training times or to select only models that have desirable properties for the problem under consideration.
348 Specifically, we used TPOT to select the best pipelines it can find from the scikit-learn library⁸⁰ combining eleven different
349 regressors and eighteen different preprocessing algorithms. This model selection process is done independently for each
350 metabolite (Table S1). After TPOT determines the optimal models associated with each metabolite, they are trained on the data
351 set of interest and are ready for use to solve equations 3 and 4. Models with the lowest 10 fold cross validated prediction root
352 mean squared error were selected. In this way the best validated models were selected for use.

353 After automated model selection via TPOT, we manually evaluated each model based on its accuracy in predicting metabolite
354 derivatives given protein and metabolite concentration at a given time point (Fig. S1). Each data set used for model fitting was
355 split into training and test sets 10 times using the shuffle split methodology implemented in scikit-learn⁸⁰. After the model was
356 fit, predictions on both the training and test sets were computed for each metabolite model and their predictive ability quantified
357 through a Pearson R coefficient (e.g. Fig. S3).

358 **Using the Model**

359 Once the models are trained, we can use them to predict metabolite concentrations by solving the following initial value
360 problem using the same function f that was learned in equations 1 and 2 :

$$\dot{\mathbf{m}} = f(\mathbf{m}, \tilde{\mathbf{p}}), \quad (3)$$

$$\mathbf{m}(t_0) = \tilde{\mathbf{m}}(t_0). \quad (4)$$

361 This problem is solved by integrating the system forward in time numerically. As a general purpose numerical integrator, we
362 used a Runga Kutta 45 implementation.

363 **Data Set Curation and Synthesis**

364 Two different data sets were used in this work. The first is an experimental data set curated from a previous publication⁶⁶,
365 comprising 3 proteomic and metabolomics time series (strains) from an isopentenol producing *E. coli* and 3 time series (strains)
366 from limonene producing *E. coli*. The second data set involves computationally simulated data from a kinetic model of the
367 limonene pathway, which is used to test how the method performance scales with the number of time series used.

368 **Description of a Real Time Series Multiomics Data Set**

369 Proteomics and metabolomics data for two different heterologous pathways engineered into *E. coli* were available from Brunk
370 *et al*⁶⁶. There are three (high, medium, and low production) variants for strains which produce isopentenol and limonene
371 respectively. All strains were derived from *E. coli* DH1. The low and high producing strain for each pathway were used to
372 predict the medium production strain dynamics by solving equations 3 and 4.

373 The isopentenol producing strains (I1, I2 & I3) were engineered to contain all of the proteins required to produce isopentenol
374 from acetyl-CoA as (Fig 3). I1 is the unoptimized strain containing the naive variants of each protein in the pathway. I2 differs
375 from the base strain I1 in that it contains a codon optimized HMGR enzyme along with the positions of PMK and MK swapped
376 on its operon. I3 uses an HMGR homologue from *Staphylococcus aureus*. Limonene producing strains (L1, L2, & L3) produce
377 limonene from acetyl-CoA (Fig 3). L1 is the unoptimized strain with the naively chosen variants for each protein in the pathway.
378 It is a two plasmid system where the lower and upper parts of the pathway are split between both constructs. L2 is a DH1
379 variant which contains the entire limonene pathway on a single plasmid. L3 is another two plasmid strain where the entire
380 pathway is present on the first plasmid, and the terpene synthases are on a second plasmid for increased expression. Starting at
381 induction, each strain had measurements taken at 7 time points during fermentation over 72 hours. At each time point pathway
382 metabolite measurements and pathway protein measurements were collected. Further details on these strains and experimental
383 design can be seen in the original publication⁶⁶.

384 **Data Augmentation through Filtering and Interpolation**

385 In the training set each time series contains 7 data points. These are too sparse to formulate accurate models. To overcome this
386 a data augmentation scheme is employed where 7 time points from the original data are expanded into 200 for each strain. This
387 is done by smoothing the data with a Savitzky-Golay filter and interpolating over the filtered curve (Figs 2 and S2). When
388 predicting the dynamics of a medium production strain from a high and low producing strain, we performed model selection by
389 scoring each model using 10 fold cross validation and a pearson R^2 metric on two data augmented training strains.

390 **Development of Realistic Kinetic Models**

391 In order to study the scaling of performance as more training sets were added, a realistic and dynamically complex model of
392 the mevalonate pathway was developed from known interactions extracted from the literature (Figs. 3 and 4). The dynamic
393 model is implemented with Michaelis-Menten like kinetics and is a 10 state coupled nonlinear system. Complete details of this
394 kinetic model are available in the supplementary material. The objective was to create a realistic model, relevant to metabolic
395 engineering, for which learning the system dynamics is a nontrivial task on par with the difficulty of learning real system
396 dynamics from experimental data.

397 **Generation of a Simulated Data Set**

398 The kinetic model described above was used to create a set of virtual data time series (strains). The kinetic model coefficients
399 were chosen to be close to values reported in the literature while maintaining a non-trivial dynamic behavior.

400 A virtual strain is created by first generating a pathway proteomic time series. This is done by randomly choosing three
401 coefficients for each protein (k_f, k_m, k_l) which specify a leaky hill function. The hill function was chosen because it models the
402 dynamics of protein expression from RNA accurately. This leaky hill function specifies the protein measurements for each time
403 point and is defined in the equation below:

$$404 \quad \tilde{p}(t) = \frac{k_f t}{k_m + t} + k_l. \quad (5)$$

405 Once all protein time series are specified, they are used in conjunction with the kinetic coefficients to solve the initial
406 value problem in equations 3 and 4 in order to determine the time series of metabolite concentrations. The resulting data set is
407 a collection of time series measurements of different strain proteomics and metabolomics. All strains use the same kinetic
408 parameters and differential equations to generate the metabolomics measurements. The code used to generate this data can be
found in the github repository, as well.

409 Fitting the Michaelis Menten kinetic model

410 To compare the hand crafted kinetic model with the data centric machine learning methodology the parameters of the kinetic
411 model were fit to strain data. To find the best fit we used a differential evolution algorithm implemented in scipy. This global
412 optimizer was chosen because its convergence is independent of the initial population choice and it tends to need less parameter
413 tuning than other methods. All kinetic parameters were constrained to be between 10^{-12} and 10^9 . This large range of acceptable
414 parameter values allowed for maximum flexibility of the kinetic model to describe the data.

415 Evaluation of Model Performance for Time Series

Dynamical prediction was tested on a held back strain that the model did not use in training. When using the experimental data
sets⁶⁶, the medium titer strains were held back for testing. When using simulated data, a random strain from the data set was
selected. For each time series, agreement between predictions and test data was assessed by calculating the root mean squared
error (RMSE) of the predicted trajectories:

$$RMSE = \sqrt{\frac{1}{n} \sum_{j=1}^n \int_{t_0}^{t_f} (\bar{m}_j(t) - \tilde{m}_j(t))^2 dt}, \quad (6)$$

416 where $\bar{m}_j(t)$ is the linear interpolation of the actual metabolite concentration of metabolite j at time t (Fig. S2), and $\tilde{m}_j(t)$ is the
417 prediction obtained from solving equations 3 and 4.

418 Biological Insight Analysis

419 In order to showcase how biological insights can be derived (Fig. 10), we trained the ML model using 50 proteomics and
420 metabolomics time series, using the Michaelis Menten kinetic model as ground truth. Another 50 proteomics time series were
421 held back as a test data set. Each metabolite time series was predicted using the machine learning model and the associated
422 proteomic time series. The final time point proteomics and final production were collected for each predicted strain. The final
423 time point proteomics data was plotted in two dimensions with a basis selected by performing a partial least squares [PLS]
424 regression between the proteomics and final production data. These first basis vector from a PLS regression is the direction that
425 explains the most covariance between the proteomics data and production data. The PLS regression was implemented by and
426 used from scikit-learn.

427 Supporting Information

428 **S1 Appendix. Details of kinetic model for the limonene pathway.** This file contains the details of the derivation and
429 sources for the model development for the limonene pathway model used in this paper, as well as supplementary figures S1, S2
430 and S3.

431 Code availability

432 Code is available at the following code repository: <https://github.com/JBEI/KineticLearning>.

433 Data availability

434 All data was obtained from Brunk *et al*⁶⁶, and is also available at the code repository: <https://github.com/JBEI/KineticLearning>.

435 Acknowledgments

436 This work was part of the DOE Agile BioFoundry (<http://agilebiofoundry.org>) and Joint BioEnergy Institute (<http://www.jbei.org>)
437 supported by the U.S. Department of Energy, Energy Efficiency and Renewable Energy, Bioenergy Technologies Office, through
438 contract DE-AC02-05CH11231 between Lawrence Berkeley National Laboratory and the U. S. Department of Energy. The
439 United States Government retains and the publisher, by accepting the article for publication, acknowledges that the United
440 States Government retains a nonexclusive, paid-up, irrevocable, worldwide license to publish or reproduce the published
441 form of this manuscript, or allow others to do so, for United States Government purposes. The Department of Energy
442 will provide public access to these results of federally sponsored research in accordance with the DOE Public Access Plan
443 (<http://energy.gov/downloads/doe-public-access-plan>). HGM was also supported by the Basque Government through the BERC
444 2014-2017 program and by Spanish Ministry of Economy and Competitiveness MINECO: BCAM Severo Ochoa excellence
445 accreditation SEV-2013-0323.

446 The authors would like to acknowledge Marcella Gomez, Jorge Alonso Gutierrez and Kevin George for valuable discussions,
447 as well as suggestions to improve the quality of this work.

448 **Competing interests**

449 The authors declare no competing financial interests.

450 **Author contributions**

451 Z.C. and H.G.M. conceived the original idea. Z.C. did all computational and mathematical work. Z.C. and H.G.M wrote the
452 paper.

453 **References**

- 454 1. WATSON, J. & CRICK, F. Genetical implications of the structure of deoxyribonucleic acid. *Nature* **171**, 964–7 (1953).
- 455 2. Russo, E. Special Report: The birth of biotechnology. *Nature* **421**, 456–457 (2003). URL <https://doi.org/10.1038%2Fnpj6921-456a>.
- 456 3. Lee, J. *et al.* Systems metabolic engineering of microorganisms for natural and non-natural chemicals. *Nat Chem Biol* **8**,
457 536–46 (2012).
- 458 4. Beller, H., Lee, T. & Katz, L. Natural products as biofuels and bio-based chemicals: fatty acids and isoprenoids. *Nat Prod*
459 *Rep* **32**, 1508–26 (2015).
- 460 5. Fortman, J. *et al.* Biofuel alternatives to ethanol: pumping the microbial well. *Trends Biotechnol* **26**, 375–81 (2008).
- 461 6. Chubukov, V., Mukhopadhyay, A., Petzold, C. J., Keasling, J. D. & Martín, H. G. Synthetic and systems biology
462 for microbial production of commodity chemicals. *npj Systems Biology and Applications* **2** (2016). URL <https://doi.org/10.1038%2Fnpjjsba.2016.9>.
- 463 7. Lienert, F., Lohmueller, J., Garg, A. & Silver, P. Synthetic biology in mammalian cells: next generation research tools and
464 therapeutics. *Nat Rev Mol Cell Biol* **15**, 95–107 (2014).
- 465 8. Ruder, W., Lu, T. & Collins, J. Synthetic biology moving into the clinic. *Science* **333**, 1248–52 (2011).
- 466 9. Slomovic, S., Pardee, K. & Collins, J. Synthetic biology devices for in vitro and in vivo diagnostics. *Proc Natl Acad Sci U*
467 *S A* **112**, 14429–35 (2015).
- 468 10. Friedman, D. C. Industrialization of Biology. A Roadmap to Accelerate the Advanced Manufacturing of Chemicals. Tech.
469 Rep. (2015). URL <https://doi.org/10.2172%2F1213471>.
- 470 11. Tang, N., Ma, S. & Tian, J. New Tools for Cost-Effective DNA Synthesis. In *Synthetic Biology*, 3–21 (Elsevier, 2013).
471 URL <https://doi.org/10.1016%2Fb978-0-12-394430-6.00001-7>.
- 472 12. Doudna, J. A. & Charpentier, E. The new frontier of genome engineering with CRISPR-Cas9. *Science* **346**, 1258096–
473 1258096 (2014). URL <https://doi.org/10.1126%2Fscience.1258096>.
- 474 13. Gardner, T. S. Synthetic biology: from hype to impact. *Trends in Biotechnology* **31**, 123–125 (2013). URL <https://doi.org/10.1016%2Fj.tibtech.2013.01.018>.
- 475 14. Stephens, Z. D. *et al.* Big data: Astronomical or genetical? *PLoS Biology* **13**, 1–11 (2015).
- 476 15. Bath, T. S. *et al.* A targeted proteomics toolkit for high-throughput absolute quantification of Escherichia coli proteins.
477 *Metabolic Engineering* **26**, 48–56 (2014). URL <https://doi.org/10.1016%2Fj.ymben.2014.08.004>.
- 478 16. Fuhrer, T. & Zamboni, N. High-throughput discovery metabolomics. *Current Opinion in Biotechnology* **31**, 73–78 (2015).
479 URL <https://doi.org/10.1016%2Fj.copbio.2014.08.006>.
- 480 17. Heinemann, J. *et al.* Real-Time Digitization of Metabolomics Patterns from a Living System Using Mass Spectrometry.
481 *Journal of The American Society for Mass Spectrometry* **25**, 1755–1762 (2014). URL <https://doi.org/10.1007%2Fs13361-014-0922-z>.
- 482 18. dot}Brien, E. J. O., Lerman, J. A., Chang, R. L., Hyduke, D. R. & Palsson, B. O. Genome-scale models of metabolism and
483 gene expression extend and refine growth phenotype prediction. *Molecular Systems Biology* **9**, 693–693 (2014). URL
484 <https://doi.org/10.1038%2Fmsb.2013.52>.
- 485 19. Karr, J. R. *et al.* A Whole-Cell Computational Model Predicts Phenotype from Genotype. *Cell* **150**, 389–401 (2012). URL
486 <https://doi.org/10.1016%2Fj.cell.2012.05.044>.
- 487 20. Tompson, J., Schlachter, K., Sprechmann, P. & Perlin, K. Accelerating Eulerian Fluid Simulation With Convolutional
488 Networks. *arXiv preprint arXiv:1607.03597* (2016).

- 493 **21.** Lewis, N. E., Nagarajan, H. & Palsson, B. O. Constraining the metabolic genotype–phenotype relationship using
494 a phylogeny of in silico methods. *Nature Reviews Microbiology* (2012). URL [https://doi.org/10.1038%](https://doi.org/10.1038%2Fnrmicro2737)
495 [2Fnrmicro2737](https://doi.org/10.1038%2Fnrmicro2737).
- 496 **22.** Henry, C. S., Broadbelt, L. J. & Hatzimanikatis, V. Thermodynamics-Based Metabolic Flux Analysis. *Biophysical Journal*
497 **92**, 1792–1805 (2007). URL <https://doi.org/10.1529%2Fbiophysj.106.093138>.
- 498 **23.** Martín, H. G. *et al.* A Method to Constrain Genome-Scale Models with ¹³C Labeling Data. *PLOS Computational Biology*
499 **11**, e1004363 (2015). URL <https://doi.org/10.1371%2Fjournal.pcbi.1004363>.
- 500 **24.** Wiechert, W. ¹³C Metabolic Flux Analysis. *Metabolic Engineering* **3**, 195–206 (2001). URL [https://doi.org/10.](https://doi.org/10.1006%2Fmben.2001.0187)
501 [1006%2Fmben.2001.0187](https://doi.org/10.1006%2Fmben.2001.0187).
- 502 **25.** Sauer, U. Metabolic networks in motion: ¹³C-based flux analysis. *Molecular Systems Biology* **2** (2006). URL <https://doi.org/10.1038%2Fmsb4100109>.
- 504 **26.** Ghosh, A. *et al.* ¹³C Metabolic Flux Analysis for Systematic Metabolic Engineering of *S. cerevisiae* for Overproduction
505 of Fatty Acids. *Frontiers in Bioengineering and Biotechnology* **4** (2016). URL [https://doi.org/10.3389%](https://doi.org/10.3389%2Ffbioe.2016.00076)
506 [2Ffbioe.2016.00076](https://doi.org/10.3389%2Ffbioe.2016.00076).
- 507 **27.** Cardenas, J. & Silva, N. A. D. Metabolic engineering of *Saccharomyces cerevisiae* for the production of triacetic acid
508 lactone. *Metabolic Engineering* **25**, 194–203 (2014). URL [https://doi.org/10.1016%2Fj.ymben.2014.07.](https://doi.org/10.1016%2Fj.ymben.2014.07.008)
509 [008](https://doi.org/10.1016%2Fj.ymben.2014.07.008).
- 510 **28.** Xu, P., Ranganathan, S., Fowler, Z. L., Maranas, C. D. & Koffas, M. A. Genome-scale metabolic network modeling results
511 in minimal interventions that cooperatively force carbon flux towards malonyl-CoA. *Metabolic Engineering* **13**, 578–587
512 (2011). URL <https://doi.org/10.1016%2Fj.ymben.2011.06.008>.
- 513 **29.** Lin, F. *et al.* Improving Fatty Acid Availability for Bio-Hydrocarbon Production in *Escherichia coli* by Metabolic
514 Engineering. *PLoS ONE* **8**, e78595 (2013). URL <https://doi.org/10.1371%2Fjournal.pone.0078595>.
- 515 **30.** Khodayari, A., Chowdhury, A. & Maranas, C. D. Succinate Overproduction: A Case Study of Computational Strain
516 Design Using a Comprehensive *Escherichia coli* Kinetic Model. *Frontiers in Bioengineering and Biotechnology* **2** (2015).
517 URL <https://doi.org/10.3389%2Ffbioe.2014.00076>.
- 518 **31.** Matsuoka, Y. & Shimizu, K. Current status and future perspectives of kinetic modeling for the cell metabolism with
519 incorporation of the metabolic regulation mechanism. *Bioresources and Bioprocessing* **2** (2015). URL [https://doi.](https://doi.org/10.1186%2Fs40643-014-0031-7)
520 [org/10.1186%2Fs40643-014-0031-7](https://doi.org/10.1186%2Fs40643-014-0031-7).
- 521 **32.** Fundamentals of Enzyme Kinetics revised edition. Athel Cornish-Bowden, Portland Press, London, 1995, 343 pp., \$29.00.
522 *Analytical Biochemistry* **231**, 275 (1995). URL <https://doi.org/10.1006%2Fbio.1995.1537>.
- 523 **33.** Heinrich, R. & Schuster, S. *The Regulation of Cellular Systems* (Springer US, 1996). URL [https://doi.org/10.](https://doi.org/10.1007%2F978-1-4613-1161-4)
524 [1007%2F978-1-4613-1161-4](https://doi.org/10.1007%2F978-1-4613-1161-4).
- 525 **34.** Costa, R. S., Machado, D., Rocha, I. & Ferreira, E. C. Hybrid dynamic modeling of *Escherichia coli* central metabolic
526 network combining Michaelis–Menten and approximate kinetic equations. *Biosystems* **100**, 150–157 (2010). URL
527 <https://doi.org/10.1016%2Fj.biosystems.2010.03.001>.
- 528 **35.** Horn, F. & Jackson, R. General mass action kinetics. *Archive for Rational Mechanics and Analysis* **47** (1972). URL
529 <https://doi.org/10.1007%2Fbfb00251225>.
- 530 **36.** Hatzimanikatis, V. & Bailey, J. E. Effects of spatiotemporal variations on metabolic control: Ap-
531 proximate analysis using (log)linear kinetic models. *Biotechnology and Bioengineering* **54**, 91–104
532 (1997). URL [https://doi.org/10.1002%2F%28sici%291097-0290%2819970420%2954%3A2%3C91%](https://doi.org/10.1002%2F%28sici%291097-0290%2819970420%2954%3A2%3C91%3A%3Aaid-bit1%3E3.0.co%3B2-q)
533 [3A%3Aaid-bit1%3E3.0.co%3B2-q](https://doi.org/10.1002%2F%28sici%291097-0290%2819970420%2954%3A2%3C91%3A%3Aaid-bit1%3E3.0.co%3B2-q).
- 534 **37.** Heijnen, J. J. Approximative kinetic formats used in metabolic network modeling. *Biotechnology and Bioengineering* **91**,
535 534–545 (2005). URL <https://doi.org/10.1002%2Fbit.20558>.
- 536 **38.** Savageau, M. A. & Voit, E. O. Power-law approach to modeling biological systems: I. Theory. *Journal of fermentation*
537 *technology* **60**, 221–228 (1982).
- 538 **39.** Tran, L. M., Rizk, M. L. & Liao, J. C. Ensemble Modeling of Metabolic Networks. *Biophysical Journal* **95**, 5606–5617
539 (2008). URL <https://doi.org/10.1529%2Fbiophysj.108.135442>.
- 540 **40.** Rizk, M. L. & Liao, J. C. Ensemble Modeling for Aromatic Production in *Escherichia coli*. *PLoS ONE* **4**, e6903 (2009).
541 URL <https://doi.org/10.1371%2Fjournal.pone.0006903>.

- 542 **41.** Tan, Y. & Liao, J. C. Metabolic ensemble modeling for strain engineers. *Biotechnology Journal* **7**, 343–353 (2011). URL
543 <https://doi.org/10.1002%2Fbiot.201100186>.
- 544 **42.** Contador, C. A., Rizk, M. L., Asenjo, J. A. & Liao, J. C. Ensemble modeling for strain development of l-lysine-producing
545 *Escherichia coli*. *Metabolic Engineering* **11**, 221–233 (2009). URL [https://doi.org/10.1016%2Fj.ymben.](https://doi.org/10.1016%2Fj.ymben.2009.04.002)
546 [2009.04.002](https://doi.org/10.1016%2Fj.ymben.2009.04.002).
- 547 **43.** Dean, J. T., Rizk, M. L., Tan, Y., Dipple, K. M. & Liao, J. C. Ensemble Modeling of Hepatic Fatty Acid Metabolism with a
548 Synthetic Glyoxylate Shunt. *Biophysical Journal* **98**, 1385–1395 (2010). URL [https://doi.org/10.1016%2Fj.](https://doi.org/10.1016%2Fj.bpj.2009.12.4308)
549 [bpj.2009.12.4308](https://doi.org/10.1016%2Fj.bpj.2009.12.4308).
- 550 **44.** Khodayari, A., Zomorodi, A. R., Liao, J. C. & Maranas, C. D. A kinetic model of *Escherichia coli* core metabolism
551 satisfying multiple sets of mutant flux data. *Metabolic Engineering* **25**, 50–62 (2014). URL [https://doi.org/10.](https://doi.org/10.1016%2Fj.ymben.2014.05.014)
552 [1016%2Fj.ymben.2014.05.014](https://doi.org/10.1016%2Fj.ymben.2014.05.014).
- 553 **45.** Khodayari, A. & Maranas, C. D. A genome-scale *Escherichia coli* kinetic metabolic model k-ecoli457 satisfying flux
554 data for multiple mutant strains. *Nature Communications* **7**, 13806 (2016). URL [https://doi.org/10.1038%](https://doi.org/10.1038%2Fncoms13806)
555 [2Fncoms13806](https://doi.org/10.1038%2Fncoms13806).
- 556 **46.** Chakrabarti, A., Miskovic, L., Soh, K. C. & Hatzimanikatis, V. Towards kinetic modeling of genome-scale metabolic
557 networks without sacrificing stoichiometric thermodynamic and physiological constraints. *Biotechnology Journal* **8**,
558 1043–1057 (2013). URL <https://doi.org/10.1002%2Fbiot.201300091>.
- 559 **47.** Savoglidis, G. *et al.* A method for analysis and design of metabolism using metabolomics data and kinetic models:
560 Application on lipidomics using a novel kinetic model of sphingolipid metabolism. *Metabolic Engineering* **37**, 46–62
561 (2016). URL <https://doi.org/10.1016%2Fj.ymben.2016.04.002>.
- 562 **48.** Gerosa, L. *et al.* Pseudo-transition Analysis Identifies the Key Regulators of Dynamic Metabolic Adaptations from Steady-
563 State Data. *Cell Systems* **1**, 270–282 (2015). URL <https://doi.org/10.1016%2Fj.cels.2015.09.008>.
- 564 **49.** Hackett, S. R. *et al.* Systems-level analysis of mechanisms regulating yeast metabolic flux. *Science* **354**, aaf2786–aaf2786
565 (2016). URL <https://doi.org/10.1126%2Fscience.aaf2786>.
- 566 **50.** Daran-Lapujade, P. *et al.* The fluxes through glycolytic enzymes in *Saccharomyces cerevisiae* are predominantly regulated
567 at posttranscriptional levels. *Proceedings of the National Academy of Sciences* **104**, 15753–15758 (2007). URL <https://doi.org/10.1073%2Fpnas.0707476104>.
- 569 **51.** Abernathy, M. H., He, L. & Tang, Y. J. Channeling in native microbial pathways: Implications and challenges for metabolic
570 engineering. *Biotechnology Advances* (2017). URL [https://doi.org/10.1016%2Fj.biotechadv.2017.06.](https://doi.org/10.1016%2Fj.biotechadv.2017.06.004)
571 [004](https://doi.org/10.1016%2Fj.biotechadv.2017.06.004).
- 572 **52.** Noor, E. *et al.* Pathway Thermodynamics Highlights Kinetic Obstacles in Central Metabolism. *PLoS Computational*
573 *Biology* **10**, e1003483 (2014). URL <https://doi.org/10.1371%2Fjournal.pcbi.1003483>.
- 574 **53.** Digel, M., Ehehalt, R., Stremmel, W. & Füllekrug, J. Acyl-CoA synthetases: fatty acid uptake and metabolic
575 channeling. *Molecular and Cellular Biochemistry* **326**, 23–28 (2008). URL [https://doi.org/10.1007%](https://doi.org/10.1007%2Fs11010-008-0003-3)
576 [2Fs11010-008-0003-3](https://doi.org/10.1007%2Fs11010-008-0003-3).
- 577 **54.** Thrun, S. Toward robotic cars. *Communications of the ACM* **53**, 99 (2010). URL [https://doi.org/10.1145%](https://doi.org/10.1145%2F1721654.1721679)
578 [2F1721654.1721679](https://doi.org/10.1145%2F1721654.1721679).
- 579 **55.** Silver, D. *et al.* Mastering the game of Go with deep neural networks and tree search. *Nature* **529**, 484–489 (2016). URL
580 <https://doi.org/10.1038%2Fnature16961>.
- 581 **56.** Wu, Y. *et al.* Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv*
582 *preprint arXiv:1609.08144* (2016).
- 583 **57.** Kosinski, M., Stillwell, D. & Graepel, T. Private traits and attributes are predictable from digital records of human behavior.
584 *Proceedings of the National Academy of Sciences* **110**, 5802–5805 (2013). URL [https://doi.org/10.1073%](https://doi.org/10.1073%2Fpnas.1218772110)
585 [2Fpnas.1218772110](https://doi.org/10.1073%2Fpnas.1218772110).
- 586 **58.** The Data That Turned the World Upside Down. URL [https://motherboard.vice.com/en_us/article/](https://motherboard.vice.com/en_us/article/mg9vvn/how-our-likes-helped-trump-win)
587 [mg9vvn/how-our-likes-helped-trump-win](https://motherboard.vice.com/en_us/article/mg9vvn/how-our-likes-helped-trump-win).
- 588 **59.** Alipanahi, B., Delong, A., Weirauch, M. T. & Frey, B. J. Predicting the sequence specificities of DNA- and RNA-binding
589 proteins by deep learning. *Nature Biotechnology* **33**, 831–838 (2015). URL [https://doi.org/10.1038%2Fnbt.](https://doi.org/10.1038%2Fnbt.3300)
590 [3300](https://doi.org/10.1038%2Fnbt.3300).

- 591 **60.** Esteva, A. *et al.* Dermatologist-level classification of skin cancer with deep neural networks. *Nature* **542**, 115–118 (2017).
592 URL <https://doi.org/10.1038%2Fnature21056>.
- 593 **61.** Poplin, R. *et al.* Creating a universal snp and small indel variant caller with deep neural networks (2016). URL
594 <https://doi.org/10.1101%2F092890>.
- 595 **62.** Paeng, K., Hwang, S., Park, S., Kim, M. & Kim, S. A unified framework for tumor proliferation score prediction in breast
596 histopathology. *arXiv preprint arXiv:1612.07180* (2016).
- 597 **63.** Aguirre, L. A. & Billings, S. A. Dynamical effects of over parametrization in nonlinear models. *Physica D* **80**, 26–40
598 (1995).
- 599 **64.** Ljung, L. Approaches to identification of nonlinear systems. *Control Conference (CCC), 2010 29th Chinese* 1–5 (2010).
600 URL <http://ieeexplore.ieee.org/xpls/abs{ }all.jsp?arnumber=5572936>.
- 601 **65.** Villaverde, A. F. & Banga, J. R. Reverse engineering and identification in systems biology: strategies, perspectives and
602 challenges. *Journal of the Royal Society Interface* **11**, 20130505 (2013).
- 603 **66.** Brunk, E. *et al.* Characterizing Strain Variation in Engineered E. coli Using a Multi-Omics-Based Workflow. *Cell Systems*
604 **2**, 335–346 (2016). URL <http://dx.doi.org/10.1016/j.cels.2016.04.004>.
- 605 **67.** Weaver, L. J. Towards predictive metabolic engineering: kinetic modeling and experimental analysis of a heterologous
606 mevalonate pathway in e. coli (2013).
- 607 **68.** Van Dien, S. From the first drop to the first truckload: commercialization of microbial processes for renewable chemicals.
608 *Current opinion in biotechnology* **24**, 1061–1068 (2013).
- 609 **69.** Alonso-Gutierrez, J. *et al.* Principal component analysis of proteomics (pcap) as a tool to direct metabolic engineering.
610 *Metabolic engineering* **28**, 123–133 (2015).
- 611 **70.** Ishii, N. *et al.* Multiple high-throughput analyses monitor the response of e. coli to perturbations. *Science* **316**, 593–597
612 (2007).
- 613 **71.** Ma, Q. *et al.* Integrated proteomic and metabolomic analysis of an artificial microbial community for two-step production
614 of vitamin c. *PloS one* **6**, e26108 (2011).
- 615 **72.** Yang, S. *et al.* Clostridium thermocellum atcc27405 transcriptomic, metabolomic and proteomic profiles after ethanol
616 stress. *Bmc Genomics* **13**, 336 (2012).
- 617 **73.** Doerfler, H. *et al.* Granger causality in integrated gc–ms and lc–ms metabolomics data reveals the interface of primary and
618 secondary metabolism. *Metabolomics* **9**, 564–574 (2013).
- 619 **74.** Dyar, K. A. & Eckel-Mahan, K. L. Circadian metabolomics in time and space. *Frontiers in neuroscience* **11**, 369 (2017).
- 620 **75.** Patel, V. R., Eckel-Mahan, K., Sassone-Corsi, P. & Baldi, P. Circadiomics: integrating circadian genomics, transcriptomics,
621 proteomics and metabolomics. *Nature methods* **9**, 772 (2012).
- 622 **76.** Arkin, A. P. *et al.* The doe systems biology knowledgebase (kbase). *bioRxiv* 096354 (2016).
- 623 **77.** Morrell, W. C. *et al.* The experiment data depot: a web-based software tool for biological experimental data storage,
624 sharing, and visualization. *ACS synthetic biology* **6**, 2248–2259 (2017).
- 625 **78.** George, K. *et al.* Correlation analysis of targeted proteins and metabolites to assess and engineer microbial isopentenol
626 production. *Biotechnol Bioeng* **111**, 1648–58 (2014).
- 627 **79.** George, K. W. *et al.* Metabolic engineering for the high-yield production of isoprenoid-based C5 alcohols in E. coli.
628 *Scientific Reports* **5** (2015). URL <https://doi.org/10.1038%2Fsrep11128>.
- 629 **80.** Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* **12**, 2825–2830
630 (2011).
- 631 **81.** Pan, S. J. & Yang, Q. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering* **22**, 1345–1359
632 (2010).
- 633 **82.** Gerber, G. The dynamic microbiome. *FEBS Lett* **588**, 4131–9 (2014).
- 634 **83.** Price, N. D. *et al.* A wellness study of 108 individuals using personal, dense, dynamic data clouds. *Nature biotechnology*
635 **35**, 747 (2017).
- 636 **84.** Chen, R. & Snyder, M. Promise of personalized omics to precision medicine. *Wiley Interdisciplinary Reviews: Systems*
637 *Biology and Medicine* **5**, 73–82 (2013).

- 638 **85.** Heintz-Buschart, A. *et al.* Integrated multi-omics of the human gut microbiome in a case study of familial type 1 diabetes.
639 *Nature Microbiology* **2**, 16180 (2016). URL <https://doi.org/10.1038%2Fnmicrobiol.2016.180>.
- 640 **86.** Narayanasamy, S., Muller, E. E. L., Sheik, A. R. & Wilmes, P. Integrated omics for the identification of key functionalities
641 in biological wastewater treatment microbial communities. *Microbial Biotechnology* **8**, 363–368 (2015). URL <https://doi.org/10.1111%2F1751-7915.12255>.
642
- 643 **87.** Muller, E. E. L. *et al.* Community-integrated omics links dominance of a microbial generalist to fine-tuned resource usage.
644 *Nature Communications* **5**, 5603 (2014). URL <https://doi.org/10.1038%2Fncomms6603>.
- 645 **88.** Shah, P. *et al.* A microfluidics-based in vitro model of the gastrointestinal human–microbe interface. *Nature Communica-*
646 *tions* **7**, 11535 (2016). URL <https://doi.org/10.1038%2Fncomms11535>.
- 647 **89.** Link, H., Fuhrer, T., Gerosa, L., Zamboni, N. & Sauer, U. Real-time metabolome profiling of the metabolic switch between
648 starvation and growth. *Nature Methods* (2015). URL <https://doi.org/10.1038%2Fnmeth.3584>.
- 649 **90.** The Rise And Fall Of The Company That Was Going To Have Us All Us-
650 ing Biofuels. [https://www.fastcompany.com/3000040/rise-and-fall-company-was-going-](https://www.fastcompany.com/3000040/rise-and-fall-company-was-going-have-us-all-using-biofuels)
651 [have-us-all-using-biofuels.](https://www.fastcompany.com/3000040/rise-and-fall-company-was-going-have-us-all-using-biofuels) URL [https://www.fastcompany.com/3000040/](https://www.fastcompany.com/3000040/rise-and-fall-company-was-going-have-us-all-using-biofuels)
652 [rise-and-fall-company-was-going-have-us-all-using-biofuels.](https://www.fastcompany.com/3000040/rise-and-fall-company-was-going-have-us-all-using-biofuels) Accessed on Wed, Octo-
653 ber 11, 2017.
- 654 **91.** Savitzky, A. & Golay, M. J. Smoothing and differentiation of data by simplified least squares procedures. *Analytical*
655 *chemistry* **36**, 1627–1639 (1964).
- 656 **92.** Olson, R. S. *et al.* *Applications of Evolutionary Computation: 19th European Conference, EvoApplications 2016, Porto,*
657 *Portugal, March 30 – April 1, 2016, Proceedings, Part I*, chap. Automating Biomedical Data Science Through Tree-Based
658 Pipeline Optimization, 123–137 (Springer International Publishing, 2016). URL [http://dx.doi.org/10.1007/](http://dx.doi.org/10.1007/978-3-319-31204-0_9)
659 [978-3-319-31204-0_9](http://dx.doi.org/10.1007/978-3-319-31204-0_9).
- 660 **93.** Friedman, J., Hastie, T. & Tibshirani, R. *The elements of statistical learning*, vol. 1 (Springer series in statistics New York,
661 2001).