

## RESEARCH ARTICLE SUMMARY

## SYSTEMS BIOLOGY

# Simultaneous cross-evaluation of heterogeneous *E. coli* datasets via mechanistic simulation

Derek N. Macklin, Travis A. Ahn-Horst, Heejo Choi, Nicholas A. Ruggero, Javier Carrera, John C. Mason, Gwanggyu Sun, Eran Agmon, Mialy M. DeFelice, Inbal Maayan, Keara Lane, Ryan K. Spangler, Taryn E. Gillies, Morgan L. Paull, Sajja Akhter, Samuel R. Bray, Daniel S. Weaver, Ingrid M. Keseler, Peter D. Karp, Jerry H. Morrison, Markus W. Covert\*

**INTRODUCTION:** The generation of biological data is presenting us with one of the most demanding analysis challenges the world has ever faced, not only in terms of storage and accessibility, but more critically in terms of its extensive heterogeneity and variability. Although issues associated with heterogeneity and variability each represent major analysis problems on their own, the challenges posed by both in combination are even more difficult but also present greater opportunities. The problems arise because assessing the data's veracity means not only determining whether the data are reproducible but also, and perhaps more deeply, whether they are cross-consistent, meaning that the interpretations of multiple heterogeneous datasets all point to the same conclusion. The opportunities emerge because seemingly discrepant results across multiple studies and measurement modalities may not be due simply to the errors associated with particular techniques, but also to the complex, nonlinear, and highly interconnected nature of biology. Therefore, what is required are analysis methods that can integrate and evaluate multiple data types simultaneously and in the context of biological mechanisms.

**RATIONALE:** Here, we present a large-scale, integrated modeling approach to simultaneously cross-evaluate millions of heterogeneous data against themselves, based on an extensive computer model of *Escherichia coli* that accounts for the function of 1214 genes (or 43% of the well-annotated genes). The model incorporates an extensive set of diverse measurements compiled from thousands of reports and accounting for many decades of research performed in laboratories around the world. Curation of these data led to the identification

of >19,000 parameter values, which we integrated by creating a computational model that brings molecular signaling and regulation of RNA and protein expression together with carbon and energy metabolism in the context of balanced growth. A major



**Integrating experimental and computational components, scientists constructed a model of *E. coli*.** Although the model described here resides as software (freely available on GitHub), the model depicted in the photo above is composed of Corning plasticware and filter tips, network cables, and Mac accessories.

advantage of this modeling approach is that heterogeneous data are linked mechanistically through the simulated interaction of cellular processes, providing the most natural, intuitive interpretation of an integrated dataset. Thus, this model enabled us to assess the cross-consistency

of all of these datasets as an integrated whole.

**RESULTS:** We assessed the cross-consistency of the parameter set and identified areas of inconsistency by populating our model with the literature-derived parameters and by running detailed simulations of cellular life cycles. Although analysis of these simulations showed that most of the data were in fact cross-consistent, we also identified critical areas in which the data incorporated in our model were not. These inconsistencies led to readily observable consequences, including that the total output of the ribosomes and RNA polymerases described by the data are not sufficient for a cell to reproduce measured doubling times, that measured metabolic parameters are neither fully compatible with each other nor with overall growth, and that essential proteins are absent during the cell cycle—and the cell is robust to this absence. After correcting for these inconsistencies, the model is capable of validatable predictions compared with previously withheld data. Finally, considering these data as a whole led to successful predictions in vitro, in this case protein half-lives.

**CONCLUSION:** Construction of a highly integrative and mechanistic mathematical model provided us with an opportunity to integrate and cross-validate a vast, heterogeneous dataset in *E. coli*, a process we now call “deep curation” to reflect the multiple layers of curation that we perform (analogous to “deep learning” and “deep sequencing”). By highlighting areas in which studies in *E. coli* contradict each other, our work suggests lines of fruitful experimental inquiry that may help to resolve discrepancies, leading to both new biological insights and a more coherent understanding of this critical model organism. We hope that this work, by demonstrating the value of a large-scale integrative approach to understanding, interpreting, and cross-validating large datasets, will inspire further efforts to comprehensively characterize other organisms of interest. ■

**ON OUR WEBSITE**  
Read the full article at <https://dx.doi.org/10.1126/science.aav3751>

The list of author affiliations is available in the full article online.  
\*Corresponding author. Email: [mcovert@stanford.edu](mailto:mcovert@stanford.edu)  
Cite this article as D. N. Macklin *et al.*, *Science* **369**, eaav3751 (2020). DOI: [10.1126/science.aav3751](https://doi.org/10.1126/science.aav3751)

## RESEARCH ARTICLE

## SYSTEMS BIOLOGY

# Simultaneous cross-evaluation of heterogeneous *E. coli* datasets via mechanistic simulation

Derek N. Macklin<sup>1,2,\*†</sup>, Travis A. Ahn-Horst<sup>1,2,\*</sup>, Heejo Choi<sup>1,2,\*</sup>, Nicholas A. Ruggero<sup>2,3,\*†</sup>, Javier Carrera<sup>1,2,\*§</sup>, John C. Mason<sup>1,2,\*¶</sup>, Gwanggyu Sun<sup>1,2</sup>, Eran Agmon<sup>1,2</sup>, Mialy M. DeFelice<sup>1,2</sup>, Inbal Maayan<sup>1,2</sup>, Keara Lane<sup>1,2,#</sup>, Ryan K. Spangler<sup>1,2</sup>, Taryn E. Gillies<sup>1,2</sup>, Morgan L. Paull<sup>1,\*\*</sup>, Sajia Akhter<sup>1</sup>, Samuel R. Bray<sup>1</sup>, Daniel S. Weaver<sup>4,†</sup>, Ingrid M. Keseler<sup>4</sup>, Peter D. Karp<sup>4</sup>, Jerry H. Morrison<sup>2</sup>, Markus W. Covert<sup>1,2,††</sup>

The extensive heterogeneity of biological data poses challenges to analysis and interpretation. Construction of a large-scale mechanistic model of *Escherichia coli* enabled us to integrate and cross-evaluate a massive, heterogeneous dataset based on measurements reported by various groups over decades. We identified inconsistencies with functional consequences across the data, including that the total output of the ribosomes and RNA polymerases described by data are not sufficient for a cell to reproduce measured doubling times, that measured metabolic parameters are neither fully compatible with each other nor with overall growth, and that essential proteins are absent during the cell cycle—and the cell is robust to this absence. Finally, considering these data as a whole leads to successful predictions of new experimental outcomes, in this case protein half-lives.

The generation of biological data is rapidly presenting us with one of the most demanding data analysis challenges the world has ever faced (1), not only in terms of storage and accessibility, but perhaps more critically, in terms of its extensive heterogeneity and variability (2). With respect to heterogeneity, study of a biological system of interest typically involves many diverse measurements, from lower-throughput blotting techniques to high-throughput sequence- and spectrometry-based technologies and beyond. In terms of variability, it is often the case that studies produced independently from each other report results that seem to be at odds with one another. This is most readily apparent when studies of the same system perform the same measurements but obtain different results, an issue that has led high-profile journals to question the reproducibility of results in multiple scientific fields (3, 4).

Although issues associated with heterogeneity and variability each represent major analysis problems on their own, the challenges posed by both in combination are even

more difficult—but also present greater opportunities for discovery. The problems arise because assessing the data's veracity means not only determining whether the data are reproducible (i.e., does a repeated study produce the same measured outcomes?) but also, and perhaps more deeply, whether they are cross-consistent, meaning that the interpretation of multiple heterogeneous datasets all points to the same conclusion. The opportunities emerge because seemingly discrepant results across multiple studies and measurement modalities may be due, not just to the error associated with a technique or the human hands performing it, but also to the complex, nonlinear, and highly interconnected nature of biology. In such cases, the identification of data discrepancy would be a strong indicator for future insight and discovery.

To this end, the goal of this project was to cross-evaluate a massive, heterogeneous set of measurements that have been reported in the model organism *Escherichia coli* in thousands of studies and by hundreds of laboratories over the past several decades. Determining the cross-consistency between these various measurements requires an understanding of the known or presumed biological relationships that connect them. Thus, we adopted a mathematical approach that can represent these relationships mechanistically while simultaneously accommodating many millions of heterogeneous data points. Efforts to model cell behavior at the cell scale span several decades (5–12). We previously reported a modeling approach that was capable of integrating all of the known functions in the simplest culturable bacterium, *Mycoplasma genitalium* (13).

A major advantage of this “whole-cell” modeling approach is that heterogeneous data are linked mechanistically through the simulated interaction of cellular processes, providing the most natural, intuitive interpretation of an integrated dataset (14). The *M. genitalium* model successfully reproduced many measured data and even predicted previously unmeasured parameters that were subsequently verified experimentally (15). Construction of this model also enabled us to cross-evaluate data and identify discrepancies. As a relatively simple but illustrative example, the DNA concentration per cell measured in *M. genitalium* was only a fraction of the DNA mass required to make up the genome sequence (13). This led us to favor the genome sequence data in determining the parameters governing DNA concentration.

*E. coli* has nearly 10 times more genes than *M. genitalium*, comprises ~50 times as many molecules, can readily grow in a wide variety of environmental conditions, and exhibits extensive self-regulation and control, all of which pose challenges to whole-cell modeling. The model described in this report only accounts for a subset of these genes, environments, and functions. However, one of the most exciting aspects of modeling *E. coli* on a large scale is the enormous effort in data generation that has already been performed. Thus, whereas only 27.5% of the parameter values in our *M. genitalium* model were actually derived from measurements using that organism, 100% of the values incorporated into the model that we describe here were derived directly from *E. coli*. This provided us with an unprecedented opportunity to assess the literature against itself.

Our overall approach is depicted in Fig. 1 and movie S1. We compiled an extensive set of high- and low-throughput measurements from databases and published reports to identify datasets that characterize mRNA and protein expression under a variety of environmental conditions (some of which we generated for this study; see the materials and methods), mRNA and protein half-lives, ribonuclease kinetics, gene locations, transcription factor-binding sites, dissociation constants for proteins bound to DNA-binding sites or other cellular and environmental ligands, translational efficiencies of mRNA transcripts, chemical reaction stoichiometry, enzyme kinetic and substrate transport rates, internal metabolite concentrations, ribosome and RNA polymerase concentrations and elongation rates, the rate of DNA initiation and other cell cycle parameters, and other physiological properties (e.g., growth rates and chemical composition of the cell) (see the supplementary materials for a complete description of included data).

Curation of these data led to the identification of >19,000 parameter values, which are listed

<sup>1</sup>Department of Bioengineering, Stanford University, Stanford, CA 94305, USA. <sup>2</sup>Allen Discovery Center at Stanford University, Stanford University, Stanford, CA 94305, USA.

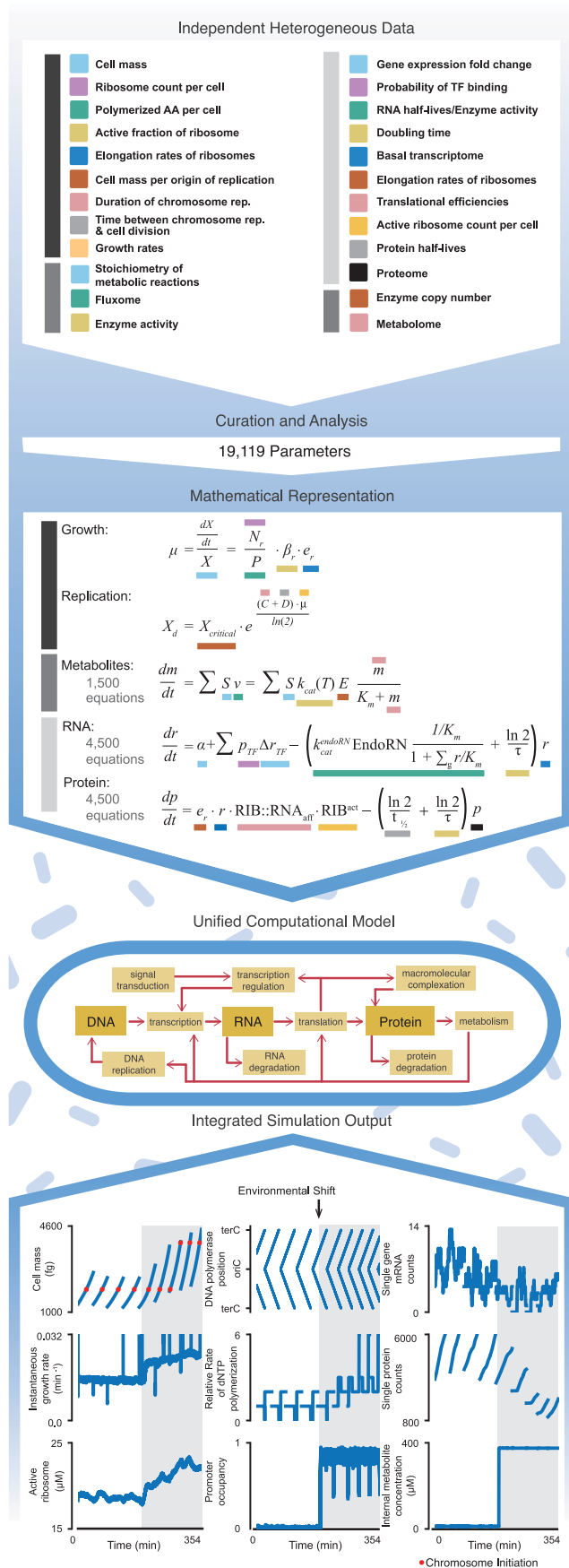
<sup>3</sup>Department of Chemical Engineering, Stanford University, Stanford, CA 94305, USA. <sup>4</sup>SRI International, Menlo Park, CA 94025, USA.

\*These authors contributed equally to this work. †Present address: Grand Rounds, Inc., San Francisco, CA 94107, USA. ‡Present address: X, Mountain View, CA 94030, USA. §Present address: Zymergen, Emeryville, CA 94608, USA. ¶Present address: Intrexon, South San Francisco, CA 94080, USA. #Present address: Department of Molecular Biosciences, Northwestern University, Evanston, IL 60208, USA. \*\*Present address: BridgeBio Pharma, Palo Alto, CA 94301, USA.

††Corresponding author. Email: mcovert@stanford.edu

**Fig. 1. Large-scale, integrated modeling approach to simultaneously cross-evaluate millions of heterogeneous data.**

The data were collected from the primary literature and key databases, and in some cases were also generated as part of this study. Subsequent data curation and analysis led to the determination of 19,119 parameter values. We then incorporated these data into a large-scale computational model of *E. coli* gene expression, metabolism, and growth based on a foundation of >10,000 interdependent mathematical equations that were then transformed into appropriate computational representations of biological processes. Color coding is used to connect terms in these equations to the data that produced their parameter values. This unified model was then used to produce fully integrated simulations, with output as shown at the bottom. See fig. S1, movies S1 and S2, and the supplementary materials for more details. Full details of the analysis required to generate this figure, as well as a pointer to the generating code, can be found in the supplementary materials, section 1.2.



Downloaded from https://www.science.org at Stanford University on May 12, 2023

by category in table S1 and described in detail in the GitHub repository for our model (<https://github.com/CovertLab/WholeCellEcoliRelease>). To compile these values, we created a computational model that brings RNA and protein expression together with carbon and energy metabolism in the context of balanced growth. These datasets are integrated mathematically, beginning with a system of >10,000 mathematical equations schematically illustrated in Fig. 1 [we used ordinary differential equations (ODEs) here as a reduced representation of the actual model, the implementation of which is more complex; for details, see the supplementary materials]. Functionally, 1214 genes (or 43% of the well-annotated genes) were included to represent these processes, which required several major improvements over our previous work in *M. genitalium* not only in terms of modeling but also software improvements in run time and accessibility (for details, see the supplementary materials). For this study, the model–data comparisons were examined under conditions of exponential growth in three experimentally characterized environments: minimal medium (M9 salts plus glucose under aerobic conditions), rich medium (minimal medium plus all amino acids), and minimal anaerobic medium.

We assessed the cross-consistency of the parameter set as a whole and identified areas of inconsistency by populating our model with these literature-derived parameters and by running detailed simulations of cellular life cycles. In the analysis of these simulations, we identified several critical areas in which the data contributing to these models were not cross-consistent. These inconsistencies led to readily observable consequences. Moreover, by incorporating these findings, we constructed a functional and predictive model that produced simulation output, as shown in Fig. 1, fig. S1, and movie S2.

The first inconsistency we identified was that the total output of the ribosomes and RNA polymerases (as derived from the integrated datasets) was not sufficient for the simulated cell to reproduce measured growth rates. The overall growth of the cell depends on the production of protein, which in turn is largely governed by these two major complexes, the cell's mRNA and protein synthesis machinery. The ribosomal content of the cell has been measured or estimated for different growth rates, as have the expression and half-lives of the ribosomal RNA and protein components (11, 16, 17), their associated translational efficiencies (18), and the stoichiometry of the functional complex (19). The expression and half-lives of the RNA polymerase subunits has also been measured or estimated (11, 16). When these measurements were integrated into our simulation, the resulting median doubling time for our simulations was 125 min,

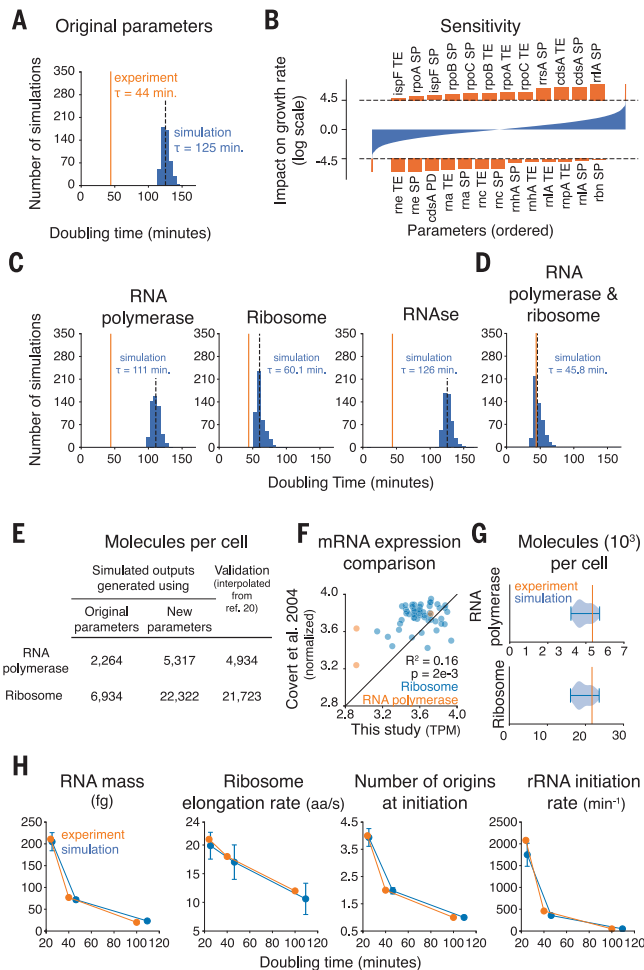
compared with the 44 min measured experimentally for cells growing on glucose minimal medium (Fig. 2A). Thus, the doubling time measurements and the measurements related to ribosomal and/or RNA polymerase output appeared to be inconsistent.

To further dissect this inconsistency, we performed a sensitivity analysis to determine which parameters were most likely to have an impact on the doubling time. We ran 20,000 simulations, each for 10 s of simulation time, in which 10% of the parameter values were randomly chosen and their value increased or decreased by fivefold (also chosen at random). To cause an observable impact, it was necessary to vary many parameter values at once because there are so many interaction effects between parameters. After the growth rate was determined at the end of each simulation, the effect of a particular parameter on growth rate was determined by finding the average growth difference between the cases in which the parameter was raised and when it was lowered and then assessing each parameter's individual effect in the context of the total distribution of parameter effects. The top hits from our analysis involved parameters related to ribosomal and RNA polymerases, RNAses, and a metabolic enzyme encoded by the *cdsA* gene (Fig. 2B).

On the basis of these findings, we first considered changing parameters related to the expression of ribosomes, RNA polymerases, and RNAses (the enzyme *cdsA* is considered in more detail below). When increasing the expression of one protein, the expression of all other genes must be decreased to maintain the total amount of mRNA and protein per cell at their experimentally measured values. Thus, we used an iterative parameter estimation approach based on ODEs that calculates the amount of protein produced from the ribosomal and RNA polymerase content at a given growth rate (see the materials and methods). Our results showed that increasing the RNA polymerase, ribosomal, or RNase expression alone was not sufficient to lower the doubling time to measured values (Fig. 2C). However, an increase in the expression for both RNA polymerases and ribosomes did enable us to simulate an accurate doubling time (Fig. 2D). The new polymerase and ribosome calculations matched well with estimates of expression [compiled in (20)] that were not used to create our model (Fig. 2E).

Although these results supported the hypothesis that the expression of RNA polymerases and ribosomes was not adequately captured by the initial parameters fed into the model, it was not clear which parameters were most likely to be problematic. Thus, we evaluated each parameter contained in our RNA polymerase and ribosomal expression

equations, grading them on the following three criteria: (i) literature reproducibility, meaning that the parameter value could be supported by independent measurements; (ii) whether changing the parameters would lead to an adequate change in the simulated doubling time; and (iii) whether the simulations performed in (ii) also matched the abundances of ribosomes and RNA polymerases from Fig. 2E (20). This analysis (detailed in figs. S2A and S2B) revealed that the transcript synthesis probabilities of genes that produce subunits of RNA polymerases and ribosomes were the most favorable parameters to change because they were relatively variable between experiments (Fig. 2F) and had a strong enough effect on the doubling time (Fig. 2D) and protein abundances (Fig. 2G). Thus, we calculated new gene transcription probabilities for RNA polymerase and ribosomal subunits based on the measured doubling time rather than on global mRNA measurements; these new transcription probabilities are the only changes to the data that continue to the rest of this study (table S2A). In total, the production of all RNA polymerase genes had to be increased by roughly twofold to recapitulate measured growth rates (see table S2B for all changes to expression parameters). Ribosomal gene expression was more complex; although some genes required an increase in the production rate greater than threefold, the expression of other subunits was actually decreased. Accommodating these changes further required a global decrease in production rate (for all other nonribosome and non-RNA polymerase genes) to ~89% of their original values to maintain the overall RNA mass in the cell. These adjustments led to simulated doubling times that were consistent with measurements on the glucose minimal aerobic medium (Fig. 2D). Similar analyses were performed for the other two environments; the final simulations in all three simulated environments were consistent, not only with doubling times (fig. S2C), but also with other measurements including RNA mass per cell, ribosome elongation rates, stable RNA synthesis rates, and the average number of DNA replication origins per cell at the time of replication initiation (Fig. 2H) (21). The final simulations could also reproduce the linear relationship between the RNA/protein mass ratio and the growth rate that was previously observed for cells growing in different environments (22) (fig. S2D). Finally, the simulation output also showed that in fast-growing cells, the cell mass added over the life cycle was uncorrelated with the initial cell mass (a phenomenon referred to as “adder” behavior), whereas for slower-growing cells, the added and initial cell masses were correlated (“sizer” behavior) (fig. S2E), in agreement with recent reports (23–26). We concluded that



**Fig. 2. Ribosomal and RNA polymerase output must be increased to support measured doubling times.** (A) Histogram comparing simulated doubling times (blue) with the experimentally determined doubling time for aerobic growth on glucose minimal medium (orange line) with the model's original parameter values taken directly from the literature. Median simulated doubling time is 125 min (dashed black line). (B) Sensitivity analysis outcome reported as the z-score (log-scale) of the difference in growth rate for all simulations in which a given parameter was adjusted higher and all simulations in which a given parameter was adjusted lower. Horizontal dashed lines represent a z-score cutoff for a  $P$  value  $< 0.05$  that has been adjusted for multiple hypothesis testing of each of the parameters that were adjusted (93% of the total parameters; see the supplementary materials for more details). Parameters are ordered by their impact on the simulated cells' growth rate along the x-axis; those having a significant z-score are highlighted in orange and shown in more detail above and below the plot of all parameters. Parameters with the largest positive correlation with model growth are listed across the top, and parameters with the largest negative correlation are listed across the bottom. TE, translational efficiency; SP, RNA synthesis probability; PD, protein degradation rate. (C and D) Histograms comparing simulated doubling times (blue) with the experimentally determined doubling time for aerobic growth on glucose minimal medium (orange line), with RNA polymerase, ribosome, and RNase expression calculated from the known doubling time as independent experiments (C) and with both RNA polymerase and ribosome expression calculated from the known doubling time (D). Median simulated doubling times are shown as dashed black lines. (E) RNA polymerase and ribosome abundances per cell as generated by the model in this study using the original (Fig. 2A) and new (Fig. 2D) transcript synthesis probabilities compared with experimental data withheld from the model's original parameterization from (20). (F) Comparison of mRNA expression as measured by RNA-sequencing in this study (TPM, transcripts per million) and from a previous microarray study (51). (G) Violin plots showing distributions of RNA polymerase and ribosome cellular abundances from the simulations shown in Fig. 2D compared with expected values determined experimentally (orange lines) (20). (H) Cellular properties calculated from the simulations for three different environmental conditions compared with their counterpart measurements reported in the literature (21). Error bars indicate SDs of each property calculated over the 1024 cells that were simulated for each medium. Full details of the analysis required to generate this figure, as well as a pointer to the generating code, can be found in the supplementary materials, section 1.2.

modifying the parameters related to the expression of certain ribosomal subunits, together with a global increase in RNA polymerase expression, caused our simulations to better reflect multiple physiological observations.

The second major discrepancy we found concerned the parameter values that determine the activity and output of *E. coli*'s metabolic network. These are the kinetic parameters of each biochemical reaction, as well as the parameters related to gene expression for each metabolic enzyme. Taken as a whole, these parameter values must be consistent with each other such that the metabolic network can support mass and energy demands without unstable pooling or depletion of intermediate metabolites. In practical terms, this means that the chemical composition of a cell or the metabolic demand on the cell has to be balanced with the supply provided by the metabolic network.

Metabolism is probably the most thoroughly characterized network in *E. coli* (8, 27, 28). In our model, a metabolic network model derived from the EcoCyc database (29) is represented using an expansion of flux balance analysis (FBA), which uses an optimization strategy to predict metabolic network behavior even when few parameters are known (30). To add kinetic information to this model, we searched through the literature—thousands of studies in all—and identified 639 relevant kinetic parameters governing the activity of 404 biochemical reactions in the metabolic network. Whereas traditional FBA is based on an objective function that serves to maximize biomass concentration in fixed relative proportions, our method uses an objective function that is both more flexible [and thus better suited to dynamic simulations (31)] and explicitly incorporates kinetic parameters, as well as metabolite and enzyme concentrations. Specifically, we implemented a two-term objective that penalizes unbalanced growth or depletion of intermediate metabolite concentrations (the metabolic cost function) while also encouraging the flux through the network to match that predicted using the kinetic parameters described above (the kinetic cost function). These two terms are related by a weighting factor, which we set to optimize a trade-off between including kinetic data in the model while not compromising cell growth (see fig. S3, A to F, and the supplementary materials for complete details).

During this process, we noticed three areas of inconsistency with regard to metabolism. First, low expression of enzyme-encoding genes could overconstrain the biochemical capacity of the metabolic network. The only example of this we found concerned the *cdsA* gene product; in particular, we found that a significant fraction of simulations would not produce an adequate number of phospholipids unless *cdsA* expression was artificially increased

in the model (Fig. 3A). We investigated the low expression of this gene further in the context of the RNA-Seq (18), proteomics (32), and gene essentiality datasets (33), the latter two of which were not used in the construction of the model. This comparison confirmed that mRNA expression of *cdsA* was indeed low (which was further confirmed by quantitative polymerase chain reaction; fig. S3G) but was detectable at the protein level, and also that it was an essential gene. That this essential protein, identified in Fig. 2B as one of the most important effectors of simulation doubling time, was so lowly expressed that its count dropped to zero in the simulations was a puzzling contradiction within our data (this is investigated further below).

Having considered the constraints that low *cdsA* expression imposed on the metabolic network, we then turned to the constraints imposed by kinetics and found that the kinetic parameter set in its initial form was also inconsistent with (i.e., unable to produce) known cellular growth rates. Preliminary comparisons between the simulations with and without kinetics specifically identified the constraints on succinate dehydrogenase and fumarate reductase as preventing cell growth because of inefficient carbon source utilization (Fig. 3B). The constraints imposed on these enzymes by their parameter values were therefore initially removed from the model. However, comparing our simulated metabolic flux outputs with a metabolic flux validation dataset (34) that was not originally used to create or parameterize the model, we found that the simulation and data were highly correlated, with the exception of two fluxes in the citric acid cycle: those mediated by succinate and isocitrate dehydrogenase (Fig. 3C). The identification of succinate dehydrogenase as problematic in both analyses, even with its kinetic constraint removed, indicated that the kinetic parameters for other reactions might also be responsible for our observations. Thus, we performed a global analysis in which every kinetic constraint was tested individually to determine whether perturbing its value affected the flux pathways through either succinate or isocitrate dehydrogenase. This analysis identified six additional reactions as having potentially problematic kinetic parameter values, for a total of nine: nicotinamide adenine dinucleotide (NADH) dehydrogenase, inorganic pyrophosphatase, cytosine deaminase, glutathione reductase, phosphoserine aminotransaminase, citrate synthase (Fig. 3D), succinate dehydrogenase, fumarate reductase (Fig. 2B), and isocitrate dehydrogenase (fig. S3H). A deeper review of the literature revealed that isocitrate dehydrogenase is part of a more complex control circuit, also involving glyoxylate reductase (35), which has not been completely specified. Because the full behavior of this circuit cannot

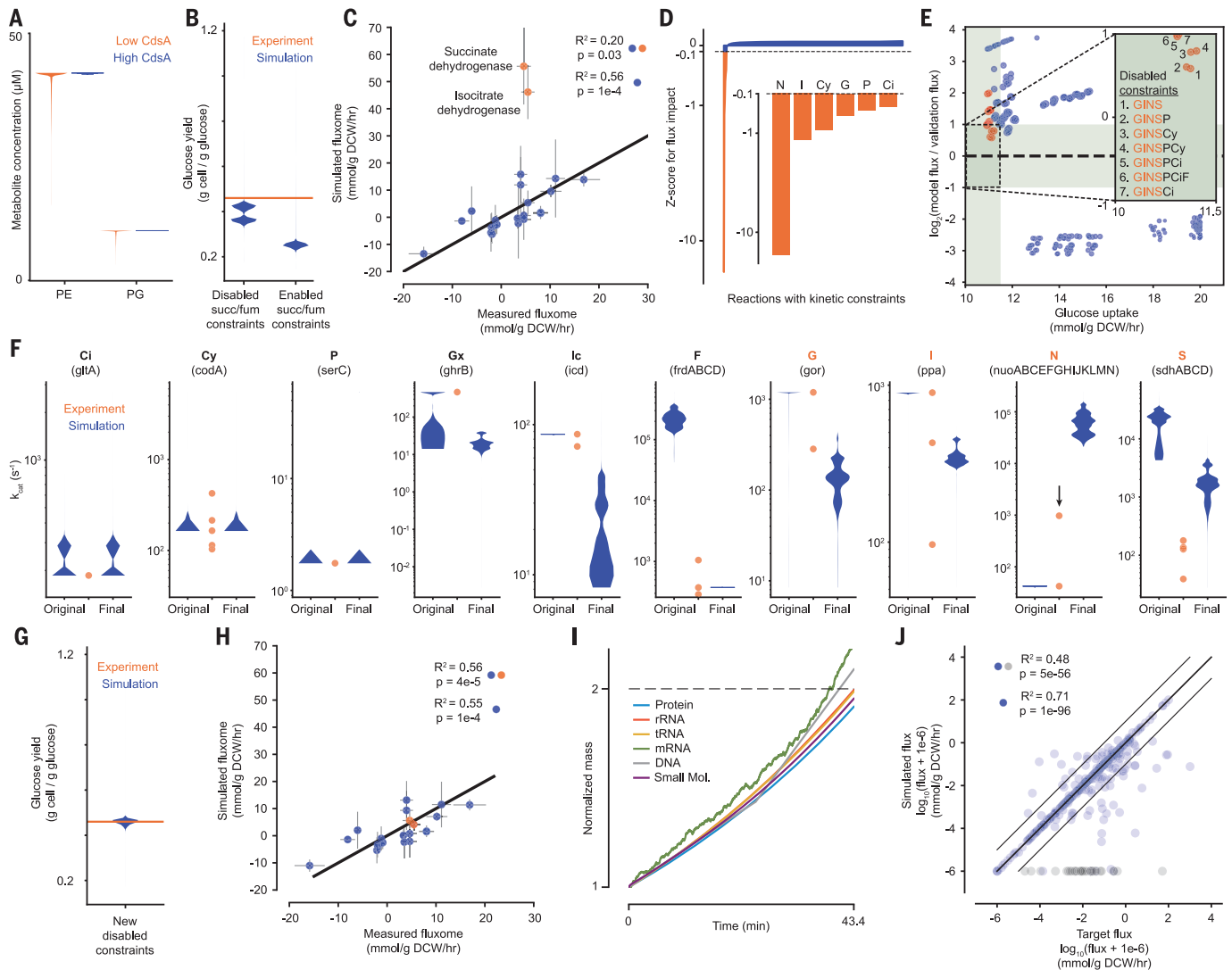
be described, the isolated kinetic constraints for these reactions were removed from the final model, leaving us with eight reactions to consider in more depth.

To determine the main and interaction effects between the eight remaining kinetic constraints, we performed a full factorial, two-level experimental design, with 256 ( $2^8$ ) sets of simulation runs. These runs simulated the result of removing or including all of our identified kinetic constraints in every possible combination. The combinations of constraints that produced simulation outputs with strong agreement with the fluxome (34), as well as with the growth yield on glucose, were always missing at least succinate dehydrogenase, NADH dehydrogenase, inorganic pyrophosphatase, and glutathione reductase (Fig. 3E), which indicated that the values for these kinetic constraints are inconsistent with the rest of the data. We therefore removed the constraints associated with these final four reactions for a new round of simulations and used the simulated fluxes and enzyme expression data to calculate a new estimated distribution for each  $k_{cat}$ . Figure 3F shows the  $k_{cat}$  distributions for all 10 of the reactions mentioned above in both our original and final model, together with kinetic parameters identified from the literature. In the cases of citrate synthase, cytosine deaminase, and phosphoserine transaminase, the distributions were similar in the original and final model and were acceptably close to measured values. The remaining cases showed stronger differences between the original and final model. We expected these differences in the cases of glyoxylate reductase and isocitrate dehydrogenase because of the complexity of these enzymes' regulation. By contrast, for the cases of fumarate reductase, glutathione reductase, and inorganic pyrophosphatase, the new  $k_{cat}$  distributions were a better reflection of the measured values. Finally, NADH dehydrogenase and succinate dehydrogenase are both membrane proteins, which are notoriously difficult to characterize kinetically. For NADH dehydrogenase in particular, a new kinetic measurement not used in the construction of the model was derived from a recent and more sensitive technology (36). The resulting  $k_{cat}$  (highlighted with an arrow in Fig. 3F) was ~23-fold higher than previous measurements and was closer to our new model  $k_{cat}$  value distribution. This supports others' assertions that the effective  $k_{cat}$  values for membrane-bound enzymes may be much higher in vivo than measurements have reported (37), and may therefore explain the discrepancies between experimental measurements and our model distributions for these enzymes. In total, we found that our new version of the model was better at matching the measured  $k_{cat}$  values (Fig. 3F), as well as the growth

yield (Fig. 3G) and fluxome (Fig. 3H), and was also able to reproducibly simulate balanced growth (Fig. 3I).

Finally, beyond our growth and fluxome comparisons, we also wanted to test the global cross-consistency of all kinetic parameters governing the activity of metabolic pathways in the overall network. One way to achieve this is by comparing the target flux values included in the kinetic component of the objective function (calculated from both the curated kinetic parameter values and the simulation's enzyme and metabolite concentrations) with the simulation output flux values. As shown in Fig. 3J, 215 of 380 fluxes were within 5% of the target flux. However, there were also 33 fluxes with values of zero in the simulations but with nonzero target values; these values reflect the fact that the model is not yet functionally complete, so the resulting metabolites would be left unused. We found that we could obtain higher, but never perfect, consistency between the flux values by increasing the weight on the kinetic component of the objective function, but this resulted in slower and less steady growth (fig. S3B). Barring the above exceptions, the strong agreement between these two sets of flux values indicates a high level of cross-consistency between the kinetic parameter values themselves.

Our third finding was that the production of cellular protein can only be met by the overall capacities of the cell (in terms of building block resources as well as cell size and mass) if most of the genes are transcribed less than once per cell cycle, including a number of essential genes. This observation was preceded by a comparison between our model simulations of protein expression with a validation dataset (>2000 points) that was also withheld from the creation and parameterization of the model (32). We found a strong correlation between the predicted and observed protein abundances at higher expression levels; at lower levels, we did not see a correlation, which can be explained by the detection limits of high-throughput mRNA and protein measurement technologies (Fig. 4A) (38). The overall correlation for protein abundances  $\geq 30$ , although significant, has an interesting consequence at the level of individual genes: We found that although many genes are transcribed multiple times as a typical cell grows and divides under these conditions (aerobic glucose minimal medium), a clear majority of the genes in *E. coli* are transcribed at a rate of less than once per cell cycle (Fig. 4, B and C). Such subgenerational gene expression has been observed both theoretically (39, 40) and experimentally (41–43), but our model led us to two insights: (i) subgenerational transcription affects >50% of the genes in *E. coli* and (ii) 72 essential genes are among those that are subgenerationally



**Fig. 3. Evaluating metabolic parameter values against each other and in the context of cellular growth.** (A) Violin plot of concentrations at each simulation time point for downstream metabolites of the reaction catalyzed by CdsA, phosphatidylethanolamine (PE), and phosphatidylglycerol (PG), when the concentration of CdsA is low (orange, indicating the original, short protein half-life) or high (blue, indicating the new, longer protein half-life) (see the main text). (B) Violin plot for glucose yield at each simulation time point for simulations with succinate dehydrogenase and fumarate reductase kinetics constraints disabled or enabled. Experimental value is 0.46 g cell/g glucose at  $\mu = 0.900 \text{ hour}^{-1}$  (52). (C) Comparison of the average fluxes from simulations with succinate dehydrogenase and fumarate reductase constraints disabled for a set of reactions in central carbon metabolism with experimental measurements (34). Orange points indicate outlier fluxes, which are discussed in more detail in the text. Correlation is shown for all data points (blue and orange) and when excluding outliers (blue). (D) Impact of individually disabling each kinetic reaction constraint on the succinate dehydrogenase flux in simulations, shown as a z-score representing the average change in flux for removing one constraint compared with the distribution of the average change in flux for removing each constraint. Constraints that have a z-score of  $< -0.1$  are highlighted in orange and shown in more detail. Highlighted reaction constraints are part of the reactions that are further explored in (E). (E) Comparison of the average metrics for simulations from a two-level full factorial design to test the effects of removing up to eight kinetic constraints of interest. Inset shows the target region where the simulated glucose uptake rate is close to the expected glucose uptake rate and simulation succinate dehydrogenase flux is within a factor of 2 of the

experimental flux (green region). Disabled constraint combinations are enumerated for each point in the target region. Orange points indicate simulations run with combinations of disabled constraints that included G, I, N, and S; blue points indicate simulations run with at least one of these constraints enabled. Abbreviations are listed below in (F). (F) Distributions of predicted  $k_{\text{cat}}$  value at each simulation time step (blue) and curated kinetic parameters (orange) for each reaction identified: citrate synthase (Ci), cytosine deaminase (Cy), phosphoserine aminotransaminase (P), glyoxylate reductase (Gx), isocitrate dehydrogenase (Ic), fumarate reductase (F), glutathione reductase (G), inorganic pyrophosphatase (I), NADH dehydrogenase (N), and succinate dehydrogenase (S). Original is from simulations without constraints for S and F; final is from simulations without constraints for Gx, Ic, G, I, N, and S. The black arrow for N indicates a newly curated  $k_{\text{cat}}$  parameter that was not used in the model. (G and H) Similar to (B) and (C) but based on data from simulations with the new set of disabled constraints. (I) Representative output from simulations with the new set of disabled constraints, showing the increase in mass (normalized to initial mass and over a single life cycle) of six key cellular mass fractions. (J) Comparison between the metabolic fluxes calculated directly from the kinetic parameters (target) and the fluxes computed by simulations with the new set of disabled constraints, as summarized by the  $R^2$  value. Gray points correspond to reactions with no simulated flux despite having a target flux. Correlations are shown for all data points (blue and gray) and with gray points excluded (blue only). Full details of the analysis required to generate this figure, as well as a pointer to the generating code, can be found in the supplementary materials, section 1.2.

transcribed (Fig. 4D) (see the supplementary materials for essentiality criteria).

How might cells survive and grow when some of their essential genetic content is not transcribed at all during a typical division cycle? One possibility is that although the mRNAs may be rarely present in the cell, the corresponding proteins and protein complexes produced are numerous and stable enough that the cell never experiences their functional absence (i.e., a period of time in which the protein is completely absent from the cell). In fact, this accounted for ~1000 of the functional protein units (including complexes and functional monomers) of subgenerationally transcribed genes in our simulations, leaving just over 1400 protein products that are completely absent from the cell at least part of the time, including 23 proteins that are considered products of essential genes (Fig. 4E) (table S4).

This result suggests that certain proteins believed to be required for cell viability are likely to be absent from single cells for at least short periods of time. In the case of an essential protein, how does the cell compensate for its temporary loss caused by very low expression rates? To answer this question, we turned to our integrative modeling framework, which enables us to investigate the loss of these proteins as part of a unified system. A representative example is 4-amino-4-deoxychorismate synthase, a heterodimeric enzyme involved in folate biosynthesis. The genes encoding this enzyme, *pabA* and *pabB*, are each transcribed with a frequency of 0.94 and 0.66 times per cell cycle, respectively (Fig. 4F), producing an average of 34 PabA proteins and 101 PabB proteins per generation. The enzyme is only active as a heterodimer (PabAB) in our model, for which the average count of active complex in our simulations is 43.3, with an SD of 35.3, and we readily observed periods of time in which no heterodimer existed (Fig. 4F, gray region). During the periods in which the PabAB dimer was completely absent, the internal pool of 5,10-dimethylene tetrahydrofolate (methylene-THF) was reduced over time; however, after a new round of *pabA* or *pabB* expression, methylene-THF was rapidly resynthesized. We further confirmed that the parameter value for the synthesis probability of *pabB* mRNA is causal for PabAB and methylene-THF depletion, because lowering the value exacerbated it (fig. S4). Supporting this proposed mechanism, others have shown that bacterial metabolite pools display a much wider dynamic range than protein concentrations and can change by 50- to 170-fold over time, including by almost complete depletion of certain metabolites (44). We conclude that internal metabolite pools, replenished by rapid enzyme kinetics, can provide a literal buffer to make cell growth robust to intermittent loss of key enzymes.

The fourth finding of this study was that the data we compiled, when considered as a unified whole, can lead to successful predictions in vitro, in this case protein half-lives. As shown in Fig. 1, the equations that govern mRNA and protein expression incorporate many types of available data and, once populated in our model, were able to successfully predict protein abundance measurements that were previously withheld from the model (Fig. 4A). Not all proteins display such consistency, however, so we performed a further analysis in which the previously withheld proteomics data (32) were also taken into account to identify and understand the causes of discrepancy for these proteins. We first noted that cells for which the entire division cycle occurs in the log or exponential phase of growth may be considered to be operating at a steady state in terms of maintaining mRNA and protein concentrations. This can be represented mathematically by setting the derivative terms in Fig. 1 to zero and substituting the solution for the mRNA concentration into the equation governing protein concentration. If the experimental data that populate these equations are consistent, then the average rate of protein production should equal the average rate of protein loss (where the loss rate includes loss by dilution as well as by degradation). This proved to largely be the case, with 85% of the production rates within an order of magnitude of the corresponding loss rate (Fig. 5A).

However, the flip side of this result is that ~15% of the protein production rates differed from the loss rates by more than an order of magnitude. In considering the cases where the production and loss rates were discrepant, we considered that one likely source of discrepancy is due to the “N-end” rule, which uses the amino acid sequence of a protein to predict its half-life (16). The N-end rule is usually accurate, but in the discrepant cases we noted, we wondered whether the rest of the data populating the model could provide a better estimate of protein half-lives. To test this hypothesis, we identified six outlier proteins from this analysis, three of which were predicted by our analysis to have longer half-lives and three that were predicted to have shorter half-lives. Measurement of the actual half-lives of these proteins experimentally confirmed that our predictions were correct (Fig. 5B). We then replaced the N-end rule-based parameter values with these new measurements (which also preserved the proteomics data as a validation dataset). This result caused us to revisit our analysis of *cdsA* expression (Figs. 2B and 3A), because the N-end rule assigns the CdsA protein a short half-life, which if incorrect could cause the simulation to have an erroneously low CdsA concentration. Our steady-state analysis supported the idea that the CdsA protein may have a longer half-life (Fig. 5A).

CdsA is a membrane protein, which makes protein extraction and traditional Western blotting difficult (45). As a result, we used immunofluorescence of overexpressed CdsA to measure the presence of protein over time, and found abundant expression of CdsA, but not RpoH (which has a short half-life; see fig. S5A), after 24 hours (Fig. 5C and fig. S5, C and D). This is consistent with a half-life on the order of 10 hours for CdsA (Fig. 5B), which was included in the finalized model. The resulting simulations (i.e., the simulations shown in Figs. 1 to 4) had a higher protein count and predicted normal growth, resolving our questions regarding *cdsA*. Our steady-state analysis thus confirmed that the N-end rule holds in most cases, but also identified the points that were most likely to be discrepant and even calculated estimates of protein decay rates that were predictive of new experimental data.

In sum, construction of a highly integrative and mechanistic mathematical model provided us with an opportunity to integrate and cross-validate a vast and heterogeneous set of data in *E. coli*, a process we now call “deep curation” to reflect the multiple layers of curation that we perform (analogous to “deep learning” and “deep sequencing”) (Fig. 1). These layers include: (i) a data layer, (ii) a layer of parameters derived from the data, (iii) a layer of equations that encapsulate the parameters and also describe the underlying biological mechanisms (which notably must also be curated from the literature), (iv) a layer that contains the unified model, and (v) a layer of the simulation output, which is executable and can be used for automated comparison with any future data that are generated. By highlighting those areas in which studies in *E. coli* contradict each other, our work suggests lines of fruitful experimental inquiry for the future that may help to resolve discrepancies, leading to both new biological insights and a more coherent understanding of this critical model organism.

We found that most of the data are in fact cross-consistent with themselves. This means that the data generated by this scientific community are reliable on the whole and may be particularly interesting given how many of these measurements were performed in vitro rather than in vivo. Moreover, the model shows that these data are capable of validatable predictions, not only on previously withheld data (fig. S2E and Figs. 3G and 4A), but also on experimental results obtained later (Fig. 5B). This strongly suggests that the model is a good representation of the overall dataset and is a starting point from which we can build toward a whole-cell model that includes many more functionalities, such as mechanisms of DNA replication initiation (46), response to nitric oxide stress (47), the formation of colonies (48), the dynamics of division site selection (49),



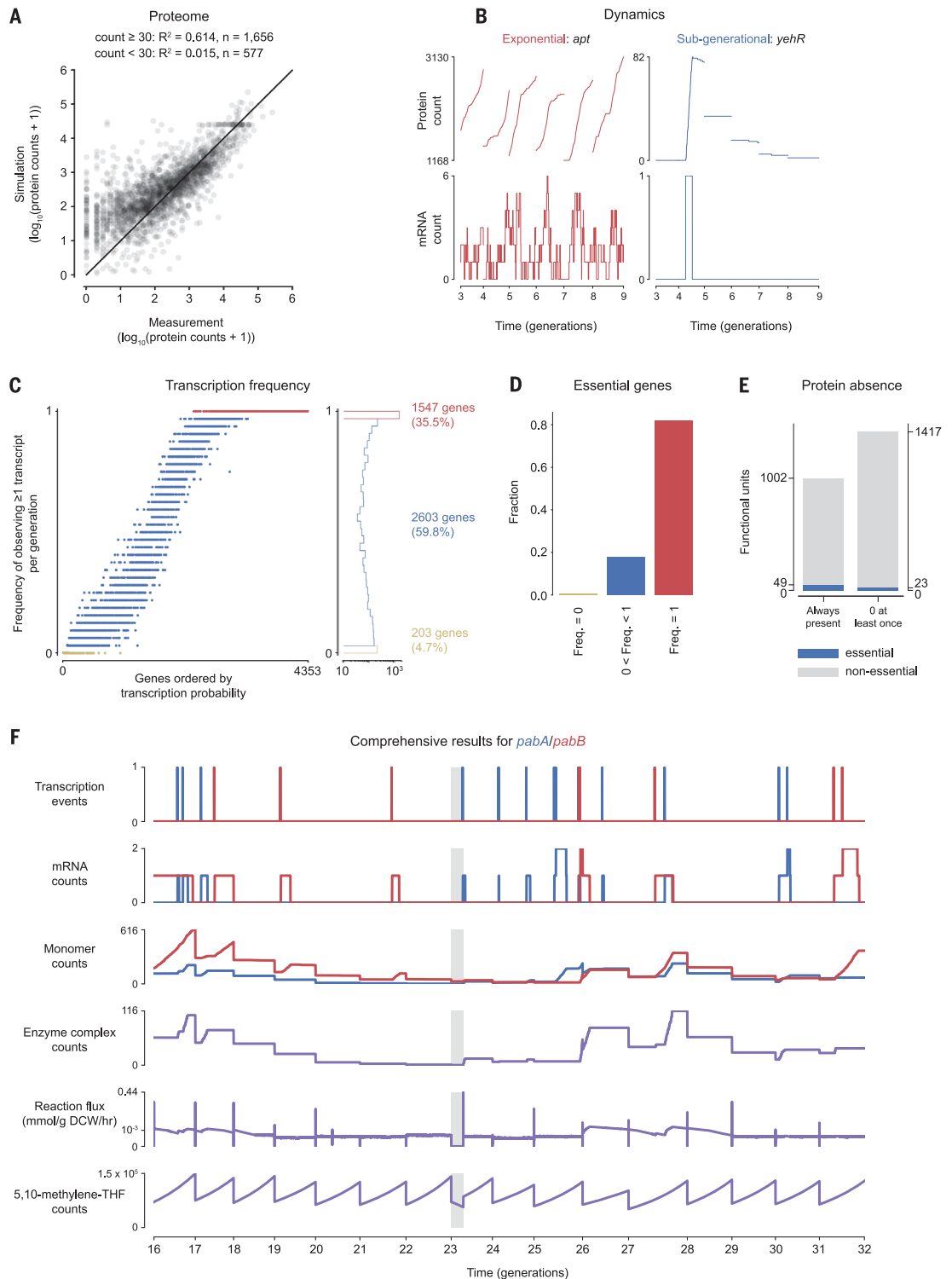
#### Fig. 4. A large fraction of *E. coli* genes are transcribed less than once per cell cycle.

**(A)** A comparison of simulation and experimental results (32) with regard to the number of proteins expressed per cell for each gene. The proteins are grouped as being highly abundant if the measured count per cell is  $\geq 30$  and otherwise lowly abundant. The  $R^2$  statistic is computed separately for each group on the log-transformed data.

**(B)** Simulations of mRNA and protein expression over multiple generations for genes that are expressed at high (left, in red) and low (right, in blue) levels of transcriptional frequencies (note that colors are conserved throughout the figure). Counts are shown for a representative six-generation long window, with an arbitrarily chosen zeroth starting generation.

**(C)** Frequency of observing at least one gene transcript per generation over a 32-generation simulation. Histograms show that 1547 genes are transcribed at least once per cell cycle (red), 203 genes are essentially never expressed in this environment (yellow), and the remaining 2603 genes are transcribed with a frequency between 0 and 1 (blue).

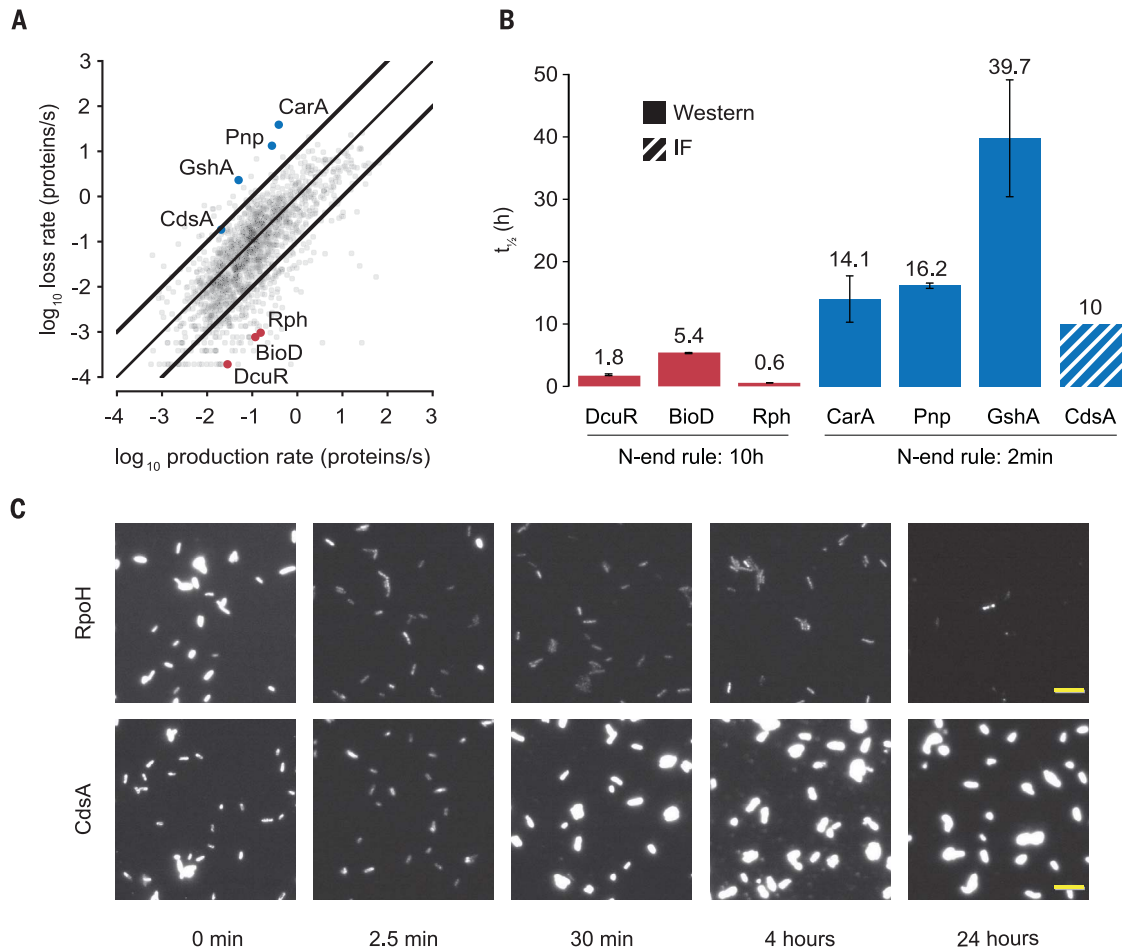
**(D)** Expression frequency analysis of known essential genes. **(E)** Division of the sub-generationally transcribed genes into those for which at least one protein is present at all times during the simulations and those for which the protein is absent for at least one time step (gray bars). Protein products of essential genes are indicated by the blue bars. Distinct protein units represent sub-generationally expressed monomers and protein complexes composed of sub-generationally expressed monomers. **(F)** Transcription, translation, complexation, and metabolic activity of the PabAB heterodimer, which catalyzes a reaction responsible for producing folates. Each new



generation is indicated with a tick mark along the x-axis; the gray area highlights a period of time in which the heterodimer is not present in the cell. All y-axes are linearly scaled except the  $[10^{-3}, 0.44]$  region of the reaction flux plot, which is log-scaled for better readability. Full details of the analysis required to generate this figure, as well as a pointer to the generating code, can be found in the supplementary materials, section 1.2.

### Fig. 5. Integrated model-data comparison leads to improved prediction of protein half-lives.

(A) Comparison of calculated protein production rates against protein loss rates for each gene. Bold lines indicate areas where the production rate and loss rate differ by more than one order of magnitude. (B) Comparison of the N-end rule with new measurements of protein half-lives for the genes highlighted in (A). The three points highlighted in red were predicted to be outliers in the steady-state analysis because their corresponding protein half-lives were much shorter than the N-end rule's prediction of 10 hours. Similarly, the proteins highlighted in blue were predicted to have much longer half-lives than the N-end rule's prediction of 2 min. Solid bars indicate half-lives that were determined by intensities on a Western blot, and the striped bar indicates an estimate (assumed from higher N-end rule value) from intensity measurements using immunofluorescence. In all seven cases, these predictions were correct. The results of control experiments (testing our protein half-life measurements against previous reports) can be found in fig. S5. (C) Images of *E. coli* MG1655 cells with either a His-tagged RpoH or CdsA plasmid that were induced for 1 hour using isopropyl- $\beta$ -D-thiogalactopyranoside followed by the addition of tetracycline to inhibit translation. At the indicated time points, aliquots of the culture were harvested and immunofluorescence was performed



and many more, all of which will in turn enable us to encapsulate many more environments and data types.

Our synthesis of heterogeneous data, along with the deep curation approach that we have described, provide a way of encapsulating and interpreting such a synthesis as a unified whole. We hope that this work, by demonstrating the value of a large-scale integrative approach with regard to understanding, interpreting, and cross-validating large datasets, will inspire further efforts to comprehensively characterize not only *E. coli* [as originally suggested by Francis Crick and Sydney Brenner (50)], but also other organisms of interest.

#### REFERENCES AND NOTES

1. Z. D. Stephens *et al.*, Big data: Astronomical or genetical? *PLoS Biol.* **13**, e1002195 (2015). doi: [10.1371/journal.pbio.1002195](https://doi.org/10.1371/journal.pbio.1002195); pmid: [26151137](https://pubmed.ncbi.nlm.nih.gov/26151137/)

2. K. Dolinski, O. G. Troyanskaya, Implications of big data for cell biology. *Mol. Biol. Cell* **26**, 2575–2578 (2015). doi: [10.1091/mbc.E13-12-0756](https://doi.org/10.1091/mbc.E13-12-0756); pmid: [26174066](https://pubmed.ncbi.nlm.nih.gov/26174066/)

3. Open Science Collaboration, PSYCHOLOGY, Estimating the reproducibility of psychological science. *Science* **349**, aac4716 (2015). doi: [10.1126/science.aac4716](https://doi.org/10.1126/science.aac4716); pmid: [26315443](https://pubmed.ncbi.nlm.nih.gov/26315443/)

4. C. G. Begley, L. M. Ellis, Drug development: Raise standards for preclinical cancer research. *Nature* **483**, 531–533 (2012). doi: [10.1038/483531a](https://doi.org/10.1038/483531a); pmid: [22460880](https://pubmed.ncbi.nlm.nih.gov/22460880/)

5. M. M. Domach, S. K. Leung, R. E. Cahn, G. G. Cocks, M. L. Shuler, Computer model for glucose-limited growth of a single cell of *Escherichia coli* B/r-A. *Biotechnol. Bioeng.* **26**, 1140 (1984). doi: [10.1002/bit.260260925](https://doi.org/10.1002/bit.260260925); pmid: [18553544](https://pubmed.ncbi.nlm.nih.gov/18553544/)

6. M. L. Shuler, P. Foley, J. Atlas, "Modeling a minimal cell," in *Microbial Systems Biology*, A. Navid, Ed. (Springer, 2012), pp. 573–610.

7. M. Tomita *et al.*, E-CELL: Software environment for whole-cell simulation. *Bioinformatics* **15**, 72–84 (1999). doi: [10.1093/bioinformatics/15.1.72](https://doi.org/10.1093/bioinformatics/15.1.72); pmid: [10068694](https://pubmed.ncbi.nlm.nih.gov/10068694/)

8. J. L. Reed, B. Ø. Palsson, Thirteen years of building constraint-based in silico models of *Escherichia coli*. *J. Bacteriol.* **185**, 2692–2699 (2003). doi: [10.1128/JB.185.9.2692-2699.2003](https://doi.org/10.1128/JB.185.9.2692-2699.2003); pmid: [12700248](https://pubmed.ncbi.nlm.nih.gov/12700248/)

9. E. Roberts, A. Magis, J. O. Ortiz, W. Baumeister, Z. Luthy-Schulten, Noise contributions in an inducible genetic switch: A whole-cell simulation study. *PLoS Comput. Biol.* **7**,

e1002010 (2011). doi: [10.1371/journal.pcbi.1002010](https://doi.org/10.1371/journal.pcbi.1002010); pmid: [21423716](https://pubmed.ncbi.nlm.nih.gov/21423716/)

10. I. Thiele, N. Jamshidi, R. M. Fleming, B. Ø. Palsson, Genome-scale reconstruction of *Escherichia coli*'s transcriptional and translational machinery: A knowledge base, its mathematical formulation, and its functional characterization. *PLoS Comput. Biol.* **5**, e1000312 (2009). doi: [10.1371/journal.pcbi.1000312](https://doi.org/10.1371/journal.pcbi.1000312); pmid: [19282977](https://pubmed.ncbi.nlm.nih.gov/19282977/)

11. J. Carrera *et al.*, An integrative, multi-scale, genome-wide model reveals the phenotypic landscape of *Escherichia coli*. *Mol. Syst. Biol.* **10**, 735 (2014). doi: [10.15252/msb.20145108](https://doi.org/10.15252/msb.20145108); pmid: [24987114](https://pubmed.ncbi.nlm.nih.gov/24987114/)

12. P. Labhsetwar, J. A. Cole, E. Roberts, N. D. Price, Z. A. Luthy-Schulten, Heterogeneity in protein expression induces metabolic variability in a modeled *Escherichia coli* population. *Proc. Natl. Acad. Sci. U.S.A.* **110**, 14006–14011 (2013). doi: [10.1073/pnas.1222569110](https://doi.org/10.1073/pnas.1222569110); pmid: [23908403](https://pubmed.ncbi.nlm.nih.gov/23908403/)

13. J. R. Karr *et al.*, A whole-cell computational model predicts phenotype from genotype. *Cell* **150**, 389–401 (2012). doi: [10.1016/j.cell.2012.05.044](https://doi.org/10.1016/j.cell.2012.05.044); pmid: [22817898](https://pubmed.ncbi.nlm.nih.gov/22817898/)

14. J. Carrera, M. W. Covert, Why build whole-cell models? *Trends Cell Biol.* **25**, 719–722 (2015). doi: [10.1016/j.tcb.2015.09.004](https://doi.org/10.1016/j.tcb.2015.09.004); pmid: [26471224](https://pubmed.ncbi.nlm.nih.gov/26471224/)

15. J. C. Sanghvi *et al.*, Accelerated discovery via a whole-cell model. *Nat. Methods* **10**, 1192–1195 (2013). doi: [10.1038/nmeth.2724](https://doi.org/10.1038/nmeth.2724); pmid: [24185838](https://pubmed.ncbi.nlm.nih.gov/24185838/)

16. A. Bachmair, D. Finley, A. Varshavsky, In vivo half-life of a protein is a function of its amino-terminal residue. *Science* **234**, 179–186 (1986). doi: [10.1126/science.3018930](https://doi.org/10.1126/science.3018930); PMID: [3018930](https://pubmed.ncbi.nlm.nih.gov/3018930/)
17. P. P. Dennis, H. Bremer, Macromolecular composition during steady-state growth of *Escherichia coli* B-r. *J. Bacteriol.* **119**, 270–281 (1974). doi: [10.1128/JB.119.1.270-281.1974](https://doi.org/10.1128/JB.119.1.270-281.1974); PMID: [4600702](https://pubmed.ncbi.nlm.nih.gov/4600702/)
18. G.-W. Li, D. Burkhardt, C. Gross, J. S. Weissman, Quantifying absolute protein synthesis rates reveals principles underlying allocation of cellular resources. *Cell* **157**, 624–635 (2014). doi: [10.1016/j.cell.2014.02.033](https://doi.org/10.1016/j.cell.2014.02.033); PMID: [24766808](https://pubmed.ncbi.nlm.nih.gov/24766808/)
19. I. M. Keseler et al., EcoCyc: Fusing model organism databases with systems biology. *Nucleic Acids Res.* **41** (D1), D605–D612 (2013). doi: [10.1093/nar/gks1027](https://doi.org/10.1093/nar/gks1027); PMID: [23143106](https://pubmed.ncbi.nlm.nih.gov/23143106/)
20. H. Bremer, P. P. Dennis, Modulation of chemical composition and other parameters of the cell at different exponential growth rates. *Ecosal Plus* **3** (2008). doi: [10.1128/ecosal.5.2.3](https://doi.org/10.1128/ecosal.5.2.3); PMID: [26443740](https://pubmed.ncbi.nlm.nih.gov/26443740/)
21. H. Bremer, P. P. Dennis, “Modulation of chemical composition and other parameters of the cell by growth rate,” in *Escherichia coli and Salmonella typhimurium: Cellular and Molecular Biology*, F. C. Neidhardt, Ed. (American Society for Microbiology, 1996); pp. 1553–1569.
22. M. Scott, C. W. Gunderson, E. M. Mateescu, Z. Zhang, T. Hwa, Interdependence of cell growth and gene expression: Origins and consequences. *Science* **330**, 1099–1102 (2010). doi: [10.1126/science.1192588](https://doi.org/10.1126/science.1192588); PMID: [21097934](https://pubmed.ncbi.nlm.nih.gov/21097934/)
23. M. Wallden, D. Fange, E. G. Lundius, Ö. Baltekin, J. Elf, The synchronization of replication and division cycles in individual *E. coli* cells. *Cell* **166**, 729–739 (2016). doi: [10.1016/j.cell.2016.06.052](https://doi.org/10.1016/j.cell.2016.06.052); PMID: [27471967](https://pubmed.ncbi.nlm.nih.gov/27471967/)
24. M. Campos et al., A constant size extension drives bacterial cell size homeostasis. *Cell* **159**, 1433–1446 (2014). doi: [10.1016/j.cell.2014.11.022](https://doi.org/10.1016/j.cell.2014.11.022); PMID: [25480302](https://pubmed.ncbi.nlm.nih.gov/25480302/)
25. J. T. Sauls, D. Li, S. Jun, Adder and a coarse-grained approach to cell size homeostasis in bacteria. *Curr. Opin. Cell Biol.* **38**, 38–44 (2016). doi: [10.1016/j.cob.2016.02.004](https://doi.org/10.1016/j.cob.2016.02.004); PMID: [26901290](https://pubmed.ncbi.nlm.nih.gov/26901290/)
26. Y. Tanouchi et al., A noisy linear map underlies oscillations in cell size and gene expression in bacteria. *Nature* **523**, 357–360 (2015). doi: [10.1038/nature14562](https://doi.org/10.1038/nature14562); PMID: [26040722](https://pubmed.ncbi.nlm.nih.gov/26040722/)
27. A. Khodayari, C. D. Maranas, A genome-scale *Escherichia coli* kinetic metabolic model k-ecoli457 satisfying flux data for multiple mutant strains. *Nat. Commun.* **7**, 13806 (2016). doi: [10.1038/ncomms13806](https://doi.org/10.1038/ncomms13806); PMID: [27996047](https://pubmed.ncbi.nlm.nih.gov/27996047/)
28. H. Kurata, Y. Sugimoto, Improved kinetic model of *Escherichia coli* central carbon metabolism in batch and continuous cultures. *J. Biosci. Bioeng.* **125**, 251–257 (2018). doi: [10.1016/j.jbiosc.2017.09.005](https://doi.org/10.1016/j.jbiosc.2017.09.005); PMID: [29054464](https://pubmed.ncbi.nlm.nih.gov/29054464/)
29. D. S. Weaver, I. M. Keseler, A. Mackie, I. T. Paulsen, P. D. Karp, A genome-scale metabolic flux model of *Escherichia coli* K-12 derived from the EcoCyc database. *BMC Syst. Biol.* **8**, 79 (2014). doi: [10.1186/1752-0509-8-79](https://doi.org/10.1186/1752-0509-8-79); PMID: [24974895](https://pubmed.ncbi.nlm.nih.gov/24974895/)
30. J. D. Orth, I. Thiele, B. Ø. Palsson, What is flux balance analysis? *Nat. Biotechnol.* **28**, 245–248 (2010). doi: [10.1038/nbt.1614](https://doi.org/10.1038/nbt.1614); PMID: [20212490](https://pubmed.ncbi.nlm.nih.gov/20212490/)
31. E. Birch, M. Udell, M. W. Covert, Incorporation of flexible objectives and time-linked simulation with flux balance analysis. *J. Theor. Biol.* **345**, 12–21 (2014). doi: [10.1016/j.jtbi.2013.12.009](https://doi.org/10.1016/j.jtbi.2013.12.009); PMID: [24361328](https://pubmed.ncbi.nlm.nih.gov/24361328/)
32. A. Schmidt et al., The quantitative and condition-dependent *Escherichia coli* proteome. *Nat. Biotechnol.* **34**, 104–110 (2016). doi: [10.1038/nbt.3418](https://doi.org/10.1038/nbt.3418); PMID: [26641532](https://pubmed.ncbi.nlm.nih.gov/26641532/)
33. T. Baba et al., Construction of *Escherichia coli* K-12 in-frame, single-gene knockout mutants: The Keio collection. *Mol. Syst. Biol.* **2** (2006). doi: [10.1038/msb4100050](https://doi.org/10.1038/msb4100050); PMID: [16738554](https://pubmed.ncbi.nlm.nih.gov/16738554/)
34. Y. Toya et al., 13C-metabolic flux analysis for batch culture of *Escherichia coli* and its Pyk and Pgi gene knockout mutants based on mass isotopomer distribution of intracellular metabolites. *Biotechnol. Prog.* **26**, 975–992 (2010). doi: [10.1002/abt.20730](https://doi.org/10.1002/abt.20730); PMID: [20730757](https://pubmed.ncbi.nlm.nih.gov/20730757/)
35. G. Shinar, J. D. Rabinowitz, U. Alon, Robustness in glyoxylate bypass regulation. *PLOS Comput. Biol.* **5**, e1000297 (2009). doi: [10.1371/journal.pcbi.1000297](https://doi.org/10.1371/journal.pcbi.1000297); PMID: [19266029](https://pubmed.ncbi.nlm.nih.gov/19266029/)
36. M. L. Verkhovskaya, N. Belevich, L. Euro, M. Wikström, M. I. Verkhovsky, Real-time electron transfer in respiratory complex I. *Proc. Natl. Acad. Sci. U.S.A.* **105**, 3763–3767 (2008). doi: [10.1073/pnas.0711249105](https://doi.org/10.1073/pnas.0711249105); PMID: [18316732](https://pubmed.ncbi.nlm.nih.gov/18316732/)
37. R. Cammack, “Assay of the enzymes with insoluble or unknown substrates: The membrane-bound quinone reductases as an example,” paper presented at ESCEC, Rudesheim/Rhein, Germany, 19–23 March 2006.
38. X. Shen et al., IonStar enables high-precision, low-missing-data proteomics quantification in large biological cohorts. *Proc. Natl. Acad. Sci. U.S.A.* **115**, E4767–E4776 (2018). doi: [10.1073/pnas.1800541115](https://doi.org/10.1073/pnas.1800541115); PMID: [29743190](https://pubmed.ncbi.nlm.nih.gov/29743190/)
39. H. Bremer, P. Dennis, M. Ehrenberg, Free RNA polymerase and modeling global transcription in *Escherichia coli*. *Biochimie* **85**, 597–609 (2003). doi: [10.1016/S0300-9084\(03\)00105-6](https://doi.org/10.1016/S0300-9084(03)00105-6); PMID: [12829377](https://pubmed.ncbi.nlm.nih.gov/12829377/)
40. J. Garcia-Bernardo, M. J. Dunlop, Tunable stochastic pulsing in the *Escherichia coli* multiple antibiotic resistance network from interlinked positive and negative feedback loops. *PLOS Comput. Biol.* **9**, e1003229 (2013). doi: [10.1371/journal.pcbi.1003229](https://doi.org/10.1371/journal.pcbi.1003229); PMID: [24086119](https://pubmed.ncbi.nlm.nih.gov/24086119/)
41. A. Bartholomäus et al., Bacteria differently regulate mRNA abundance to specifically respond to various stresses. *Philos. Trans. R. Soc. Lond. A* **374**, 20150069 (2016). doi: [10.1098/rsta.2015.0069](https://doi.org/10.1098/rsta.2015.0069); PMID: [26857681](https://pubmed.ncbi.nlm.nih.gov/26857681/)
42. Y. Taniguchi et al., coli proteome and transcriptome with single-molecule sensitivity in single cells. *Science* **329**, 533–538 (2010). doi: [10.1126/science.1188308](https://doi.org/10.1126/science.1188308); PMID: [20671182](https://pubmed.ncbi.nlm.nih.gov/20671182/)
43. I. El Meouche, Y. Siu, M. J. Dunlop, Stochastic expression of a multiple antibiotic resistance activator confers transient resistance in single cells. *Sci. Rep.* **6**, 19538 (2016). doi: [10.1038/srep19538](https://doi.org/10.1038/srep19538); PMID: [26758525](https://pubmed.ncbi.nlm.nih.gov/26758525/)
44. M. Liebeke et al., A metabolomics and proteomics study of the adaptation of *Staphylococcus aureus* to glucose starvation. *Mol. Biosyst.* **7**, 1241–1253 (2011). doi: [10.1039/c0mb000315h](https://doi.org/10.1039/c0mb000315h); PMID: [21327190](https://pubmed.ncbi.nlm.nih.gov/21327190/)
45. S.-H. Lin, G. Guidotti, “Purification of membrane proteins,” in *Methods in Enzymology: Guide to Protein Purification*, R. Burgess, M. Deutscher, Eds. (Elsevier, ed. 2, 2009), vol. 463, pp. 619–629.
46. J. C. Atlas, E. V. Nikolaev, S. T. Browning, M. L. Shuler, Incorporating genome-wide DNA sequence information into a dynamic whole-cell model of *Escherichia coli*: Application to DNA replication. *IET Syst. Biol.* **2**, 369–382 (2008). doi: [10.1049/iet-syb:20070079](https://doi.org/10.1049/iet-syb:20070079); PMID: [19045832](https://pubmed.ncbi.nlm.nih.gov/19045832/)
47. J. L. Robinson, M. P. Brynildsen, Construction and experimental validation of a quantitative kinetic model of nitric oxide stress in enterohemorrhagic *Escherichia coli* O157:H7. *Bioengineering (Basel)* **3**, 9 (2016). doi: [10.3390/bioengineering3010009](https://doi.org/10.3390/bioengineering3010009); PMID: [28952571](https://pubmed.ncbi.nlm.nih.gov/28952571/)
48. J. A. Cole, Z. Luthey-Schulten, Whole cell modeling: From single cells to colonies. *Isr. J. Chem.* **54**, 1219–1229 (2014). doi: [10.1002/ijch.201300147](https://doi.org/10.1002/ijch.201300147); PMID: [26989262](https://pubmed.ncbi.nlm.nih.gov/26989262/)
49. K. C. Huang, N. S. Wingreen, Min-protein oscillations in round bacteria. *Phys. Biol.* **1**, 229–235 (2004). doi: [10.1088/1478-3967/1/4/005](https://doi.org/10.1088/1478-3967/1/4/005); PMID: [16204843](https://pubmed.ncbi.nlm.nih.gov/16204843/)
50. F. H. C. Crick, K. Project, “The complete solution of *E. Coli*.” *Perspect. Biol. Med.* **17**, 67–70 (1973). doi: [10.1353/pbm.1973.0061](https://doi.org/10.1353/pbm.1973.0061)
51. M. W. Covert, E. M. Knight, J. L. Reed, M. J. Herrgard, B. O. Palsson, Integrating high-throughput and computational data elucidates bacterial networks. *Nature* **429**, 92–96 (2004). doi: [10.1038/nature02456](https://doi.org/10.1038/nature02456); PMID: [15129285](https://pubmed.ncbi.nlm.nih.gov/15129285/)
52. P. J. Senior, Regulation of nitrogen metabolism in *Escherichia coli* and *Klebsiella aerogenes*: Studies with the continuous-culture technique. *J. Bacteriol.* **123**, 407–418 (1975). doi: [10.1128/JB.123.2.407-418.1975](https://doi.org/10.1128/JB.123.2.407-418.1975); PMID: [238954](https://pubmed.ncbi.nlm.nih.gov/238954/)

## ACKNOWLEDGMENTS

We thank Fathom Information Design for creating the two videos and K. C. Huang, L. Willis, F. Mohammad, R. Green, R. Chang, and members of the Covert laboratory for support. **Funding:** This work was supported by the Paul G. Allen Frontiers Group through the Allen Discovery Center at Stanford; the Stanford Center for Systems Biology (NIH P50GM107615); an NIH Director’s Pioneer Award (5DP1LM01150); an Allen Distinguished Investigator award to M.W.C.; a Stanford Graduate Fellowship, DOE Computational Science Graduate Fellowship (DE-FG02-97ER25308) and Siebel Scholarship to D.N.M.; an NSF Graduate Research Fellowship to N.A.R. and M.L.P.; a Stanford School of Medicine Dean’s Postdoctoral Fellowship to J.C.; a Stanford Graduate Fellowship and NSF Graduate Research Fellowship to H.C.; an Agilent Graduate Student Fellowship and Stanford NIST JIMB Training Grant Graduate Fellowship (70NANB15H192) to J.C.M.; and an NIH grant (GM077678) to P.D.K. **Author contributions:** D.N.M., N.A.R., J.C., H.C., T.A.A., J.C.M., and M.W.C. created and analyzed the model. S.A., D.S.W., I.M.K., and P.D.K. facilitated access to the data and knowledge bases. R.K.S. and J.H.M. facilitated software development. D.N.M., M.M.D., J.C., K.L., T.E.G., and I.M. performed experiments. E.A., G.S., M.M.D., M.L.P., and S.R.B. performed further analyses. D.N.M., N.A.R., J.C., H.C., T.A.A., J.C.M., G.S., M.M.D., K.L., E.A., P.D.K., and M.W.C. wrote and edited the manuscript. **Competing interests:** M.W.C. has consulted for Google LLC. **Data and materials availability:** Sequencing data are available at GEO with accession number GSE85472. Simulation source code and documentation are available on GitHub (<https://github.com/CovertLab/WholeCellEcoliRelease>).

## SUPPLEMENTARY MATERIALS

[science.sciencemag.org/content/369/6502/eaav3751/suppl/DC1](https://science.sciencemag.org/content/369/6502/eaav3751/suppl/DC1)  
Materials and Methods

Figs. S1 to S5

Tables S1 to S9

References (53–84)

Movies S1 and S2

[View/request a protocol for this paper from Bio-protocol.](#)

11 September 2018; resubmitted 28 October 2019

Accepted 26 May 2020

10.1126/science.aav3751



## Simultaneous cross-evaluation of heterogeneous *E. coli* datasets via mechanistic simulation

Derek N. Macklin, Travis A. Ahn-Horst, Heejo Choi, Nicholas A. Ruggero, Javier Carrera, John C. Mason, Gwanggyu Sun, Eran Agmon, Mialy M. DeFelice, Inbal Maayan, Keara Lane, Ryan K. Spangler, Taryn E. Gillies, Morgan L. Paull, Sajia Akhter, Samuel R. Bray, Daniel S. Weaver, Ingrid M. Keseler, Peter D. Karp, Jerry H. Morrison, and Markus W. Covert

*Science*, **369** (6502), eaav3751.

DOI: 10.1126/science.aav3751

### Testing biochemical data by simulation

Can a bacterial cell model vet large datasets from disparate sources? Macklin *et al.* explored whether a comprehensive mathematical model can be used to verify or find conflicts in massive amounts of data that have been reported for the bacterium *Escherichia coli*, produced in thousands of papers from hundreds of labs. Although most data were consistent, there were data that could not accommodate known biological results, such as insufficient output of RNA polymerases and ribosomes to produce measured cell-doubling times. Other analyses showed that for some essential proteins, no RNA may be transcribed or translated in a cell's lifetime, but viability can be maintained without certain enzymes through a pool of stable metabolites produced earlier.

*Science*, this issue p. eaav3751

### View the article online

<https://www.science.org/doi/10.1126/science.aav3751>

### Permissions

<https://www.science.org/help/reprints-and-permissions>

Use of this article is subject to the [Terms of service](#)

---

*Science* (ISSN 1095-9203) is published by the American Association for the Advancement of Science. 1200 New York Avenue NW, Washington, DC 20005. The title *Science* is a registered trademark of AAAS.

Copyright © 2020 The Authors, some rights reserved; exclusive licensee American Association for the Advancement of Science. No claim to original U.S. Government Works