

1 A foundation model for microbial growth dynamics

2

3 Zachary A. Holmes^{1,2}, Irida Shyti^{1,2}, Alexandra L. Hoffman^{1,2,*}, Katherine E. Duncker^{1,2}, Helena
4 R. Ma^{1,2}, Zhengqing Zhou^{1,2}, Dongheon Lee^{1,2}, Rohan Maddamsetti^{1,2,†}, Kyeri Kim^{1,2}, Emrah
5 Şimşek^{1,2,¶}, Grayson S. Hamrick^{1,2}, Hyein Son^{1,2}, César A. Villalobos^{1,2}, Jia Lu^{1,2}, Yuanchi Ha^{1,2},
6 Ashwini R. Shende^{1,2}, Zhixiang Yao^{1,2}, Sizhe Liu^{1,2,‡}, Daniel M. Shapiro¹, Kseniia Kholina^{1,2},
7 Harris Davis^{1,2}, Yasa Baig^{1,+}, Feilun Wu^{1,§}, Shangying Wang^{1,||}, Xiran Wang⁴, Pranam
8 Chatterjee^{1,2,3,5}, Michael Lynch^{1,2}, Allison J. Lopatkin^{6,7,8}, Lawrence David^{1,2,9}, Emma Chory^{1,2},
9 Lingchong You^{1,2}

10

- 11 1. Department of Biomedical Engineering, Duke University, Durham, NC, USA
- 12 2. Center for Quantitative Biodesign, Duke University, Durham, NC, USA
- 13 3. Department of Computer Science, Duke University, Durham, NC, USA
- 14 4. Department of Mathematics, Duke University, Durham, NC, USA
- 15 5. Department of Biostatistics and Bioinformatics, Duke University, Durham, NC, USA
- 16 6. Department of Chemical Engineering, University of Rochester, Rochester, NY, USA
- 17 7. Department of Microbiology and Immunology, University of Rochester, Rochester, NY,
18 USA
- 19 8. Department of Biomedical Engineering, University of Rochester, Rochester, NY, USA
- 20 9. Department of Molecular Genetics and Microbiology, Duke University, Durham, NC, USA

21

22 Current affiliations

- 23 * Department of Molecular, Cellular, and Developmental Biology, Yale University, New
24 Haven, CT, USA

25 † Department of Biochemistry and Microbiology, Rutgers University, New Brunswick, NJ,
26 USA

27 ¶ Department of Physics, University of Florida, Gainesville, FL, USA

28 ‡ Department of Systems, Synthetic, and Physical Biology, Rice University, Houston, TX,
29 USA

30 † Department of Bioengineering, Stanford University, Stanford, CA, USA

31 § The Ensynble, Inc., San Francisco, CA, USA

32 || Bay Area Institute of Computation, Altos Labs, Redwood City, CA, USA.

33

34 Abstract

35 Microbial growth dynamics contain rich information about microbial populations, which support
36 applications from antibiotic testing to microbiome engineering. However, the high dimensionality
37 of growth data and the scarcity of large, task-specific datasets have limited generalizable
38 modeling analysis across systems. Here, we develop a foundation model for microbial growth
39 dynamics. It is a large-scale, self-supervised representation model trained on ~370,000
40 experimental and simulated growth curves spanning diverse microbial species, environmental
41 conditions, and community contexts. The model learns lower-dimensional latent embeddings
42 that capture essential dynamical features of raw growth data and enable accurate reconstruction
43 of these data. The concise representations enhance predictive performance in diverse
44 downstream applications. Using these embedding, we achieve few-shot learning for antibiotic
45 classification and concentration prediction, accurate forecasting of simulated and experimental
46 communities, and inference of total abundance from relative-abundance data. By extracting
47 transferable representations from heterogeneous datasets, our model provides a general
48 framework for analyzing and predicting microbial community dynamics from limited
49 measurements.

51 Introduction

52 Microbial growth dynamics are rich in information about how microbial populations behave and
53 interact. For example, growth curves contain information to both identify different strains and
54 predict antibiotic resistance of clinical isolates [1]. A single optical density curve is sufficient to
55 estimate interaction parameters between two strains [2]. Additionally, by predicting microbial
56 consortia dynamics, researchers can design clinically relevant consortia with desired metabolite
57 profiles [3]. Despite their importance, growth dynamics are difficult to generalize across
58 experimental systems and conditions. Variations in strain/species background, assay
59 configuration, and environmental conditions lead to high-dimensional and heterogeneous data
60 that exhibit alternative steady states and abrupt changes to composition [4].

61 High-dimensional data like these growth dynamics often reside on a lower-dimensional
62 manifold, because many observed features are correlated or redundant [5], [6], [7]. Recent
63 studies show that microbial growth curves can be encoded into lower-dimensional vectors using
64 variational autoencoders, and these embeddings not only allow us to reconstruct the original
65 growth curves but also serve as effective inputs for downstream analyses [8], [9]. Further,
66 microbial growth dynamics are predicted by mapping inputs, such as initial cell densities, ODE
67 parameters, or experimental conditions, to these lower-dimensional vectors [8]. Building on this
68 foundation, we sought to construct a large-scale, generalizable model that learns the universal
69 statistical structure of microbial growth, analogous to how foundation models in other biological
70 domains capture universal relationships within protein sequences [10], genomes [11], or
71 single-cell transcriptomes [12].

72

73 Here, we define a foundation model for microbial growth dynamics as a large-scale,
74 self-supervised representation model trained on heterogeneous microbial time-series data to
75 learn generalizable embeddings transferable across tasks and systems. In contrast to models

76 tailored for specific experiments, a foundation model captures shared properties of microbial
77 dynamics, potentially enabling predictions where data are sparse or difficult to standardize.
78 Similar strategies have transformed other fields. For example, the protein foundation model
79 ESM-3 is used to design novel proteins, which simulates evolution [10]. Similarly, the DNA
80 foundation model Evo is trained on genomes and is used to generate functional proteins [11].
81 Further, the single-cell RNA sequencing foundation model scBERT is trained on RNA sequence
82 databases and can be used to effectively classify different cell types in new samples [12]. We
83 reason that microbial time-series data are likely amenable to such scalable representation
84 learning.

85

86 Building a foundation model requires a sufficiently large dataset. Training specific models for
87 tasks in microbial consortia can be inhibitive based on the number of combinations and
88 availability of data [13], [14]. However, there is sufficient data available to train a broadly
89 applicable foundation model. Between high-throughput robotics and collecting previously
90 created growth curves, we curated an extensive dataset on microbial growth [15], [1], [16], [17],
91 [18], [19], [20], [21], [22], [23], [24], [25], [26]. Generating longitudinal studies of microbiota is
92 becoming easier with more access to sequencing technology, and we can collect many of these
93 datasets [4], [27], [28]. Training on broad, diverse datasets leads to better predictive outcomes
94 in downstream tasks [29]. Foundation models have seen a vast array of data, such that they can
95 then be applied to smaller datasets which may be insufficient to train a new model on their own
96 [12], [29].

97

98 In this study, we compiled ~370,000 microbial growth trajectories across multiple species,
99 growth conditions, and community contexts, including both clonal and consortial dynamics. We
100 trained a variational autoencoder on this dataset to learn a compact, low-dimensional latent
101 representation of growth curves. We show that the model accurately reconstructs growth

102 dynamics, generalizes to new datasets, and improves prediction accuracy in data-sparse
103 regimes. Specifically, the model enables few-shot learning for antibiotic classification and
104 concentration prediction, forecasting of microbial consortia dynamics, and inference of total
105 abundance from relative-abundance data. By unifying heterogeneous microbial growth
106 measurements within a shared latent framework, this work establishes a scalable,
107 self-supervised model that serves as a foundation for quantitative prediction and interpretation
108 of microbial population dynamics.

109 **Results**

110 ***Experimental data collection and preprocessing***

111 We collected 306,417 experimental growth curves to build our foundation model. These data
112 consist of:

- 113 • Time courses run in a plate reader for a duration of 12 to 40 hours [Fig 1A]
- 114 • Continuous cultures run using high-throughput robotics, where the bacteria receive a
115 continuous flow of media throughout the experiment [Fig 1B]
- 116 • Longitudinal microbiome studies where the composition is sequenced on a given
117 interval, such as daily over the course of weeks or months [Fig 1C]

118

119 While *E. coli* is the most abundant species in our collected data, the dataset also includes other
120 lab strains, clinical isolates, and natural gut and saliva microbiota, as detailed in the Methods
121 section and supplemental information. We collected both published and unpublished growth
122 data, and we provide a sample of different growth trajectories from the different sources [Fig
123 1D]. For the plate reader and continuous culture experiments, most of the data contains
124 between 100 to 150 time points, so we selected an input size of 128 time points for our model.
125 Generally, most of the dynamics occur in the first 12 to 20 hours of an experiment. In many

126 cases the experiments were at a steady state by 128 time points. For many of the microbiome
127 datasets, the data were collected on a daily basis. 128 days of growth contain different
128 dynamics than the 12 to 40 hours of growth in the plate readers, so we interpolated the data to
129 bring them to similar scales, despite the diversity of the dynamics. We perform this
130 preprocessing such that our machine-learning model will have consistent input across different
131 experiment types. Details regarding data truncation, interpolation, and preprocessing are
132 included in the Methods section.

133 ***Generation of simulated data***

134 We augmented our experimental data with the simulations of three mathematical models:

- 135 ● A modified logistic equation
- 136 ● A generalized Lotka-Volterra (gLV) model with dispersal
- 137 ● An antibiotic treatment model

138

139 The modified logistic equation [Eq 1] stems from the traditional logistic equation and includes
140 additional modifying terms to better replicate the behavior observed in experimental plate reader
141 data, as shown in our previous studies [8], [9], [16]. These modifications account for both
142 empirical growth patterns and background optical density from media, which are not addressed
143 by the standard logistic model. This adaptation enables the simulated data to more closely align
144 with experimental growth curves.

145

146 The gLV model with a dispersal term [Eq 2, 3] was adopted from Hu *et al.* [30]. It has been
147 shown to readily generate highly complex dynamics, including chaotic dynamics. We chose the
148 parameters to allow the generation of such complex dynamics.

149

150 The antibiotic treatment model [Eq 4-13], based on the work of Ma *et al.*, simulates the
151 dynamics of β -lactam resistant bacteria [16]. It incorporates interactions between resistant and
152 sensitive bacterial populations, as well as the effects of β -lactam antibiotics and β -lactamase
153 inhibitors.

154

155 We provide a sample of different growth trajectories from these simulations [Fig 1E].

156 ***Training and evaluating models for universal representation of microbial growth*** 157 ***dynamics***

158 Previous work in our lab determined that microbial growth dynamics can be encoded in a
159 lower-dimensional space using machine learning models, such as autoencoders [8]. Here, we
160 developed a foundation model to learn latent representations of microbial growth dynamics. We
161 considered different model architectures, such as custom variational autoencoders, custom
162 transformers, and principal components analysis [Supp Info]. Additionally, we tested the impact
163 of different latent dimensions, ranging from 2 to 32. Based on the results from our analysis, we
164 selected a single model to focus on for the remainder of the paper: the custom variational
165 autoencoder, Autoencoder7X (A7X) with 8 latent dimensions.

166

167 A7X encodes 128-time-point curves into an 8-dimensional latent vector, and then decodes the
168 latent vector to reconstruct the original curve [Fig 2A]. Training a variational autoencoder
169 involves computing the mean squared error for the growth curve reconstruction and the
170 Kullback-Leibler divergence loss to enforce a normal latent space distribution. During training,
171 random noise is applied to the 8-dimensional latent vector, with the KL divergence ensuring a
172 Gaussian distribution.

173

174 We trained our models on a large, diverse dataset: after preprocessing the input data, we used
175 306,417 experimental and 69,674 simulated growth curves. We decided to supplement the
176 experimental with simulated curves, because including the simulated curves allowed for higher
177 reconstruction of the experimental data for most of the latent dimension sizes [Supp Fig S1]. We
178 trained the model using a training set that consists of 80% of the data for the project (300,872
179 curves), and 20% was withheld for testing (75,219 curves). Using A7X, the reconstruction of the
180 training data has a coefficient of determination of 0.991, while the testing data is 0.988 [Fig 2B].
181 The model can accurately reconstruct a broad array of the data [Fig 2C]. We have included
182 more examples in the supplemental material [Supp Fig S2].

183

184 We visualize this distribution by encoding all of the curves in our collection and applying
185 t-distributed stochastic neighbor embedding to present the vectors in a 2D display. For each of
186 the four overarching data types - clonal simulation, consortia simulation, clonal experiment, and
187 consortia experiment - we visualized where the curves from those data types exist in the 2D
188 display [Fig 2D].

189

190 We find that different model configurations can perform better on different tasks, so the best
191 architecture can depend on a specific application [Fig 2E]. We chose to use a variational
192 autoencoder, because we want accurate reconstructions of growth curves. To select the number
193 of latent dimensions, we used both established statistical methods for estimating intrinsic
194 dimensions, as well as empirical results from our training process. Using the Maximum
195 Likelihood Estimator (MLE) proposed by Levina and Bickel, we estimate the intrinsic dimension
196 of our entire dataset to be 5.6 [7]. As discussed in their paper, estimating intrinsic dimensions is
197 not exact, and the estimated dimension can vary depending on the sample size and parameters
198 used [7]. The estimated intrinsic dimensions vary greatly depending on the complexity of the
199 data, and we see estimates of 1 to 2 for simple growth curves and 8 to 10 for complex or noisy

200 datasets. Therefore, we balance this overall estimate with our empirical observations after
201 testing various dimensions across different models. We selected 8 dimensions based on the
202 accuracy of reconstruction. Latent dimensions below 8 significantly reduced the reconstruction
203 accuracy, while latent dimensions greater than 8 did not add additional benefit to the
204 reconstruction [Supp Table ST5]. The accuracy of the reconstruction is higher when the data
205 come from a less complex dataset, and the accuracy is lower when the data come from a more
206 complex dataset [Supp Fig S3].

207 ***Classifying antibiotic treatment and predicting antibiotic concentration***

208 Few-shot learning is an approach that takes a learned embedding from a foundation model and
209 applies it to a specific problem set [12], [31]. Few-shot learning is especially relevant in microbial
210 biology, because data sparsity is a common problem [13], [14]. Therefore, it is useful to be able
211 to build predictive models using small datasets. We hypothesized that at lower training sizes,
212 the foundation model representation of the curves would outperform the raw growth data. We
213 hypothesized this, because the foundation model reduces the curves to a lower dimension
214 manifold that both reduces noise in the experimental data and emphasizes the key
215 characteristics of the growth curves.

216

217 To explore this application, we used a dataset generated by antibiotic resistance work in our lab.
218 This dataset consists of 6,924 growth curves. 672 came from experiments growing *E. coli*
219 TOP10F in different media conditions and initial cell densities [Supp Fig S4]. 5,100 came from
220 experiments growing *E. coli* Keio strains in different antibiotic treatments [Supp Fig S5]. 1,152
221 came from *E. coli* MG1655 strains grown in different antibiotic treatments [Supp Fig S6]. In our
222 first application, our objective was to classify the type of antibiotic used. In our second
223 application, we predicted the concentration of the antibiotic. For each task, we used a single test

224 dataset, but we adjusted the size of the training dataset to compare the impact of data sparsity
225 on the ability to accurately classify with both the raw curves and the latent vectors.

226

227 We trained a classification model to predict which antibiotic was used based on the growth
228 curve data [Fig 3A]. As an input, we used the OD600 curves from the data. We extracted the
229 latent vectors by encoding the curves with Autoencoder7X (A7X). We trained two separate
230 classification models using either the raw curves or the latent vectors as input, and we predicted
231 the antibiotic used. The latent vectors are normalized to a maximum value of one before
232 encoding, so the raw maximum value of the growth curve is appended to the latent vector to
233 give an input vector of 9 dimensions. For this task, we used 4,890 growth curves. These were
234 the OD600 readings for the experimental conditions that had antibiotics added. We used 20%,
235 978 curves, as the test dataset. For training, we used anywhere from 1% to 80% of the data, 48
236 to 3,912 curves. We did this with fivefold cross-validation, such that each curve was part of the
237 test dataset in one of the folds. The classification accuracy for the classifier using the latent
238 vectors was higher than the raw data at every training size, and they got closer together with
239 more training data [Fig 3B]. This means that with less data the latent vectors provided more
240 accurate classifications than the raw data.

241

242 We trained a regression model to predict the concentration of antibiotic used in an experiment
243 based on the growth curve data [Fig 3C]. As an input, we used the OD600 curves from the data.
244 We extracted the latent vectors by encoding the curves with A7X. We trained two separate
245 regression models using either the raw curves or the latent vectors as input, and we predicted
246 the antibiotic used. For this task, we used 6,156 growth curves. There were more curves in this
247 task, because we included conditions with no antibiotics added. We used 20%, 1,232 curves, as
248 the test dataset. For training, we used anywhere from 1% to 80% of the data, 61 to 4,924
249 curves. We did this with fivefold cross-validation, such that each curve was part of the test

250 dataset in one of the folds. The prediction accuracy was higher using the latent vectors than the
251 raw data at every training size, and they got closer together with more training data [Fig 3D].
252 This was consistent with our hypothesis that less training data will benefit more from the use of
253 the foundation model, but large enough datasets will be sufficient for training the regression
254 model without latent representations. Additionally, we compared the performance of different
255 model architectures and latent sizes on the antibiotic prediction tasks [Supp Fig S8, Supp Table
256 ST6, Supp Table ST7].

257 ***Classifying antibiotic resistance from cell density curves***

258 Bacterial growth curves contain sufficient information to predict phenotypes such as antibiotic
259 resistance to ampicillin–sulbactam (SAM), ciprofloxacin (CIP), trimethoprim–sulfamethoxazole
260 (SXT), and gentamicin (GM), as shown in previous work from our lab [1]. Predicting these
261 phenotypes from growth data is meaningful, because 1) genotypes do not have a direct linkage
262 to phenotypes, so WGS is insufficient for making these predictions and 2) we can use a simple
263 OD600 growth curve which does not require testing different antibiotics individually [1]. In
264 previous work from our lab, Baig *et al.* trained a variational autoencoder end-to-end on a subset
265 of these curves, and by embedding the derivatives of the curves they were able to classify the
266 resistance more accurately than with the raw data [8]. Following their work, we performed a
267 similar task of classifying the resistance of bacteria based solely on their OD600 growth curves.
268 Instead of taking the derivative of the curve, we used the blank-corrected data as the input.

269

270 There were 4,432 growth curves in this dataset [Supp Fig S7]. The dataset consisted of 277
271 clinical isolates with 4 growth conditions, and each experiment had 4 replicates. Of the 277
272 strains, 42 are resistant to none of the antibiotics, 47 are resistant to 1, 114 are resistant to 2, 45
273 are resistant to 3, and 29 are resistant to all 4. We used a training size of 80% of the data
274 randomly separated, so the training dataset was 3,324 curves and the test dataset was 1,108

275 curves. We did this with fivefold cross-validation, such that each curve was part of the test
276 dataset in one of the folds. We tested 2 different input configurations to predict the resistance:
277 raw 128-time-point curves and 8-dimensional latent vector from Autoencoder7X (A7X)
278 foundation model encoding [Fig 3E]. The latent vectors are normalized to a maximum value of
279 one before encoding, so the raw maximum value of the growth curve is appended to the latent
280 vector to give an input vector of 9 dimensions. No information about the strains or growth
281 conditions were included in the model.

282

283 For each of four types of antibiotics, we ran our training pipeline to classify whether a given
284 strain was resistant or sensitive to the antibiotic. The results are reported as the F1 score, which
285 is a balance of precision, the ability to avoid false negatives, and recall, the ability to identify all
286 of the positives. We found that for two of the antibiotics, SAM and CIP, the latent vectors
287 outperformed the raw data [Fig 3F]. For the other two antibiotics, SXT and GM, the latent
288 vectors and raw data are comparable [Fig 3F]. The foundation model provides the benefit of
289 having seen many more growth curves than what is available in this study, and this is why its
290 latent vectors are able to meaningfully represent the growth curves in this sparse dataset.
291 Additionally, we compared the performance of different model architectures and latent sizes on
292 the antibiotic prediction tasks [Supp Fig S8, Supp Table ST8, Supp Table ST9, Supp Table
293 ST10, Supp Table ST11].

294 ***Forecasting simulated microbial community dynamics***

295 Microbial consortia are found everywhere in nature: from human saliva, skin, and guts to soil or
296 the ocean [4], [27], [28]. Researchers often want to predict the behavior of microbial consortia to
297 understand the impact on things like metabolite production in the gut [3]. While understanding
298 how an entire microbial consortia operates can be important, in many cases we are only
299 interested in a few key members of the consortia. For example, when examining complex

300 microbial consortia, we typically focus on a few key strains, such as the infectious gut pathogen
301 *Clostridioides difficile* and the inhibitory strains *Clostridium scindens* and *Clostridium hiranonis*
302 [32]. In this case, we can modify our dataset to focus on the key members of a consortium and
303 omit the background strains and consider their influence as background noise. Using our
304 foundation model, we can represent this community in a lower dimensional vector and predict its
305 dynamics and behavior into the future.

306

307 For this application, we are using a dataset which consists of 150 different consortia [Supp Fig
308 S9]. These consortia were generated using two different variations of the generalized
309 Lotka-Volterra model. The first is the dispersal variation [Eq 2, 3] used by Hu *et al.* [30]. The
310 second is the bounded version [Eq 14] published from our lab [33]. These models result in both
311 stable and chaotic behavior. Each of the consortia was originally run with 20 to 100 members in
312 the community. From these original simulations, 2 to 10 focal members were selected for the
313 analysis. Each consortia was simulated 10,000 times, unless stated otherwise, and we split the
314 data randomly into 8,000 curves for the full training set and 2,000 curves for the test set. In each
315 simulation, the initial quantities of the focal populations were randomized. For some of the
316 consortia, the background communities had the same initial quantities in each simulation, and in
317 others they were randomized. From the 8,000 curves in the training set, we randomly selected
318 smaller training sizes of 20, 100, 200, 1,000, and 4,000 curves. On each of our 6 training
319 datasets, we performed fivefold cross-validation to train 5 separate models, and we report the
320 mean values from the results of these models. All final metrics are measured on the original test
321 set with 2,000 curves.

322

323 To demonstrate the utility of our model, we selected two of these consortia to analyze. The first
324 community is our simple community with 3.0 estimated intrinsic dimensions. This consortium
325 was generated using 20 background strains, and we randomly selected 5 members as the focal

326 community. The second community is our complex community with 6.0 estimated intrinsic
327 dimensions. This consortium was generated using 100 background strains, and we randomly
328 selected 8 members as the focal community.

329

330 We hypothesize that these curves have sufficient content to predict future behavior of the
331 consortium. We designed a pipeline to train a regression model to use 128 time points to predict
332 the 128 time points [Fig 4A]. For each of the input configurations, the curves were stacked such
333 that the input to the regression model with M communities \times 128 or 8, respectively. We tested 2
334 different input configurations.

335 1) Raw 128-time-point curves

336 2) 8-dimensional latent vector from Autoencoder7X foundation model encoding plus
337 maximum value of raw curve for a total size of 9 dimensions

338

339 Inspired by the success of GraphCast and Lam *et al.* in forecasting weather [34], we use our
340 regression models to forecast microbial consortia behavior into the future. We use an initial seed
341 of 128 time points, and then we predict the next 128 time points by using the latent vector as the
342 target value. We then append the predicted results to the initial seed, and slide the window to
343 use the updated prediction as the input for the next segment [Fig 4B].

344

345 For the simple consortium, we are able to accurately forecast into the future [Fig 4C]. To
346 demonstrate the utility of the foundation model at lower training data sizes, we used the training
347 dataset with 800 curves. To generate the accuracy metrics, we used each of the 5 models from
348 the fivefold cross-validation and used them to forecast the 2,000 curves from the test dataset.
349 The reported metrics are the mean values of these 5 forecasts. The accuracy of the prediction is
350 high with a coefficient of determination above 0.995 and RMSE below 0.022 for both raw curves
351 and latent vectors. For the complex consortium, we forecast all the members into the future [Fig

352 4D]. The forecast accuracy decreases over time, but the latent vectors provide better accuracy
353 further into the future compared to the raw curves as inputs. We have included more examples
354 in the supplemental material [Supp Fig S10].

355 ***Forecasting experimental microbial community dynamics***

356 While predicting simulated consortia is a useful exercise, being able to predict real-world
357 microbial consortia is a much more meaningful task. For this task, we use the microbiome
358 experiments completed by Fujita *et al.* on their natural microbial communities [4]. This dataset is
359 useful in our analysis, because they have 48 unique experiments. The experiments are
360 completed using 2 different source microbiomes in 3 different media types with 8 replicates each
361 [Fig 5A].

362

363 One challenge with collecting microbiome data is the quality of the data. Species and members
364 in a consortia that are lower in abundance are harder to measure, so we only use members that
365 are above a certain abundance for this task. For each of the 6 experimental conditions, there
366 were anywhere from 6 to 77 strains present in the experiment. For each condition, we selected
367 2 to 18 members that were present in all 8 of the replicates [Supp Fig S11]. We treat these as
368 our focal communities within the larger consortia. We applied eightfold cross-validation and
369 used each of the replicates as the test set for each round of validation [Fig 5B]. This allows us to
370 measure the average performance of prediction on the community instead of relying on a single
371 sample.

372

373 Similar to the simulated consortia, we designed a pipeline to train a regression model to predict
374 the next segment of points based on the previous 128 time points [Fig 5C]. For each of the
375 configurations, the curves were stacked such that the input to the regression model with M
376 communities \times 128 or 9, respectively. We tested 2 different configurations.

377 1) Input of latent vector plus maximum raw value to predict latent vector plus maximum raw
378 value of the next 128 time points

379 2) Input of raw curve to predict raw curve of the next 128 time points

380

381 As with the simulated consortia, we forecasted the behavior of the focal communities [Fig 5D].

382 We use an initial seed of 128 time points, and then we predict the next 128 time points.

383 Forecasting the experimental consortia is more difficult than the simulated consortia, and we

384 find that despite the latent to latent prediction outperforming the raw to raw prediction, the

385 minimum RMSE value across all the experiments is 0.177 [Fig 5E]. We have all examples in the

386 supplemental material [Supp Fig S12].

387

388 We expand this application to a task that is useful in analyzing microbial consortia: predicting

389 abundance in the future. For all the communities, we used the first 128 time points to predict

390 which species would have an abundance above 10% at the final time point. In other words, we

391 use the first 14 days worth of interpolated data to predict the species abundance 96 days later

392 [Fig 5F]. We find that the latent input outperforms the raw input in 5 of the 6 experimental

393 conditions [Fig 5G].

394 ***Predicting absolute abundances from relative abundances***

395 In microbiome research, we are generally interested in understanding the composition of

396 microbial consortia [35]. A common way to derive composition is to use amplicon sequencing,

397 such as 16S rRNA, which provides relative abundance of different members of a consortium

398 [36]. However, these measurements lack information about total microbial abundance, which is

399 critical for understanding ecological dynamics, host-microbe interactions, and community

400 stability [35], [37], [38]. We hypothesize that temporal patterns in relative abundance data

401 contain sufficient information to infer total abundance dynamics. Furthermore, we propose that

402 foundation models trained on diverse microbial growth dynamics could extract features from
403 these temporal patterns, and these features enable more accurate predictions than using raw
404 relative abundance data alone.

405

406 To test our hypothesis, we generated two distinct synthetic datasets, each containing 10,000
407 simulated microbial communities with 20 species tracked over 128 time points. The first dataset
408 modeled community dynamics using a gLV variant that incorporated logistic growth, saturating
409 positive interactions, and density dependent negative interactions [33]. The second dataset was
410 based on a classical generalized Lotka-Volterra model with strong intraspecific competition and a
411 weak uniform dispersal term, designed to capture more complex and potentially chaotic
412 dynamics. Both models used randomly sampled interaction parameters to create diverse
413 community behaviors, with initial conditions drawn from a Dirichlet distribution to ensure varied
414 compositional starting points at low total biomass. From each absolute abundance trajectory
415 generated by these models, we computed relative abundances by normalizing by the total
416 community abundance at each time point, creating paired datasets where the task was to
417 predict total abundance from relative abundance profiles.

418

419 For each sample, we used the complete trajectory of relative abundances across all species (20
420 species \times 128 time points) as input to predict the corresponding total abundance at each time
421 point. We compared two input representations:

- 422 1) Raw 128-time-point curves of relative abundance data for each of the 20 members [Fig
423 6B]
- 424 2) 8-dimensional latent vector from Autoencoder7X foundation model encoding plus
425 maximum value of raw curve for a total size of 9 dimensions for each of the 20 members
426 [Fig 6A]

427

428 We implemented multi-layer perceptron (MLP) regression models to predict total abundance
429 from the input features. Both models directly predicted 128-dimensional output vectors
430 corresponding to total abundance at each time point, using mean squared error as the loss
431 function. Critically, for the latent representation approach, we bypassed the foundation model's
432 decoder: the MLP predicted total abundance directly from the compressed latent features
433 without reconstructing back to the original data space. This design allowed us to test whether
434 the latent representations captured dynamics-relevant information sufficient for prediction,
435 independent of reconstruction accuracy. We trained both models using fivefold cross-validation
436 across varying training set sizes (1% to 80% of data), implementing early stopping based on
437 validation loss to prevent overfitting [Fig 6A, Fig 6B].

438

439 The results demonstrate that foundation model latent representations consistently outperform
440 raw relative abundance data across both dynamical regimes. For the GLV without dispersal
441 dynamics, the latent representation achieved superior performance across all training set sizes
442 [Fig 6E]. The GLV with dispersal dynamics showed the same patterns [Fig 6C]. The foundation
443 model's compressed latent space (9 dimensions) capture more generalizable dynamical
444 features that are transferred effectively to the prediction task compared to the high-dimensional
445 raw input space. We observe similar trends for RMSE: latent features consistently outperform
446 raw data for both GLV with dispersal [Fig 6D] and GLV without dispersal [Fig 6F]. This finding
447 supports our hypothesis that foundation models pre-trained on diverse microbial dynamics
448 extract robust, transferable features that enable better generalization than raw temporal
449 patterns, with consistent advantages across both simple and complex dynamical regimes.

450 Discussion

451 Microbial growth trajectories are among the most fundamental and widely measured descriptors
452 of population behavior. Yet, despite decades of data, the lack of a unified analytical framework

453 has limited our ability to compare, interpret, or predict growth dynamics across species,
454 environments, and community contexts. Here, we introduce a foundation model for microbial
455 growth dynamics. This model is a large-scale, self-supervised representation model that learns
456 the shared statistical structure of microbial time series from heterogeneous datasets. The
457 resulting latent representations capture key dynamical features, which enable accurate
458 reconstruction, transfer across experimental systems, and improve performance in downstream
459 predictive tasks.

460

461 Our model builds on the recognition that microbial dynamics, while high-dimensional in their raw
462 form, lie on a low-dimensional manifold. Previous studies have established this empirically
463 within the context of specific systems and conditions [8], [9]. Our current study extends that
464 insight to a far broader scale, integrating hundreds of thousands of experimental and simulated
465 trajectories to reveal the universal structure underlying microbial population dynamics. This
466 generalization represents a conceptual advance: rather than fitting a single mechanistic model
467 to each system, we train a single data-driven model that captures the statistical essence of
468 microbial growth across systems and conditions.

469

470 Similar to protein or transcriptomic foundation models, which learn universal embeddings from
471 large-scale biological data, our model distills microbial time-series behavior into compact latent
472 variables that can be reused for diverse applications, including classification, forecasting, and
473 inference, from limited or noisy measurements.

474

475 The success of this framework reflects both biological and statistical regularities. Biologically,
476 microbial populations often exhibit stereotyped phases of lag, exponential growth, and
477 saturation, which constrain the effective dimensionality of the system. Statistically, these
478 constraints allow a deep encoder-decoder architecture to learn transferable representations that

479 summarize population-level behavior. The emergent low-dimensional structure suggests that
480 microbial growth, despite the complexity of underlying physiology, can be effectively
481 characterized by a small set of latent variables.

482

483 The model, along with the training, can be revised and expanded in multiple ways. More data
484 can be included in the training process. This would involve sourcing from across all publications
485 to include as much data as possible, and continuing to collect data from our colleagues and
486 collaborators. Different model architectures could allow for more flexibility, or may perform better
487 in different tasks. The requirement for 128 time points requires truncation or interpolation, which
488 is one potential limitation in this approach. Designing a model to take an input of any size would
489 allow for more flexible applications. Additionally, actual time between samples is not a focus in
490 this study. Future versions of this model should include a time component that allows the model
491 to interpret specific time points from experiments. This could overcome some of the limitations
492 we saw with trying to represent experiments with 5-minute gaps between data in the same
493 model as daily readings from the gut. Finally, future versions could incorporate categorical data,
494 such as strain information, media type, and experimental conditions, so the model could learn
495 how these known factors influence growth.

496

497 Beyond microbiology, our work highlights the broader potential of large-scale, self-supervised
498 learning to extract interpretable and generalizable representations of biological dynamics. The
499 ability to encode population behavior into transferable latent features provides a bridge between
500 data-driven discovery and mechanistic modeling, complementing rather than replacing explicit
501 models of growth, metabolism, and ecological interaction. By establishing a unifying latent
502 framework for microbial time series, our study offers a conceptual and practical foundation for
503 predictive microbiology, where data collected in one context can inform inference, design, and
504 control in another.

505 **Methods**

506 ***Plate reader experiments***

507 Plate readers are used to measure optical density and fluorescence over time. Optical density
508 corresponds with the biomass of the cells in the plate wells. Fluorescence is a combination of
509 fluorescence output per cell and the number of fluorescent cells. The general protocol is as
510 follows. Cells were inoculated the night before the plate reader experiment began and grew to
511 stationary phase. The next morning, the cells were diluted and added to media and chemical
512 combinations according to the specific experimental protocol. These mixtures were added to
513 plates, which were placed in the plate reader. The plate reader maintained an experimental
514 temperature, such as 37°C. Every 5 to 10 minutes, as specified by the experiment, the plates
515 were shaken and then measurements were taken. This resulted in growth curves as measured
516 by optical density and fluorescence. The exact experimental conditions, including media,
517 treatments, temperatures, inoculum densities, etc., varied depending on the data source, and
518 further details can be found in the Supplemental Information.

519 ***Continuous culture experiments***

520 Continuous culture experiments were completed using liquid handling robots. There were two
521 methods used for these experiments: turbidostats or basic kinetic experiments. The
522 high-throughput turbidostats were bacterial wells continuously fed with media. The real-time
523 feedback allowed the robot to adjust the flow to the turbidostat based on the measurements of
524 the growth of the cells, which enabled calculating the growth rate of the cells. A detailed
525 protocol on the implementation of this system can be found in the publication from Chory *et al.*
526 [39].

527 ***Microbiome experiments***

528 Microbiome experiments are conducted by sampling microbiome communities and running
529 sequencing to quantify the abundance of each of the members in the community. We used four
530 different collections of microbiome data in our study.

- 531 1) Microbiomes were collected from natural samples and experimented with different media
532 conditions [4]. We used the data directly from the published article.
- 533 2) Gut and saliva samples [27].
- 534 3) Our lab has previously developed a collection of bar-coded Keio strains, and we used
535 experimental data from communities of these strains [33], [40].

536 ***Experimental data preprocessing***

537 All original growth curves are available in the supplemental data. While many of the growth
538 curves were previously published, this study also includes new growth curves generated under
539 various configurations, as described in the supplemental material [Supp Table ST12].

540

541 For plate reader and continuous culture datasets with more than 128 time points, we truncated
542 the data after 128. For experiments with less than 128 time points, we interpolated to fill in the
543 missing data. We interpolated time courses from the plate reader using linear interpolation.

544

545 Microbiome datasets presented additional challenges due to their longer timescales. To address
546 this, we segmented the microbiome data into 16-day chunks. Missing data points were
547 interpolated using a Gaussian Process instead of linear interpolation, as Gaussian Processes
548 produced smoother, organic shapes that better captured the expected trends in these datasets.

549

550 The data were combined into a single dataset of 376,091 curves. The data were split into a
551 training set of 80% of the curves and a test set of 20% of the curves. We split the data by
552 iterating through the dataset and randomly splitting every 1,000 curves into 800 for training and
553 200 for test, so that each dataset was properly represented in each group.

554 ***Simulation using a modified logistic equation***

555 Our lab developed a modified logistic equation to address the limitations of the traditional logistic
556 equation, which does not adequately replicate growth curves from plate readers [8], [16]. Plate
557 readers are subject to background optical density from the media as well as the secondary
558 growth stages. To better fit the logistic function to experimental data, we introduced two
559 modification constants: k and θ .

560 In the following equation, the variables are as follows:

- 561 • n is the cell density
- 562 • t_{lag} is the lag time
- 563 • α is the cell growth rate
- 564 • n_{max} is the carrying capacity
- 565 • k is a modification constant
- 566 • θ is a modification constant

567 Simulation parameters are provided in the supplemental data [Supp Table ST1].

$$\frac{dn}{dt} = \begin{cases} t < t_{lag}, 0 \\ t \geq t_{lag}, \frac{\alpha}{1 + \left(\frac{n}{kn_{max}}\right)^\theta} \left(1 - \frac{n}{n_{max}}\right) \end{cases} \quad (1)$$

568

569 **Simulating community dynamics**

570 To generate dynamics with varying complexity (including chaotic dynamics), we used a
571 generalized Lotka-Volterra model with dispersal, as published by Hu *et al.* [30]. The model was
572 chosen for its ability to capture complex, dynamic interactions within microbial communities.
573 Specifically, we focus on chaotic interactions, because they provide more examples of cells with
574 changing cell density, as opposed to the examples with stable steady states.

575 In the following equations, the variables are as follows:

- 576 • n is the cell density
- 577 • A is the interaction matrix
- 578 • D is the dispersal rate
- 579 • a_{ij} is the interaction of species i with species j
- 580 • M is the number of members in the population

581 Simulation parameters are provided in the supplemental data [Supp Table ST2].

$$\frac{dn}{dt} = n(1 - An) + D \quad (2)$$

$$A = \begin{bmatrix} 1 & a_{0,1} & \dots & a_{0,M} \\ a_{1,0} & 1 & \dots & a_{1,M} \\ \vdots & \vdots & \ddots & \vdots \\ a_{M,0} & a_{M,1} & \dots & 1 \end{bmatrix} \quad (3)$$

582

583 **Simulating growth dynamics in response to β -lactam treatment**

584 To simulate antibiotic resistance, we used the β -lactam resistance model described by Ma *et al.*
585 [16]. This model employs ordinary differential equations to capture the dynamics between
586 resistant and sensitive cells and their response to β -lactam antibiotics and β -lactamase
587 inhibitors.

588 In the following equations, the variables are as follows:

- 589 • n_s sensitive cell density
- 590 • n_r resistant cell density
- 591 • s nutrient concentration
- 592 • a antibiotic concentration
- 593 • b β -lactamase concentration
- 594 • g growth rate
- 595 • l lysis rate of cells
- 596 • γ lysis rate by antibiotic
- 597 • α cost of β -lactamase production
- 598 • β benefit of β -lactamase production
- 599 • ξ nutrient recycling
- 600 • φ antibiotic degradation by living cells
- 601 • i β -lactamase inhibitor concentration
- 602 • d_b effect of β -lactamase inhibitor on β -lactamase
- 603 • c private benefit
- 604 • κ_b antibiotic degradation by β -lactamase
- 605 • d_a basal antibiotic degradation
- 606 • h_a Hill coefficient for antibiotic
- 607 • h_i Hill coefficient for inhibitor
- 608 • ι effect of β -lactamase inhibitor on β -lactamase

609 Simulation parameters are provided in the supplemental data [Supp Table ST3].

$$g = \frac{s}{1 + s} \quad (4)$$

$$l = \gamma g \frac{a_a^h}{1 + t_a^h} \quad (5)$$

$$t = \frac{i_i^h}{1 + i_i^h} \quad (6)$$

$$\beta = \beta_{min} + c(1 - \beta_{min}) t \quad (7)$$

$$\phi = \phi_{max} (1 - ct) \quad (8)$$

$$\frac{dn_s}{dt} = (g - l) n_s \quad (9)$$

$$\frac{dn_r}{dt} = (\alpha g - \beta l) n_r \quad (10)$$

$$\frac{ds}{dt} = (\xi l - g) n_s + (\xi \beta l - \alpha g) n_r \quad (11)$$

$$\frac{da}{dt} = (-k_b b - \phi n_s - d_a) a \quad (12)$$

$$\frac{db}{dt} = \beta l n_r - d_b b \quad (13)$$

610

611 Simulation using a bounded generalized Lotka-Volterra

612 To generate consortia bounded between values of 0 and 1, we used a generalized
613 Lotka-Volterra model, as published by Wu *et al.* [33]. The model was chosen for its ability to
614 generate meaningful consortia dynamics.

615 In the following equations, the variables are as follows:

- 616 • n is the cell density
- 617 • μ is the growth rate
- 618 • σ is the environmental stress
- 619 • γ^+ is the positive interaction matrix
- 620 • γ^- is the negative interaction matrix

621 Simulation parameters are provided in the supplemental data [Supp Table ST4].

$$\frac{dn}{dt} = \mu n \left(1 - n - \frac{\sigma}{1 + \gamma^+ n + \gamma^- n} \right) \quad (14)$$

622

623 **Foundation model formulation**

624 Our custom variational autoencoder, Autoencoder7X (A7X), consists of 7 encoding layers and 8
625 decoding layers. For the encoder, the first 6 layers are convolution layers which convolutes the
626 128-time-point curve to a 128 x 8 vector. The vector is then fed into the final layer of the
627 encoder, which consists of using a linear layer to reduce the vector to the latent dimension size
628 8. During training, this final layer of the encoder consists of both a mu linear layer that calculates
629 the mean values of the latent vector, and a logvar linear layer that applies random noise to the
630 mean. Therefore, each input corresponds to a continuous latent range, and not a single value.
631 This noise is regularized to a Gaussian distribution using the Kullback-Leibler divergence loss.
632 During evaluation, only the mu layer is used to compute the latent vector. The decoder
633 architecture reverses the encoder. It first takes the 8-dimensional latent vector and passes it
634 through a linear layer. Then, it performs 7 transposed convolutions to output a 128-dimensional
635 vector. The decoder is made up of 8 layers. The first layer is a linear layer which receives the
636 latent vector as the input. The remaining layers are transposed convolution layers. The last layer
637 outputs the reconstructed growth curve. The losses were calculated using both mean squared
638 error and Kullback–Leibler divergence. For each model architecture, we used Optuna to identify
639 optimal hyperparameters [41]. Optuna is a framework developed to optimize hyperparameters in
640 machine learning tasks.

641 ***Antibiotic data preprocessing***

642 For each of the three antibiotic tasks, 1) classifying treatment, 2) predicting concentration, and
643 3) classifying resistance, we want to compare how the latents and raw curves perform on
644 different size data sets. We take our initial dataset, and split it into 80% training data and 20%
645 test data. From there, we subset the training data to create smaller training datasets to see how
646 the models perform with less training data. For each type of input, we train a model to predict
647 either the specific value (regression) or category (classification). We then compare the accuracy
648 of the models on a test dataset withheld from training.

649 ***Simulated consortia data generation and preprocessing for testing ML-based forecasting***

650 150 consortia were simulated using the chaotic consortia model [30] and the bounded
651 generalized Lotka-Volterra model [33] [Supp Fig 9A]. The consortia were generated with a total
652 number of members, and from that a subset were selected as focal members [Supp Fig 9B].
653 Each consortia had 10,000 simulations. In all of the simulations the focal community initial cell
654 densities were randomized. In some of the simulations the background strain starting cell
655 densities were randomized, but in the remaining the starting cell densities were fixed. Using the
656 maximum likelihood estimator [7], we estimated the intrinsic dimensions for each of the
657 consortia and selected 2 for our analysis: a simple consortium with 3 estimated intrinsic
658 dimensions and a complex consortium with 6 estimated intrinsic dimensions [Supp Fig 9C]. The
659 simple community had 5 focal members with 15 background strains, and the data were
660 generated using the bounded generalized Lotka-Volterra model. The complex community had 8
661 focal members with 92 background strains, and the data were generated using the chaotic
662 consortia model.

663 ***Experimental consortia data preprocessing for testing ML-based forecasting***

664 We used a dataset generated by Fujita *et al.* to predict consortia behavior over time [4]. The
665 dataset had 3 different media conditions and 2 different source microbial consortia, resulting in 6
666 different experiment conditions. For each of the 6 experiment conditions, they ran 8 replicates.
667 They took daily readings for 110 days, and each of the experiment conditions had anywhere
668 from 28 to 144 strains. However, not all the strains were present in each of the replicates, so for
669 each experiment condition we selected a focal community consisting of the strains that were
670 present across all 8 replicates. This resulted in focal communities ranging from 2 to 18 members
671 [Supp]. For each of the 6 experiment conditions, we performed eightfold cross-validation to train
672 our models. With each fold, one of the replicates was used as a test dataset and the training
673 was completed on the other 7 replicates. Due to the fact this data is daily readings, we
674 interpolate the data. We linearly interpolated the 110 time points to 1,024 time points, resulting in
675 approximately 13.75 days of data per 128 time points.

676 ***Consortia prediction and forecasting using ML***

677 The pipeline for consortia prediction and forecasting was the same for both the simulations and
678 the experimental data. We generated datasets by taking 128 time point windows from the data
679 and predicting the next 128 time points. We predicted all focal members of the consortia
680 simultaneously, so the input and output vectors were of shape M number of members \times 128 time
681 points. For each member, the 128 time point window is encoded into its latent representation,
682 and we include the maximum value from the window. This results in input and output vectors of
683 shape $M \times 9$.

684

685 We trained two separate regression models: one using the raw curves and one using the latent
686 vectors. For the raw curves, we flattened the input and output to change them from shape $M \times$
687 128 to $128M \times 1$. For the latent vectors, we flattened the input and output to change them from
688 shape $M \times 9$ to $9M \times 1$. We then used these flattened vectors to train the regression models.

689

690 To forecast into the future, we used the trained regression model. We begin with a seed of the
691 first 128 time points from a simulation. Then, we iteratively predicted the next 128-time-point
692 chunk and slid the window of the input. This allows us to forecast indefinitely into the future from
693 an initial seed.

694

695 Each consortia was individually trained on its own regression models. To characterize the
696 forecast accuracy of the model we consider the root mean square error (RMSE) over time.
697 RMSE is used for the forecasting task, because its units match the underlying consortia. It is
698 considered over time, because this allows us to analyze for how long the predictions are
699 accurate.

700 ***Consortia final abundance classification using ML***

701 For the experimental consortia, we wanted to predict whether the final abundance for each of
702 the members would be above or below 0.1. To do this, we took the first 128 time points, which
703 equated to approximately 14 days, and trained a classifier to predict whether each member of
704 the consortium would be above or below 0.1 in the final time point, 96 days later. Similar to the
705 forecasting, we predicted all members of the consortia simultaneously, so the input vectors were
706 of shape M number of members \times 128 time points and the binary output vectors were of shape
707 $M \times 1$. For each member, the 128 time point window is encoded into its latent representation,
708 and we include the maximum value from the window. This results in input and output vectors of
709 shape $M \times 9$.

710

711 We trained two separate classification models: one using the raw curves and one using the
712 latent vectors. For the raw curves, we flattened the input to change them from shape $M \times 128$ to

713 $128M \times 1$. For the latent vectors, we flattened the input to change them from shape $M \times 9$ to $9M$

714 $\times 1$. We then used these flattened vectors to train the classification models.

715

716 Each consortia was individually trained on its own classification models. We used classification

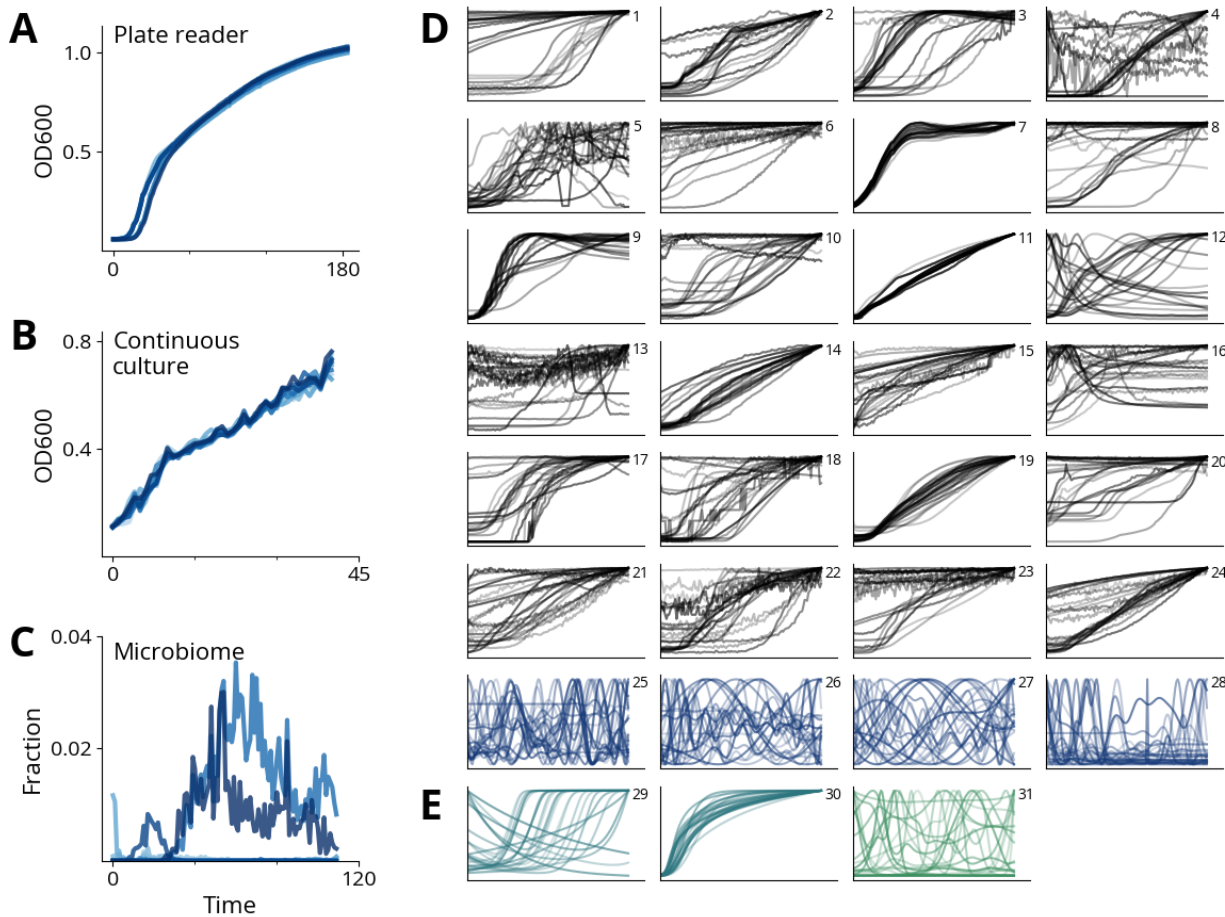
717 accuracy to report the results, which reports the fraction of the samples that were correctly

718 classified.

719

720

721 Figures



722

723 Figure 1. Overview of datasets used to train our model.

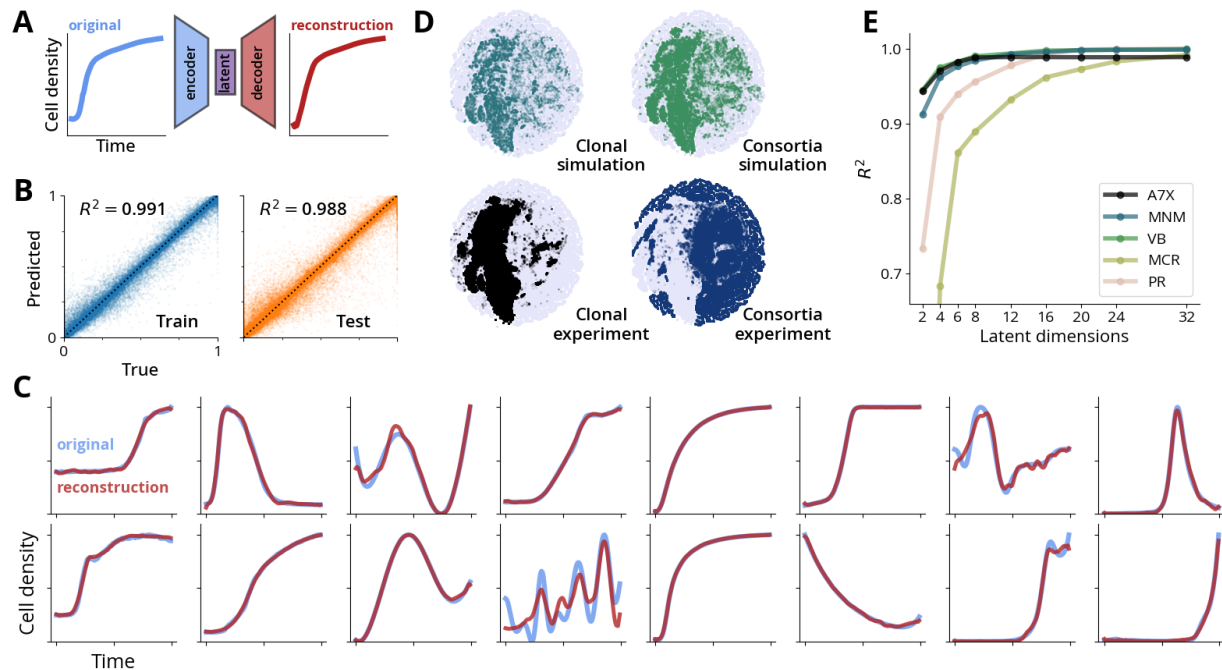
724 A) Sample growth curve using OD values from a plate reader.

725 B) Sample growth curve using OD values from a continuous culture reactor.

726 C) Sample growth curves using fractions of strains from a microbiome.

727 D) Experimental growth curves (306,417 curves) were primarily from OD plate readers,
728 automation-enabled continuous cultures, and longitudinal microbiome sequencing. Our
729 data comes from experimental clonal populations (1-24) and experimental consortia
730 (25-28).

731 E) Simulated growth curves (69,674 curves) from clonal populations (29, 30) and microbial
732 consortia (31).



733

734 **Figure 2. Encoding microbial growth dynamics in a low-dimensional latent space. Our**
735 **variational autoencoder (Autoencoder7X) compresses a 128-point growth curve into an**
736 **8-dimensional latent vector and reconstructs the original curve with high accuracy.**

737 A) Schematic of the encoder-decoder architecture used to learn latent representations.

738 B) Reconstruction accuracy for the training and test data. The training dataset is 80% of the
739 total data. The test dataset is 20% of the total data. The model has a coefficient of
740 determination of 0.991 for the training data and 0.988 for the test data. The training and
741 test datasets are 300,872 growth curves x 128 time points and 75,219 growth curves x
742 128 time points, respectively, which results in 38.5M and 9.6M data points total. For
743 visualization, we randomly selected 100,000 data points from each dataset to display in
744 the scatter plots.

745 C) Examples of original and reconstructed growth curves from our different datasets.

746 D) Visualization of the 8-dimension vectors in 2D by applying t-distributed stochastic
747 neighbor embedding. Different latent values represent different types of growth curves.

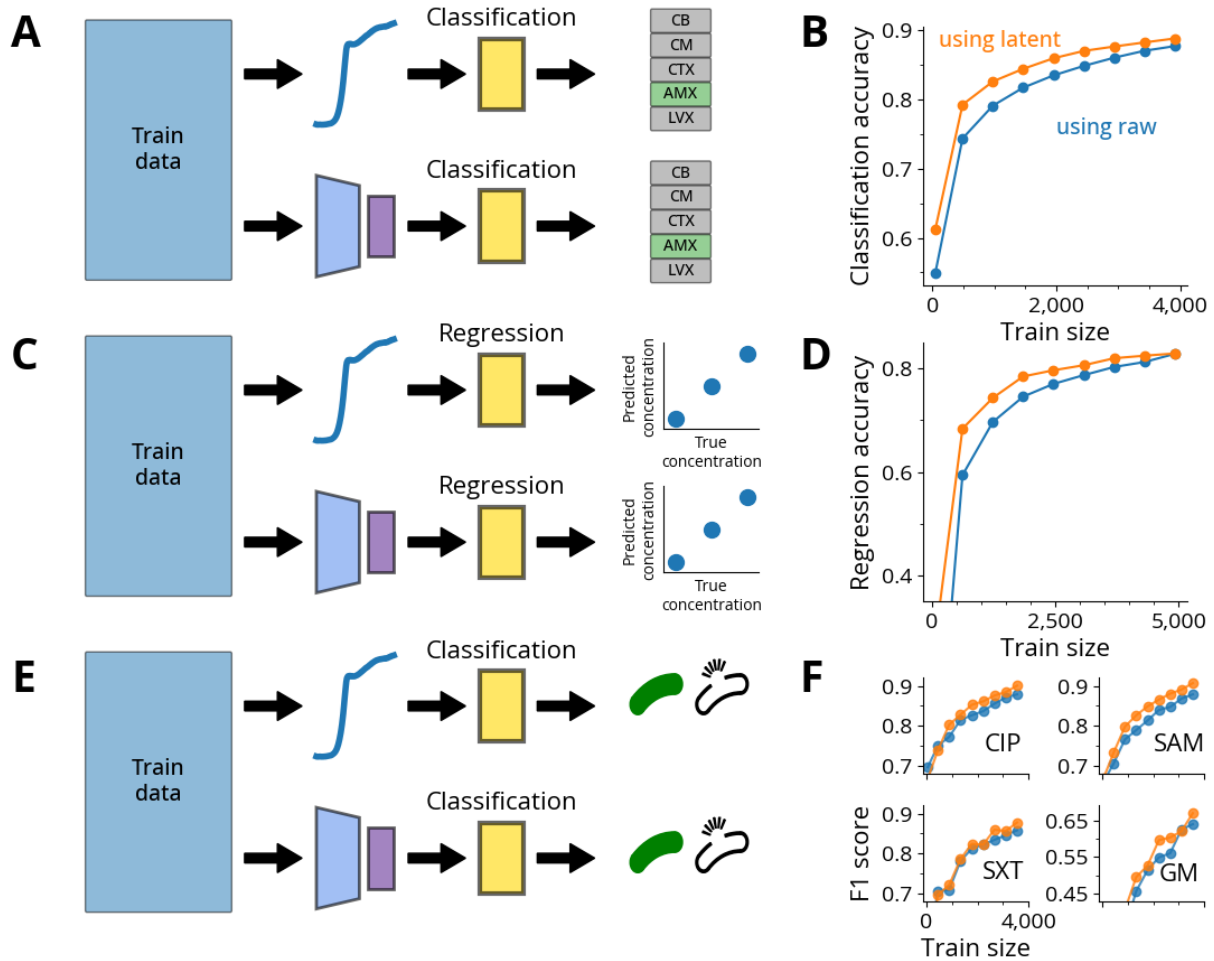
748 We visualized where the curves exist in the latent representation for each of the four

749 overarching data types: clonal simulation, consortia simulation, clonal experiment, and
750 consortia experiment.

751 E) Evaluation of different model architectures and latent dimensions. The custom variational
752 autoencoders, Autoencoder7X (A7X), Multi-layer perceptron Network Model (MNM), and
753 Variational autoencoder Bottleneck (VB), have the highest reconstruction accuracy with
754 latent dimensions below 16. Above 16, the principal components analysis reducer (PR)
755 has similar reconstruction accuracy. The transformer architecture, microBERT curve
756 reducer (MCR), does not perform as well, which is to be expected since the transformer
757 is not trained to learn reconstruction, rather it is trained to learn missing data.

758

759



760

761 **Figure 3. Latent representations enable accurate prediction of antibiotic responses even**
762 **with a small training sample size.**

763 A) Dataset of 6,924 growth curves from *E. coli* strains exposed to different antibiotics and
764 media. In this example, we predicted which type of antibiotic was used to treat the
765 bacteria using only the growth curve as the input.

766 B) Classification of antibiotic type. The latent vectors outperformed the raw curves at every
767 training size, especially under limited data.

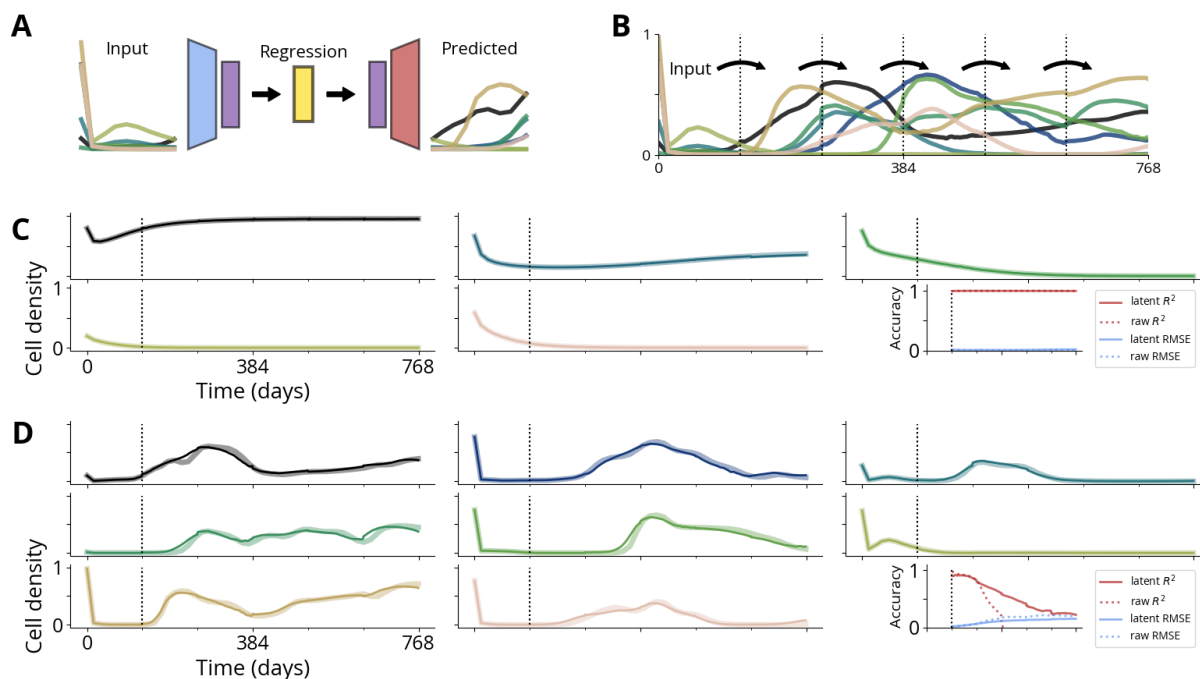
768 C) Using the same collection, we also predicted the concentration of antibiotics used to
769 treat the bacteria.

770 D) Prediction of antibiotic concentration by regression. The latent vectors outperformed the
771 raw curves at every training size, especially under limited data.

772 E) Classifying strains. Previous studies in the You lab compared the classification of
773 antibiotic resistance by comparing phenotypes, such as growth curves, against
774 genotypes. Generally, they found using growth curves provided more accurate
775 classifications of strains than the genotypes, even when the growth curves were
776 embedded to a lower dimension [1], [8]. The input data for this task is OD600 readings
777 without any added antibiotics for 277 different clinical isolates, with 4 replicates each in 4
778 different media conditions, resulting in 4,432 curves. For each antibiotic, we train a
779 classifier on either the raw curves or latent vectors, and we predict whether the cell is
780 resistant to the given antibiotic. Unless otherwise specified, we used the extra-trees
781 regressor and classifier from scikit-learn [42].

782 F1 score for the classifier on each of the antibiotics. With ciprofloxacin (CIP), the latent
783 vectors outperform the raw curves at 7 of the 9 training sizes. With ampicillin–sulbactam
784 (SAM), the latent vectors outperformed the raw curves at every training size. With
785 trimethoprim–sulfamethoxazole (SXT) and gentamicin (GM), the latent vectors
786 outperform the raw curves at 6 of the 9 training sizes. With gentamicin (GM), the latent
787 vectors outperform the raw curves at lower training dataset sizes, but both cases
788 struggle to identify all the GM samples.

789



790

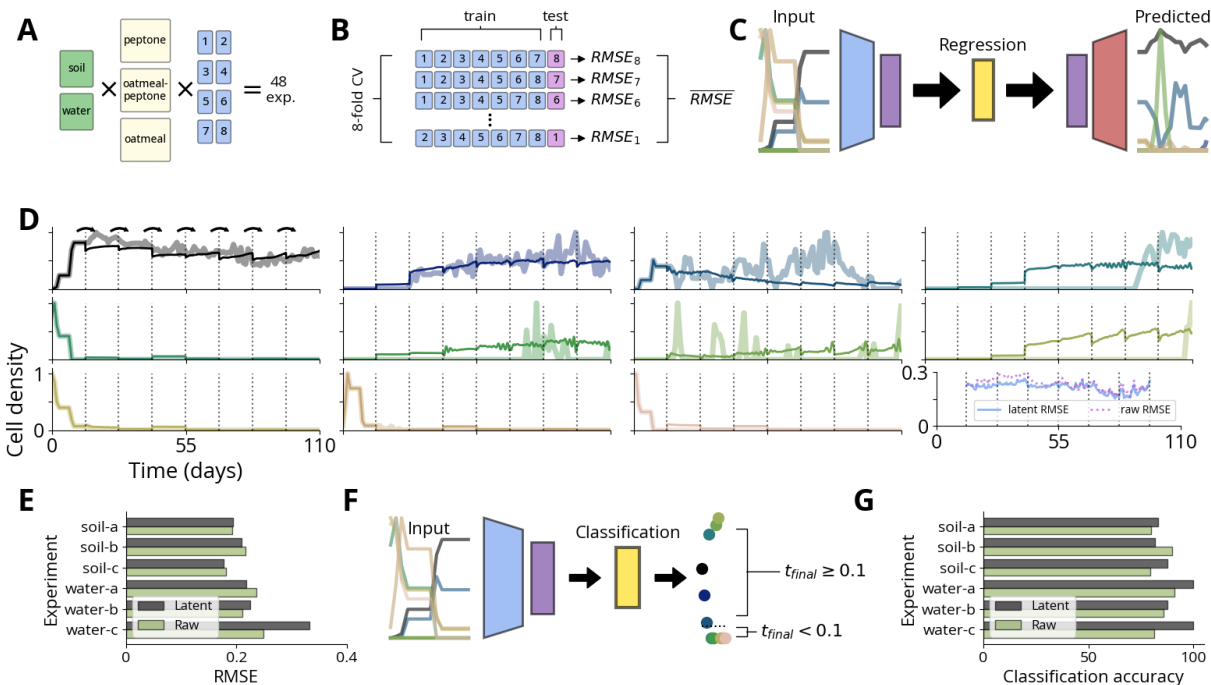
791 **Figure 4. Forecasting simulated microbial community dynamics using latent**
792 **representations.**

793 A) Workflow for training a regression model for forecasting. We defined 150 consortia and
794 simulated 10,000 curves for each consortium. The consortia each contained 20 to 100
795 members. For each consortium, we selected M ($= 2$ to 10) members for our analysis
796 (see Methods for details). From the simulated data, we use $128 \times M$ time points to predict
797 the next $128 \times M$ time points. We tested two different conditions for this task: 1) using the
798 raw data to predict the raw data and 2) using the latent representation of the input to
799 predict the latent representation of the output. For the latter, our pipeline consists of
800 training the regression model by providing the input curves, encoding to the latent
801 representation, training the regression model to predict the latent representation, which
802 is then decoded into the future 128 time points. We split the 10,000 curves into 8,000 for
803 the training dataset and 2,000 for the test dataset.

804 B) Iterative forecasting using the trained regression model. Using a single input of the first
805 128 time points from the simulation, we predict the remaining 640 time points. The
806 regression model is trained to take 128 time points and predict the next 128 time points.
807 Therefore, we use points 1 to 128 to predict 129 to 256. Then, we use the newly
808 predicted 129 to 256 to predict 257 to 384. We perform this iteratively until we have
809 predicted the entire time course.

810 C) Forecasts for a simple consortium (estimated intrinsic dimension ≈ 3) using latent
811 representations. The displayed results are from using 10% of the training data, or 800
812 simulations, to train a regression model. This community was simulated using 20 total
813 members, and we selected 5 as target members. The wider, faint lines are the ground
814 truth and the skinnier, solid lines are the forecasted behavior. Each color is a different
815 member of the community. The displayed example is a random sample from the 2,000
816 test simulations. Regression accuracy of R^2 (red) and RMSE (blue) for both using raw
817 (dotted) and latent (solid). These metrics are calculated by taking the mean values from
818 fivefold cross-validation from training regression models on the 800 curves. For this
819 simple community, both using raw and latent perform very well, and there is no
820 difference in the accuracy between the two.

821 D) Forecasts for a complex community (estimated intrinsic dimension ≈ 6). These are 8
822 members from a total community of 100 members. The displayed results are from using
823 10% of the training data, or 800 simulations, to train a regression model. The displayed
824 example is a random sample from these 2,000. Regression accuracy of R^2 (red) and
825 RMSE (blue) for both using raw (dotted) and latent (solid). For the complex community,
826 we see that using latent outperforms using raw significantly in the latter half of the
827 forecast. These metrics are calculated by taking the mean values from fivefold
828 cross-validation from training regression models on the 800 curves.



829

830 Figure 5: Forecasting and abundance prediction in experimental microbial consortia.

831 A. Experimental data, from Fujita *et al.* 48 consortia derived from two environmental
 832 sources were cultured in 3 different media types, each with 8 replicates and measured
 833 daily for 110 days [4]. Each of the source/media combinations had 28 to 144 members,
 834 however not every member was present across all 8 replicates. For each source/media
 835 combination, we selected the members (2 to 18) that were present across all 8 replicates
 836 as our focal members.

837 B. Eightfold cross-validation treats each replicate as an independent test set to quantify
 838 predictive performance. We report our results by averaging the output of the model, such
 839 as taking the mean RMSE for predicting the consortia behavior.

840 C. Regression pipeline for forecasting. The model uses $128 \times M$ time points as input to
 841 predict the next $128 \times M$ time points, where M is the number of selected members. We
 842 consider 2 different methods: 1) using raw curves to predict raw curves and 2) using

843 latent representations to predict latent representations. For each iteration of the
844 cross-validation, we are using 7 replicates for the training. The original data consists of
845 110 daily readings. We interpolate this data to a total length of 768 time points, which
846 results in about 14 days per 128 time points. During training, we take every slice from
847 these training curves: we use time points 1 to 128 to predict 129 to 256, 2 to 129 to
848 predict 130 to 257, 3 to 130 to predict 131 to 258, etc. This results in a training set of 512
849 samples per replicate, or 3,584 samples total.

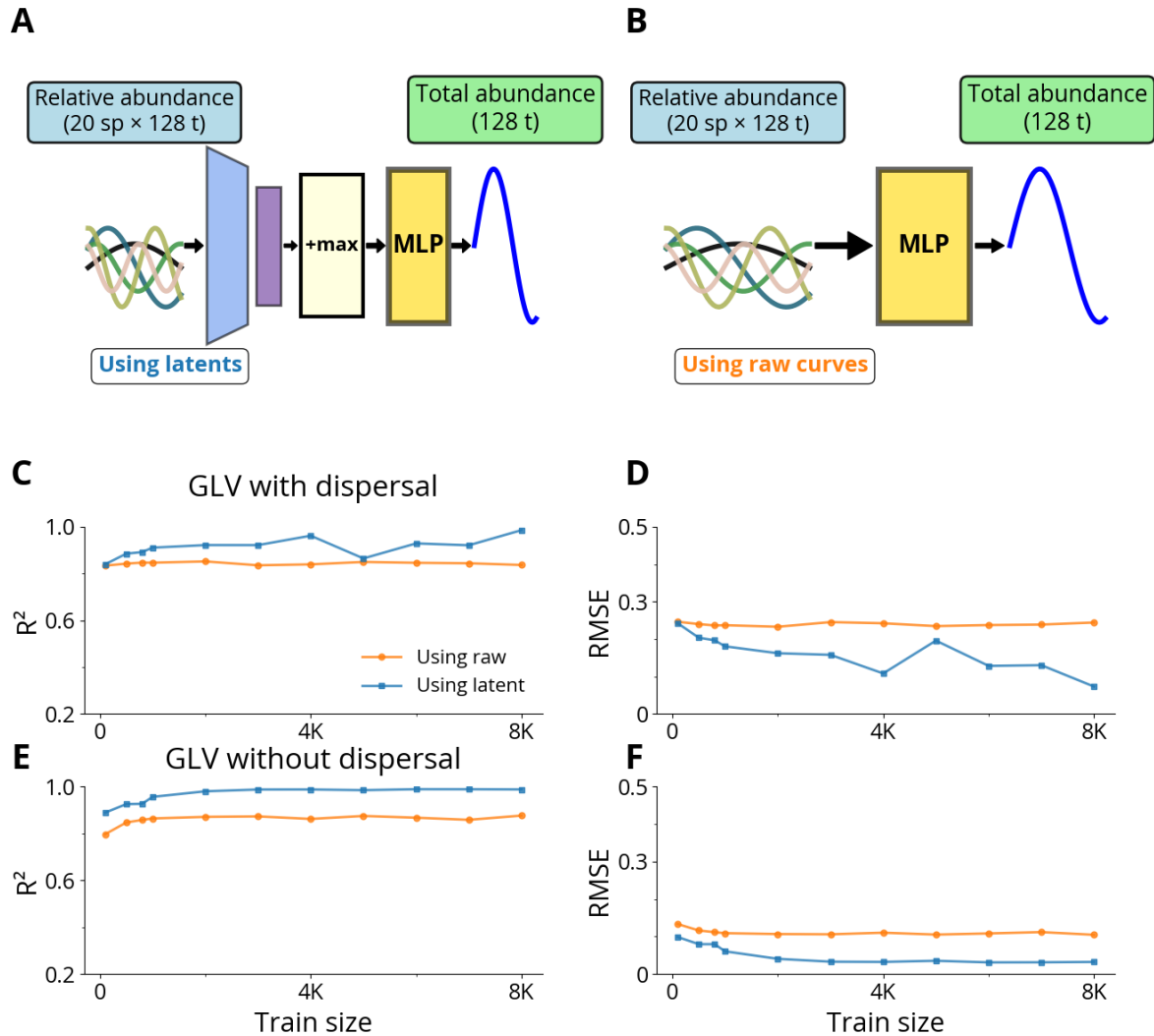
850 D. Example forecasts for one consortium. Each panel is one of the 11 members from the
851 consortium. RMSE for these data (final panel) is much higher than the simulated data
852 from the previous figure. For this consortium, latent forecasts (solid) are slightly better
853 than those based on the raw curves (dotted).

854 E. Across six source/medium combinations, latent-based forecasting yields lower mean
855 RMSE than raw-based forecasting in half of the conditions. Overall, the two approaches
856 perform similarly.

857 F. Predicting end-point dominance using early-time-window data. We take our initial 128
858 time points, which equates to 14 days, and predict whether the final abundance on day
859 110 will exceed 10%. Similar to the forecasting, we use eightfold cross-validation and
860 report the mean accuracy from the results.

861 G. Latent vector representation of the input data leads to higher classification accuracy in 5
862 of the 6 experimental conditions.

863



864

865 **Figure 6. Predicting total abundances from relative abundances in simulated**
866 **communities**

867 A. Relative abundance trajectories (20 species × 128 time points) are encoded by the
868 foundation model (Autoencoder7X) into an 8-dimensional latent space. This compressed
869 representation is concatenated with the maximum relative abundance value to create a
870 9-dimensional feature vector, which is fed into an MLP regressor to predict total
871 abundance dynamics (128 time points). Example curves illustrate diverse temporal
872 patterns in species abundances and the predicted output trajectory.

- 873 B. The same relative abundance input data is directly fed into an MLP regressor without
874 dimensionality reduction or feature extraction (2,560 total input dimensions: 20 species ×
875 128 time points). The model outputs predicted total abundance trajectories using the
876 same prediction framework as the latent pathway.
- 877 C. R^2 for both approaches evaluated across increasing training set sizes on complex
878 community dynamics generated from a generalized Lotka-Volterra framework. Latents
879 consistently outperform raw data across all training sizes.
- 880 D. RMSE values for both prediction approaches show that using the latents achieves lower
881 prediction error than using the raw data across all training set sizes, with the
882 performance gap particularly pronounced at smaller training set sizes.
- 883 E. For stable GLV dynamics, both raw relative abundance data and foundation model latent
884 representations achieved nearly identical performance, suggesting that the relatively
885 simple, predictable dynamics could be captured effectively by either representation.
- 886 F. RMSE values for the stable dynamics regime show comparable prediction errors
887 between the latent and raw pathways across all training set sizes.

888 **Acknowledgements**

889 This work was supported by National Institute of Health (L.Y.: R01AI125604, R01GM098642,
890 R01EB031869), DARPA (L.Y. HR0011-23-2-0008), and the National Science Foundation (L. Y.:
891 Cooperative Agreement No. EEC-2133504, and a graduate fellowship to Z. A. H.).

892 **Author contributions**

893 Z.A.H. and L.Y. conceived the research and designed the research framework. Z.A.H. and I.S.
894 performed all computation, model simulation, neural network training, and machine learning

895 analysis with input from A.L.H., X.W., P.C., and, L.Y.. Z.A.H., A.L.H., K.E.D., H.R.M., D.L., R.M.,
896 K.Kim, E.Ş., G.S.H., H.S., C.A.V., J.L., Y.H., A.R.S., Z.Y., S.L., D.M.S., H.D., F.W., S.W., and
897 E.C. generated experimental data. M.L., A.J.L., and L.D. provided experimental data. Z.Z.
898 generated numerical simulations for microbial consortia tasks. K.Kholina developed the
899 transformer-based model with input from Z.A.H. and P.C.. Z.A.H., I.S., and L.Y. wrote the
900 manuscript with input from H.R.M., Z.Z., D.L., R.M., E.Ş., J.L., Y.B., and S.W..

901 **Data availability**

902 All data are available in the main text, supplemental, in our GitHub repository
903 (<https://github.com/youlab/foundation>), or in our Hugging Face pages for models
904 (<https://huggingface.co/you-lab/models>) and datasets (<https://huggingface.co/you-lab/datasets>).
905 Instructions for use of the Hugging Face collection are included in the GitHub README files.

906 **Code availability**

907 All code for the project, including data collection, model training, and figure generation, is
908 available in our GitHub repository (<https://github.com/youlab/foundation>).

909 **References**

- 910 [1] C. Zhang *et al.*, “Temporal encoding of bacterial identity and traits in growth dynamics,”
911 *Proc. Natl. Acad. Sci.*, vol. 117, no. 33, pp. 20202–20210, Aug. 2020, doi:
912 10.1073/pnas.2008807117.
913 [2] Y. Ram *et al.*, “Predicting microbial growth in a mixed culture from growth curve data,”
914 *Proc. Natl. Acad. Sci.*, vol. 116, no. 29, pp. 14698–14707, July 2019, doi:
915 10.1073/pnas.1902217116.
916 [3] M. Baranwal, R. L. Clark, J. Thompson, Z. Sun, A. O. Hero, and O. S. Venturelli,
917 “Recurrent neural networks enable design of multifunctional synthetic human gut
918 microbiome dynamics,” *eLife*, vol. 11, p. e73870, June 2022, doi: 10.7554/eLife.73870.
919 [4] H. Fujita *et al.*, “Alternative stable states, nonlinear behavior, and predictability of
920 microbiome dynamics,” *Microbiome*, vol. 11, no. 1, Art. no. 1, Dec. 2023, doi:
921 10.1186/s40168-023-01474-5.
922 [5] G. Armstrong *et al.*, “Applications and Comparison of Dimensionality Reduction Methods
923 for Microbiome Data,” *Front. Bioinforma.*, vol. 2, Feb. 2022, doi:
924 10.3389/fbinf.2022.821861.

- 925 [6] M. A. Carreira-Perpinan, “A Review of Dimension Reduction Techniques”.
- 926 [7] E. Levina and P. J. Bickel, “Maximum Likelihood Estimation of Intrinsic Dimension,” in
927 *Advances in Neural Information Processing Systems*, MIT Press, 2004. [Online]. Available:
928 [https://papers.nips.cc/paper_files/paper/2004/hash/74934548253bcab8490ebd74afed7031](https://papers.nips.cc/paper_files/paper/2004/hash/74934548253bcab8490ebd74afed7031-Abstract.html)
929 [-Abstract.html](#)
- 930 [8] Y. Baig, H. R. Ma, H. Xu, and L. You, “Autoencoder neural networks enable low
931 dimensional structure analyses of microbial growth dynamics,” *Nat. Commun.*, vol. 14, no.
932 1, p. 7937, Dec. 2023, doi: 10.1038/s41467-023-43455-0.
- 933 [9] Z. Zhou *et al.*, “Linear scaling reveals low-dimensional structure in observable microbial
934 dynamics,” June 19, 2025, *bioRxiv*. doi: 10.1101/2025.06.13.659614.
- 935 [10] “Simulating 500 million years of evolution with a language model | Science.” Accessed:
936 Feb. 27, 2025. [Online]. Available:
937 <https://www-science-org.proxy.lib.duke.edu/doi/10.1126/science.ads0018>
- 938 [11] E. Nguyen *et al.*, “Sequence modeling and design from molecular to genome scale with
939 Evo,” *Science*, vol. 386, no. 6723, p. eado9336, Nov. 2024, doi: 10.1126/science.ado9336.
- 940 [12] “scBERT as a large-scale pretrained deep language model for cell type annotation of
941 single-cell RNA-seq data | Nature Machine Intelligence.” Accessed: Jan. 21, 2025.
942 [Online]. Available: <https://www.nature.com/articles/s42256-022-00534-z>
- 943 [13] A. B. George and K. S. Korolev, “Ecological landscapes guide the assembly of optimal
944 microbial communities,” *PLOS Comput. Biol.*, vol. 19, no. 1, p. e1010570, Jan. 2023, doi:
945 10.1371/journal.pcbi.1010570.
- 946 [14] S. Arya, A. B. George, and J. P. O’Dwyer, “Sparsity of higher-order landscape interactions
947 enables learning and prediction for microbiomes,” *Proc. Natl. Acad. Sci.*, vol. 120, no. 48,
948 p. e2307313120, Nov. 2023, doi: 10.1073/pnas.2307313120.
- 949 [15] “Enabling high-throughput biology with flexible open-source automation | Molecular
950 Systems Biology.” Accessed: Jan. 22, 2025. [Online]. Available:
951 <https://www.embopress.org/doi/full/10.15252/msb.20209942>
- 952 [16] H. R. Ma, H. Z. Xu, K. Kim, D. J. Anderson, and L. You, “Private benefit of β -lactamase
953 dictates selection dynamics of combination antibiotic treatment,” *Nat. Commun.*, vol. 15,
954 no. 1, p. 8337, Sept. 2024, doi: 10.1038/s41467-024-52711-w.
- 955 [17] J. N. Hennigan, R. Menacho-Melgar, P. Sarkar, M. Golovsky, and M. D. Lynch, “Scalable,
956 robust, high-throughput expression & purification of nanobodies enabled by 2-stage
957 dynamic control,” *Metab. Eng.*, vol. 85, pp. 116–130, Sept. 2024, doi:
958 10.1016/j.ymben.2024.07.012.
- 959 [18] S. Li, Z. Ye, E. A. Moreb, R. Menacho-Melgar, M. Golovsky, and M. D. Lynch, “2-Stage
960 microfermentations,” *Metab. Eng. Commun.*, vol. 18, p. e00233, June 2024, doi:
961 10.1016/j.mec.2024.e00233.
- 962 [19] S. Li *et al.*, “Dynamic control over feedback regulatory mechanisms improves NADPH flux
963 and xylitol biosynthesis in engineered *E. coli*,” *Metab. Eng.*, vol. 64, pp. 26–40, Mar. 2021,
964 doi: 10.1016/j.ymben.2021.01.005.
- 965 [20] R. Menacho-Melgar *et al.*, “Scalable, two-stage, autoinduction of recombinant protein
966 expression in *E. coli* utilizing phosphate depletion,” *Biotechnol. Bioeng.*, vol. 117, no. 9, pp.
967 2715–2727, 2020, doi: 10.1002/bit.27440.
- 968 [21] M. Ahmad *et al.*, “Tradeoff between lag time and growth rate drives the plasmid acquisition
969 cost,” *Nat. Commun.*, vol. 14, no. 1, p. 2343, Apr. 2023, doi: 10.1038/s41467-023-38022-6.
- 970 [22] S. V. Aduru *et al.*, “Sub-inhibitory antibiotic treatment selects for enhanced metabolic
971 efficiency,” *Microbiol. Spectr.*, vol. 12, no. 2, pp. e03241-23, Jan. 2024, doi:
972 10.1128/spectrum.03241-23.
- 973 [23] A. Palomino *et al.*, “Metabolic genes on conjugative plasmids are highly prevalent in
974 *Escherichia coli* and can protect against antibiotic treatment,” *ISME J.*, vol. 17, no. 1, pp.
975 151–162, Jan. 2023, doi: 10.1038/s41396-022-01329-1.

- 976 [24] H. Prenskey, A. Gomez-Simmonds, A. Uhlemann, and A. J. Lopatkin, "Conjugation
977 dynamics depend on both the plasmid acquisition cost and the fitness cost," *Mol. Syst.*
978 *Biol.*, vol. 17, no. 3, p. e9913, Mar. 2021, doi: 10.15252/msb.20209913.
- 979 [25] R. Maddamsetti *et al.*, "Duplicated antibiotic resistance genes reveal ongoing selection and
980 horizontal gene transfer in bacteria," *Nat. Commun.*, vol. 15, no. 1, p. 1449, Feb. 2024, doi:
981 10.1038/s41467-024-45638-9.
- 982 [26] Z. Ye *et al.*, "Two-stage dynamic deregulation of metabolism improves process robustness
983 & scalability in engineered *E. coli.*," *Metab. Eng.*, vol. 68, pp. 106–118, Nov. 2021, doi:
984 10.1016/j.ymben.2021.09.009.
- 985 [27] L. A. David *et al.*, "Host lifestyle affects human microbiota on daily timescales," *Genome*
986 *Biol.*, vol. 15, no. 7, p. R89, July 2014, doi: 10.1186/gb-2014-15-7-r89.
- 987 [28] J. Schluter *et al.*, "The gut microbiota is associated with immune cell dynamics in humans,"
988 *Nature*, vol. 588, no. 7837, pp. 303–307, Dec. 2020, doi: 10.1038/s41586-020-2971-8.
- 989 [29] C. V. Theodoris *et al.*, "Transfer learning enables predictions in network biology," *Nature*,
990 vol. 618, no. 7965, pp. 616–624, June 2023, doi: 10.1038/s41586-023-06139-9.
- 991 [30] J. Hu, D. R. Amor, M. Barbier, G. Bunin, and J. Gore, "Emergent phases of ecological
992 diversity and dynamics mapped in microcosms," *Science*, vol. 378, no. 6615, pp. 85–89,
993 Oct. 2022, doi: 10.1126/science.abm7841.
- 994 [31] F. Yang, F. Wang, L. Huang, L. Liu, J. Huang, and J. Yao, "Reply to: Deeper evaluation of a
995 single-cell foundation model," *Nat. Mach. Intell.*, vol. 6, no. 12, pp. 1447–1450, Dec. 2024,
996 doi: 10.1038/s42256-024-00948-x.
- 997 [32] J. E. Sulaiman *et al.*, "Elucidating human gut microbiota interactions that robustly inhibit
998 diverse *Clostridioides difficile* strains across different nutrient landscapes," *Nat. Commun.*,
999 vol. 15, no. 1, p. 7416, Aug. 2024, doi: 10.1038/s41467-024-51062-w.
- 1000 [33] F. Wu *et al.*, "Modulation of microbial community dynamics by spatial partitioning," *Nat.*
1001 *Chem. Biol.*, vol. 18, no. 4, pp. 394–402, Apr. 2022, doi: 10.1038/s41589-021-00961-w.
- 1002 [34] R. Lam *et al.*, "Learning skillful medium-range global weather forecasting," *Science*, vol.
1003 382, no. 6677, pp. 1416–1421, Dec. 2023, doi: 10.1126/science.adi2336.
- 1004 [35] "Complementing 16S rRNA Gene Amplicon Sequencing with Total Bacterial Load To Infer
1005 Absolute Species Concentrations in the Vaginal Microbiome | mSystems." Accessed: Oct.
1006 21, 2025. [Online]. Available: <https://journals.asm.org/doi/10.1128/msystems.00777-19>
- 1007 [36] G. B. Gloor, J. M. Macklaim, V. Pawlowsky-Glahn, and J. J. Egozcue, "Microbiome
1008 Datasets Are Compositional: And This Is Not Optional," *Front. Microbiol.*, vol. 8, Nov. 2017,
1009 doi: 10.3389/fmicb.2017.02224.
- 1010 [37] J. T. Barlow, S. R. Bogatyrev, and R. F. Ismagilov, "A quantitative sequencing framework
1011 for absolute abundance measurements of mucosal and luminal microbial communities,"
1012 *Nat. Commun.*, vol. 11, no. 1, p. 2590, May 2020, doi: 10.1038/s41467-020-16224-6.
- 1013 [38] "Metagenomic estimation of absolute bacterial biomass in the mammalian gut through
1014 host-derived read normalization | mSystems." Accessed: Oct. 21, 2025. [Online]. Available:
1015 <https://journals.asm.org/doi/10.1128/msystems.00984-25>
- 1016 [39] "Enabling high-throughput biology with flexible open-source automation | Molecular
1017 Systems Biology." Accessed: Jan. 09, 2025. [Online]. Available:
1018 <https://www.embopress.org/doi/full/10.15252/msb.20209942>
- 1019 [40] T. Baba *et al.*, "Construction of *Escherichia coli* K-12 in-frame, single-gene knockout
1020 mutants: the Keio collection," *Mol. Syst. Biol.*, vol. 2, no. 1, p. 2006.0008, Feb. 2006, doi:
1021 10.1038/msb4100050.
- 1022 [41] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, "Optuna: A Next-generation
1023 Hyperparameter Optimization Framework," in *Proceedings of the 25th ACM SIGKDD*
1024 *International Conference on Knowledge Discovery & Data Mining*, in KDD '19. New York,
1025 NY, USA: Association for Computing Machinery, July 2019, pp. 2623–2631. doi:
1026 10.1145/3292500.3330701.

1027 [42] F. Pedregosa *et al.*, “Scikit-learn: Machine Learning in Python,” *J. Mach. Learn. Res.*, vol.
1028 12, no. 85, pp. 2825–2830, 2011.

1029