



# Sparsity of higher-order landscape interactions enables learning and prediction for microbiomes

Shreya Arya<sup>a,1</sup> , Ashish B. George<sup>b,c,d,1</sup> , and James P. O'Dwyer<sup>b,d,2</sup>

Edited by Nils Stenseth, Universitetet i Oslo, Oslo, Norway; received May 2, 2023; accepted October 16, 2023

Microbiome engineering offers the potential to leverage microbial communities to improve outcomes in human health, agriculture, and climate. To translate this potential into reality, it is crucial to reliably predict community composition and function. But a brute force approach to cataloging community function is hindered by the combinatorial explosion in the number of ways we can combine microbial species. An alternative is to parameterize microbial community outcomes using simplified, mechanistic models, and then extrapolate these models beyond where we have sampled. But these approaches remain data-hungry, as well as requiring an a priori specification of what kinds of mechanisms are included and which are omitted. Here, we resolve both issues by introducing a mechanism-agnostic approach to predicting microbial community compositions and functions using limited data. The critical step is the identification of a sparse representation of the community landscape. We then leverage this sparsity to predict community compositions and functions, drawing from techniques in compressive sensing. We validate this approach on *in silico* community data, generated from a theoretical model. By sampling just  $\sim 1\%$  of all possible communities, we accurately predict community compositions out of sample. We then demonstrate the real-world application of our approach by applying it to four experimental datasets and showing that we can recover interpretable, accurate predictions on composition and community function from highly limited data.

microbial ecology | compressive sensing | microbiome | theoretical ecology

Our planet is host to a multitude of microbial communities, also known as microbiomes, which perform an enormous range of functions in shaping biogeochemical processes, agricultural productivity, and animal and human health (1–3). In recent years, there have been concerted efforts to modify such communities in order to alter plant, animal, human, and environmental health for the better (4–14). The complex interspecific interactions present in real communities lead to diverse steady-state communities, but these interactions also mean that the final composition and function of a microbiome may be very different from its initial composition. Moreover, for large numbers of potential taxa to include in a community, there is a combinatorially large number of distinct initial compositions. Putting these two issues together makes a brute-force approach to cataloging potential microbiomes impossible: with as few as ten species from which to build an initial community,  $2^{10}$ , or around 1,000 experiments would be needed to survey the full range of possible outcomes arising as the dynamics of communities play out. With 100 potential initial species, that number becomes  $10^{30}$  (15).

Building mechanistic models of microbial interspecific interactions has the potential to alleviate this issue. If we can write down and parameterize a model that accurately represents interactions, but using only a limited amount of experimental data, we might be able to use such a model to predict outcomes beyond those experiments. A longstanding approach to modeling interspecific interactions, inspired by Robert May's seminal work (16) on complex systems, has been to use pairwise interactions among species, with the strength of this pairwise dependence encoded in a community interaction matrix (17, 18). There is certainly no conceptual obstacle to pairwise interactions providing a potential description of microbial communities—theoretical work demonstrates that pairwise interactions can lead to diverse, realistic communities (19–26). But there is also a recent realization that microbial communities may be infused with higher-order interspecific interactions, where the influence of one species on another is dependent on the presence or absence of a third, fourth, or fifth species (27–31). Even for pairwise interspecific interactions, having  $S^2$  parameters to fit means that mechanistic models could just bring us back to a similar experimental load of the combinatorial problem we started with (18, 27, 32, 33). Higher-order interspecific interactions would then only

## Significance

Engineered microbiomes can hugely benefit human, plant, and animal health. However, the diversity and complexity of microbiomes hinder a full understanding, and hence, prediction, of community assembly outcomes, and experimental efforts are limited by the exponential number of combinations required to be designed and tested. We consider ecological landscapes of microbial abundances, which are maps from input species to output steady-state species abundance. Using tools from signal processing, a field that focuses on information acquisition and reconstruction, we find that species abundances in both real and simulated microbiomes have a sparse representation, which translates to relative rarity of higher-order landscape interactions. We then use this sparsity to learn and predict entire landscapes from highly limited experimental data using compressive sensing.

Author contributions: S.A., A.B.G., and J.P.O. designed research; S.A., A.B.G., and J.P.O. performed research; S.A. and A.B.G. contributed new reagents/analytic tools; S.A. and A.B.G. analyzed data; and S.A., A.B.G., and J.P.O. wrote the paper.

The authors declare no competing interest.

This article is a PNAS Direct Submission.

Copyright © 2023 the Author(s). Published by PNAS. This article is distributed under [Creative Commons Attribution-NonCommercial-NoDerivatives License 4.0 \(CC BY-NC-ND\)](https://creativecommons.org/licenses/by-nc-nd/4.0/).

<sup>1</sup>S.A. and A.B.G. contributed equally to this work.

<sup>2</sup>To whom correspondence may be addressed. Email: [jodwyer@illinois.edu](mailto:jodwyer@illinois.edu).

This article contains supporting information online at <https://www.pnas.org/lookup/suppl/doi:10.1073/pnas.2307313120/-/DCSupplemental>.

Published November 22, 2023.

make this problem of fitting parameters even more challenging. Finally, deep learning has been deployed to address this question (32, 34), but at the expense of making results and predictions challenging to interpret.

Developing a method that requires only limited experimental data, is agnostic to any particular ecological model, and provides ecologically interpretable, accurate results would add an important tool to the microbial ecologist's toolbox. Here, we will address this need, introducing an approach that reveals a previously hidden sparsity in the landscape describing the map from initial composition to final community states, for both simulated and experimental microbial communities. The sparsity of this landscape means that the outcomes of microbial community assembly are highly constrained, in effect containing much less information than the full combinatorial set would suggest. We then leverage this sparsity, using algorithms from compressive sensing to accurately recover the sparse representation. Thus, we predict the late-time outcomes for microbial community composition, from highly limited input data, in a way that is readily interpretable in terms of the sparsity of landscape interactions.

## Framework

**The Challenge of Predicting Microbial Community Composition.** We first explicitly state the problem to address: Given a set of  $S$  species to draw from, there are  $2^S - 1$  possible combinations, or seed communities, that can be formed, based on the initial presence-absence of species. This initial condition naturally does not capture the full subsequent behavior of the community, which could be extremely complex (35–38). Here, we will make a simplifying assumption that, from a given initial condition, all species will approach an equilibrium at late times, such that each species ends up with a relative abundance that only depends on the seed community composition. This excludes more general dynamics, including chaos (39), limit cycles (40), or priority effects (41). We justify this assumption by noting, as in ref. 15, that this simple behavior is what has been frequently observed in experimental communities. Furthermore, this behavior is predicted by many mechanistic models (see table 1 in ref. 15). More complex behavior is possible, both in theoretical models, where the Lotka-Volterra equations can exhibit multiple attractors (40), and in experimental communities, where both dynamical attractors or a single stable equilibrium can occur (42). But making progress in the case of systems with a unique steady state may be an important step toward tackling these more general cases. Even with this assumption of a unique steady state, if there are any interspecific interactions at all, the abundance of a species at steady state will in general depend in a complex way on which other species are present.

We will label each possible subcommunity one can obtain from a pool of  $S$  species using a binary vector  $\vec{\sigma}$ , with ones (zeros) denoting the presence (absence) of each species in the initial community. There are  $2^S$  distinct values of  $\vec{\sigma}$  corresponding to the different species combinations possible, with one ecologically trivial case of all species being absent initially. The steady-state abundance of a species  $i$  in one of the subcommunities,  $\vec{\sigma}$ , can then be written as  $N_i(\vec{\sigma})$ , since the abundance by assumption depends on the presence of species in the seed community. Further, since species  $i$  is absent from half of the  $2^S$  possible combinations,  $N_i(\vec{\sigma})$  can take up to  $2^{S-1}$  nonzero, and potentially distinct, values. Thus, representing species abundances in the  $\vec{\sigma}$  basis requires a binary-ordered vector

$\vec{a}_i(\vec{\sigma})$ , of up to  $2^{S-1}$  nonzero coefficients,  $N_i(\vec{\sigma})$ . For each species, we have thus defined a map between the  $\vec{\sigma}$  space and the space of steady-state abundances. This formalism, shown in Fig. 1, defining maps based on composition, has recently been explored in the context of composition and community functions, where the maps are referred to as ecological landscapes (43–45). Here, we define a distinct ecological landscape for each species in the community, where the mapping of interest is between the initial composition of the community, and the final, steady-state abundance of each focal species. Many functions and services of microbial communities depend on species abundances, and thus, reliably predicting late-time abundances is a key step on the way to predicting community function.

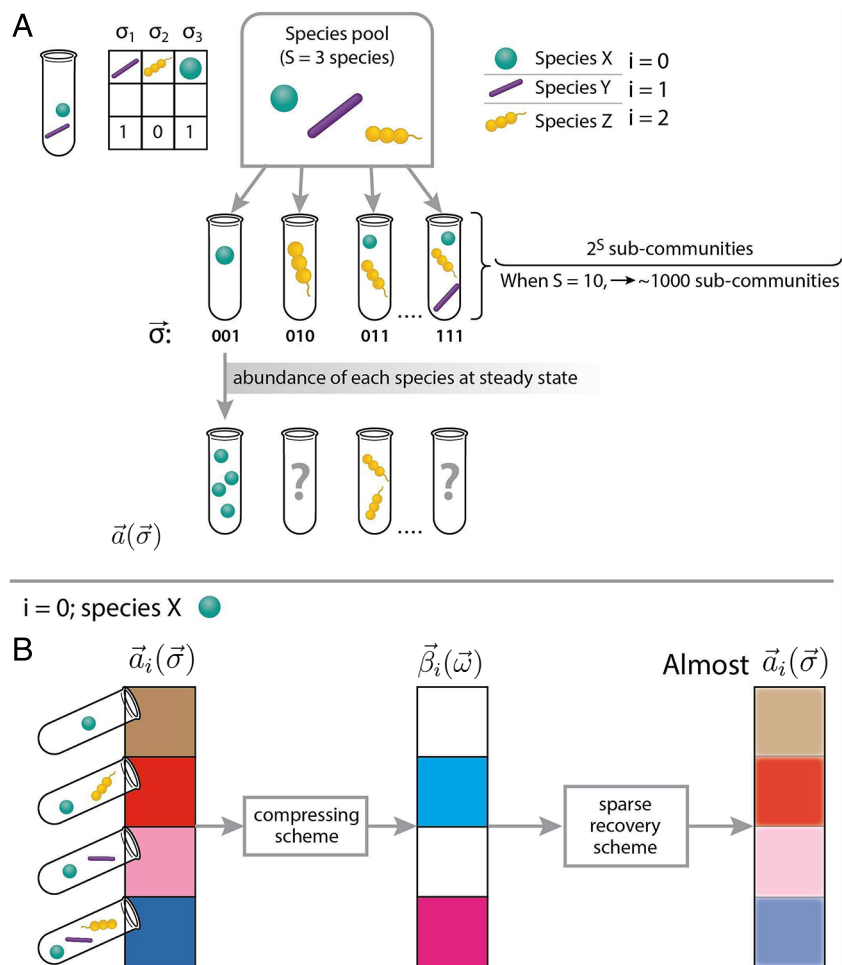
A promising approach to overcoming this combinatorial challenge would be if there were a more efficient representation of steady-state species abundances using fewer coefficients. In the field of signal processing, we would say that this is possible for a given signal if there is a representation that compresses it. We are not guaranteed that such a representation exists for microbiomes, but below we will show that, viewed in the right basis, microbiomes are highly compressible.

## Results

**Species Abundances Can be Represented Using Only a Few Coefficients in a Transformed Basis.** The question of mapping initial community composition to species abundances at late times is related to the problem of mapping genotypes to fitness in evolutionary genetics, where the effect of interactions between species, which we term “landscape interactions,” is analogous to the effect of epistasis among distinct genetic mutations. This connection can be made more concrete by considering a combinatorially complete fitness landscape of a genome with  $S$  positions. The different states of this genome can be enumerated by the presence or absence of a mutation at each of the positions. For each of these  $2^S$  genotypes, there exists a scalar value of interest, the fitness. The landscape can be represented as a vector  $\vec{f}$ , with length,  $2^S$ . Each element is a fitness value and corresponds to a genotype ordered by an  $S$ -bit binary number whose digits 1 and 0, respectively, signal the presence or absence of the mutation at the corresponding position (46, 47). In the case of microbiomes and under the assumptions we outlined in the previous section, the various genotypes are analogous to the distinct subcommunities that are possible when combining species based on initial presence-absence, and the fitness vectors are analogous to the steady-state abundances, with one vector for each species.

Here, we draw upon recent progress in the field of evolutionary genetics to inspire a related approach to microbiome prediction. Specifically, some empirical fitness landscapes are sparse when represented in a type of high-dimensional Fourier basis, called the Walsh-Hadamard basis (47–52). Moreover, in the field of signal processing, it has been noted that most natural images and many time-series signals are stored efficiently in the same kind of basis (53). Identification of these sparse representations, or sparse coding, has then been used to learn entire datasets from sampling only limited data using algorithms from the field of compressive sensing (48, 52, 54, 55).

These breakthroughs in other fields motivated us to consider the representation of species abundances in various subcommunities in a similarly transformed perspective. Specifically, we used a weighted Walsh-Hadamard basis (*Methods* and *SI Appendix*) (52, 56, 57). Up to multiplication by a diagonal matrix, this



**Fig. 1.** A sparse recovery algorithm can be used to predict microbial community end points. (A) Given a species pool with  $S$  members, there are  $2^S - 1$  subcommunities, based on different choices of initial presence-absence of the members. This problem of exponential scale means that only a few subcommunities can be sampled in the laboratory. In this paper, we focus on predicting end point community compositions of unseen subcommunities, given limited data. (B) We found that the relative abundances of any species in the pool when stacked according to the subcommunities in which it was initially present is sparse in a weighted Walsh–Hadamard basis ( $\vec{\beta}(\vec{\omega})$ ). This means that even though a chosen focal species (green, round) may have numerous, different steady-state abundances in different subcommunities, a weighted Walsh–Hadamard representation of this vector will have only a few components. The numerous, different possible steady-state abundances are represented by the colors in the column corresponding to  $\vec{a}_i(\vec{\sigma})$ . Sparsity of the representation,  $\vec{\beta}_i(\vec{\omega})$ , is visualised by the number of transparent components, or boxes, which indicate that a lot of the coefficients in this representation are insignificant. This sparsity may be leveraged by using a sparse recovery technique, called compressive sensing, which prescribes that a much smaller, generic sampling of subcommunities is required to efficiently predict steady-state abundances of the species in all other, unsampled subcommunities.

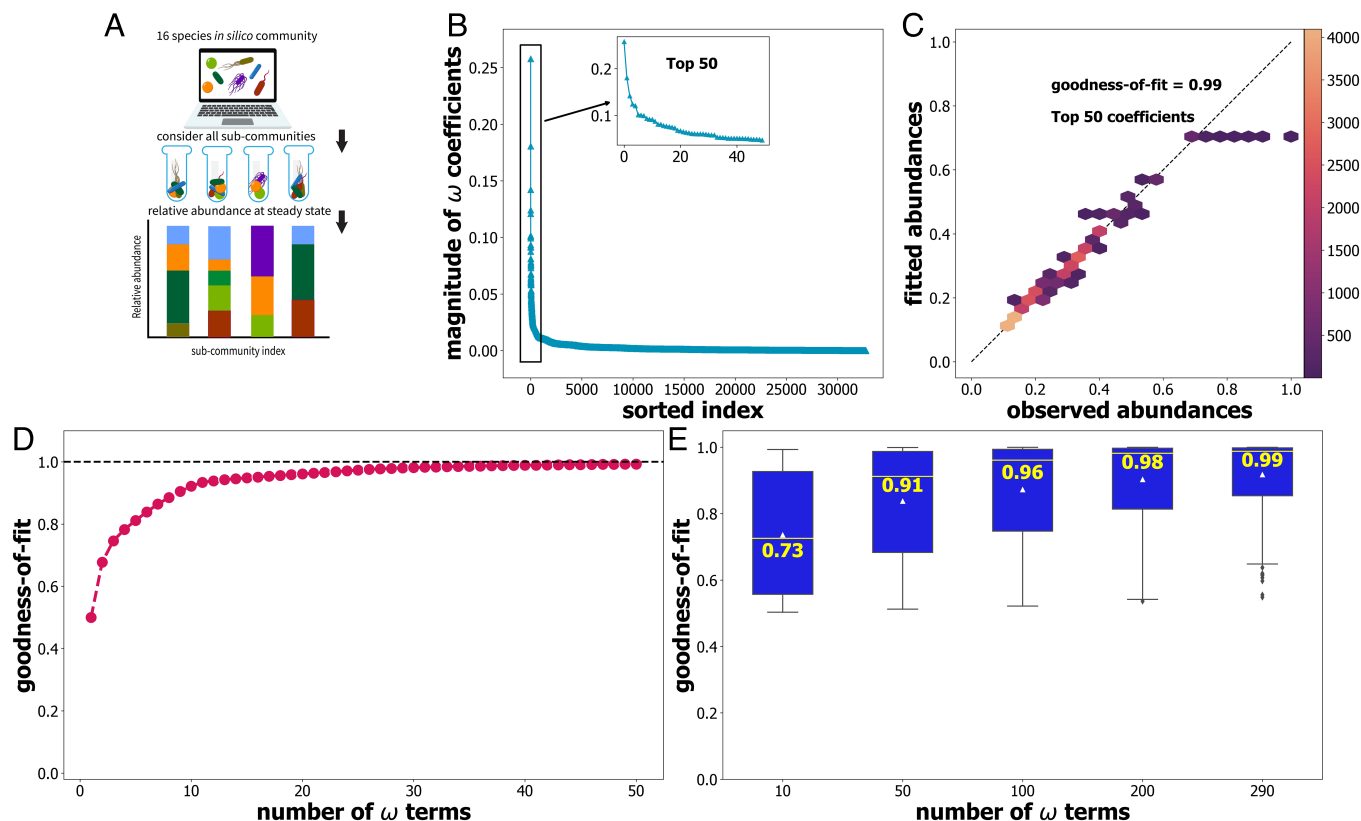
transform is the well-known Walsh–Hadamard (WH) transform, an analog of Fourier transforms in this discrete basis. It has several useful properties as discussed in *SI Appendix*, including orthogonality and symmetry (46, 47, 58). In this work, we denote this weighted Walsh–Hadamard basis as  $\vec{\omega}$ . This is different from the presence–absence basis, which we denote using  $\vec{\sigma}$ . In terms of the abundance of a species in the presence–absence basis, which is given by  $\vec{a}(\vec{\sigma})$ , the transformed abundances are given by  $\vec{\beta}(\vec{\omega})$  with the following relation between the two representations:

$$\vec{\beta}(\vec{\omega}) = \Omega \vec{a}(\vec{\sigma}), \quad [1]$$

where  $\Omega$  is the transformation matrix defined as  $\Omega = VH$ ,  $H$  is the Hadamard matrix, and  $V$  is a diagonal weighting matrix, both described in *Methods*. Multiplication by the Hadamard matrix  $H$  transforms the representation into an orthogonal basis, while multiplication by the weighting matrix  $V$  implements an ecologically motivated assumption that down-weights higher-order landscape interactions. This is based on the assumption that while higher-order landscape interactions

are allowed, lower-order interactions are relatively more likely, and the weighting encoded in  $V$  captures this expectation (see *Methods* and *SI Appendix*, Figs. S3 and S7 for more details).

To test whether microbiomes could plausibly be sparse in the  $\vec{\omega}$  basis, we first generated abundance data from simulated microbial communities using microbial consumer–resource (MiCRM) models (59, 60). We consider a community of 16 species, with interspecific interactions mediated by consumption of resources. As elaborated in *Methods*, the MiCRM is a realistic choice for studying in silico microbiomes since its dynamical equations include terms for both cross-feeding and competition. In these in silico communities, species are distinguished by the differences in their ability to consume and secrete various resources. The advantage of starting with simulated communities is that we are able to comprehensively generate all possible species combinations and thus test rigorously to what extent sparsity does or does not hold in the  $\vec{\omega}$  basis. In particular, from a pool of 16 species, we considered the stationary points of the MiCRM corresponding to each of the total possible  $2^{16} = 65,536$  species



**Fig. 2.** Relative abundances of a simulated community are sparse in the weighted Walsh-Hadamard ( $\vec{w}$ ) basis. (A) We simulated a pool of 16-species in silico using microbial consumer–resource models (59, 60). Simulations started from all 65,535 possible subcommunities corresponding to distinct species presence-absences ( $\vec{\sigma}$ ) to obtain species abundances in the different steady-state communities. The relative abundances of a species, in the 32,678 subcommunities in which it was initially present, constitute the coefficients of  $\vec{a}(\vec{\sigma})$ . By using a weighted Walsh-Hadamard transform, species abundance can be represented in an orthogonal  $\vec{w}$  basis. (B) The coefficients used to represent abundances of a typical species in the 32,678 subcommunities using the  $\vec{w}$  basis are shown sorted by absolute size. The vast majority of the 32,678 coefficients are small; the inset plots the largest 50 coefficients. (C) Using only 50 coefficients, we were able to explain most (99%) of the variation in species abundances. The color bar indicates the density of points in the hexplot. Panel (D) demonstrates that a small number of  $\vec{w}$  coefficients is sufficient to fit abundances almost perfectly. (E) Box-plots of goodness-of-fit, aggregated for all 16 species in all 10 replicates, demonstrate that only a small number of  $\vec{w}$  coefficients are required to explain most of the observed variation in abundances. We display medians in the text and by yellow lines, while mean statistics is denoted by white triangles. Coefficients in panel (B) are ordered by magnitude. The goodness-of-fit is calculated via Eq. 2.

assemblages, obtaining the late-time, steady-state abundances arising from each initial species combination. To ensure the robustness of our conclusions, we considered 10 different sets of 16-species communities.

Fig. 2B shows the coefficients in the  $\vec{w}$  basis,  $\beta(\vec{w})$ , calculated from the steady-state abundances of a typical species from this simulated pool. A small fraction of the 32,768 ( $=2^{16-1}$ ) coefficients are significantly larger than the others, and using just the 50 largest coefficients to estimate the observed abundances in each community, we explain 99% of the observed variance (Fig. 2C). Fig. 2D shows the goodness of fit as we change the number of  $\beta(\vec{w})$  coefficients used to fit species abundance,  $\vec{a}(\vec{\sigma})$ . The goodness of fit is measured as a prediction score (52, 61),

$$\text{PS} = \frac{1}{2 - R^2}, \quad [2]$$

where  $R^2$  is the coefficient of determination. The prediction score always lies between 0 and 1. For a model that could predict only the mean abundance (over all subcommunities in which a species was initially present) regardless of the input subcommunity, this prediction score would be 0.5 (corresponding to  $R^2 = 0$ ). A score of 1 indicates perfect prediction. The high value of this goodness-of-fit measure indicates that we are able to capture

most of the variation in species abundance using a very small number of coefficients, reinforcing that this representation of species abundances is highly sparse. This finding is consistent for all species across the 10 species pools we simulated (Fig. 2E).

**Compressive Sensing Predicts Simulated Community Compositions from Limited Data.** In simulated consumer–resource models, we have established that there is a sparse representation of species abundances. But we obtained this sparse representation from a complete knowledge of all species abundances at steady-state,  $\vec{a}(\vec{\sigma})$ , from all initial species combinations. Even if sparsity does also hold in real microbiomes, we will need a method to identify the largest, most significant coefficients of the vector  $\vec{\beta}(\vec{w})$  from only limited experimental data. Compressive sensing (CS) provides this method. Assuming that a sparse representation *does* exist for a given signal, compressive sensing will recover the signal in its entirety from only limited, generic sampling (48, 52, 54).

In the context of microbial ecology, results shown in Fig. 2 conclusively establish sparsity for our in silico communities: The steady-state abundance vector,  $\vec{a}_i(\vec{\sigma})$ , of most species is afforded a sparse representation in the  $\vec{w}$  basis. Compressive sensing may then allow us to recover these  $\vec{w}$  coefficients from limited sampling of abundance data,  $\vec{a}_{i,\text{sampled}}(\vec{\sigma})$ . Specifically, using

compressive sensing translates to implementing a constrained optimization algorithm that finds a solution,  $\vec{\beta}_{i,\text{bestguess}}(\vec{\omega})$ , which has the lowest-possible  $l_1$ -norm while possessing an inverse Walsh–Hadamard transform that optimally fits the sampled abundances. In this work, we implemented the Basis Pursuit Denoising (BPDN) algorithm (52, 55, 62) to find the  $\vec{\beta}_{i,\text{bestguess}}(\vec{\omega})$ . The inverse Walsh–Hadamard transform then yields the model prediction for abundances:

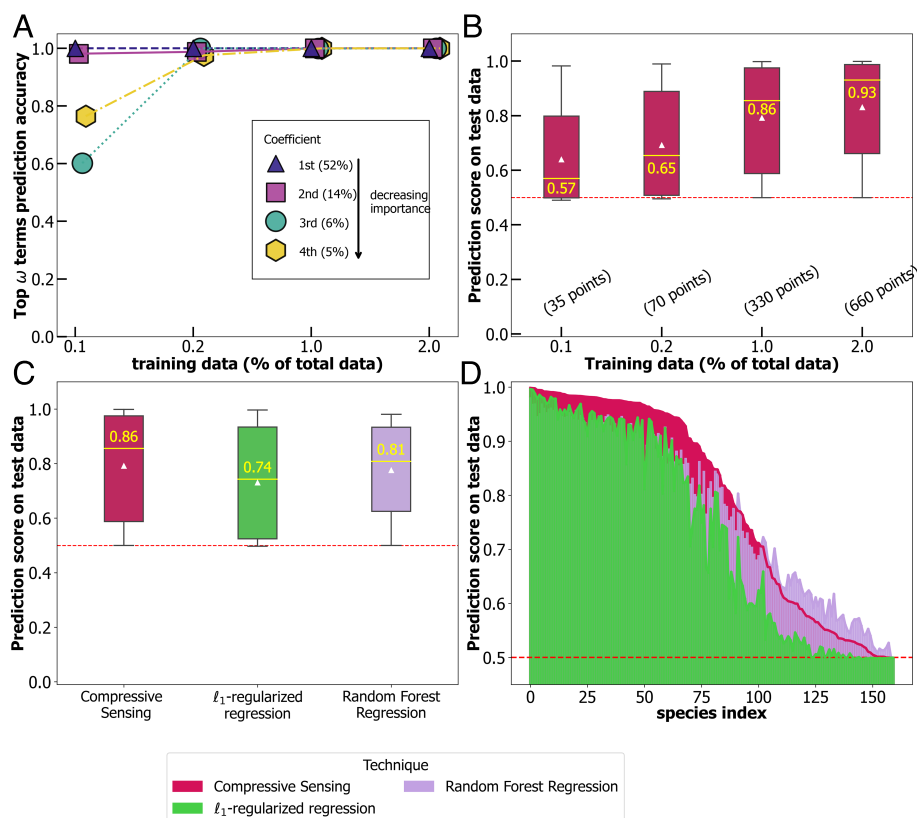
$$\vec{a}_{\text{bestguess}}(\vec{\sigma}) = \Omega^{-1} \vec{\beta}_{\text{bestguess}}(\vec{\omega}). \quad [3]$$

We tested the predictive power of compressive sensing for the *in silico* communities, for which we have access to all data in practice, and where we have already shown that a sparse basis exists. Specifically, we restricted ourselves to observing/sampling 1% or 327 subcommunities of the 32,678 possible subcommunities with a given species present initially. We also considered lower sampling percentages of 0.1%, 0.2%, and a higher percentage of 2%. With these limited “training” data, we used CS algorithms to estimate coefficients in the  $\vec{\omega}$  basis. Fig. 3A shows that CS is able

to infer the most significant coefficients of the ground truth  $\vec{\beta}(\vec{\omega})$  accurately from limited training data, and the accuracy increases as more data are obtained for training.

Next, we compare the recovered abundances with the ground-truth abundances. We compute the recovered abundance for a species using Eq. 3, where  $\vec{\beta}_{\text{bestguess}}(\vec{\omega})$  is the Walsh–Hadamard vector inferred by the BPDN algorithm when trained on the limited abundance data for the species. Fig. 3B shows the resulting prediction score on out-of-sample abundance data.

We see that the algorithm’s predictive power improves with more training data and appears to asymptote at as little as 1% of all possible data included in the training set. We compare our method with two alternative approaches: an  $l_1$ -regularized regression on the  $\vec{\sigma}$  basis that tries to estimate  $\vec{a}_i(\vec{\sigma})$  directly and random forest regression (see *Methods* for detailed description of these algorithms). The  $l_1$ -regularized regression is adopted since the number of observations is smaller than the number of possible coefficients, i.e., the problem is underdetermined. We found that compressive sensing outperforms both regularized regression and random forest regression, as can be seen from



**Fig. 3.** Sparsity in the  $\vec{\omega}$  coefficients implies abundances can be practically recovered from limited data in simulated communities. (A) The accuracy of predicting the 4 most significant  $\omega$  coefficients by compressive sensing (CS) increases as the size of training data is increased. The largest coefficients, which explain most of the variance in abundance, are more easily predicted than smaller coefficients from limited data. The percentages in the parenthesis next to the coefficient symbol and index denote the percentage of the variance explained by each coefficient taken individually. This means, for example, that the first  $\omega$  term explains 52% of the total variance. (B) The ability of compressive sensing to predict species abundances in new subcommunities when trained on a small fraction of all possible subcommunities is shown using the Prediction Score (Eq. 2). Note that 1% of the total number of possible subcommunities corresponds to 327 subcommunities, which would require four 96-well plates. The median score, shown by the yellow line is annotated. The number of subcommunities sampled at each training level is indicated within parentheses. (C) The prediction score of two alternative methods,  $l_1$ -regularized regression and random forest regression, is shown for comparison. We controlled for data so that all the methods were trained on the same 327 (1%) communities. Box plots denote the performance across the 160 species (10 replicates of pools with 16 species). A model that outputs only the mean abundances of a species in all subcommunities in which the species was initially present, regardless of the input subcommunities, would get a prediction score of 0.5. We see that compressive sensing does objectively better than such a model (prediction score shown using red dashed line) while outperforming random forest regression and  $l_1$ -regularized regression. (D) Comparing the performance of CS,  $l_1$ -regularised regression, and random forest on a species-by-species basis, we find that some species are easier to predict than others. Nevertheless, CS outperforms both regressions for most species. The permutation test  $P$ -value for the difference in the performance at 1% sampling between compressive sensing and random forest regression is statistically significant at  $10^{-4}$  and that between compressive sensing and  $l_1$ -regularized regression is  $5 \times 10^{-5}$ . All results are averages over 5 runs.

the median values of the predictions in Fig. 3C. Finally, Fig. 3D compares the prediction score of the three methods for each species, which is a matched comparison at the species level of the three approaches, with training data held at 1%. This analysis reveals that some species are easier to predict, with a prediction score close to 1, while some others that are harder to predict. In *SI Appendix*, we show that species which are easier to predict are dominated by low-order landscape interactions and have a smoother abundance landscape (*SI Appendix*, Fig. S4 and ref. 58). We can thus conclude that prevalence of low-order interactions leads to a more predictable landscape. Despite this variation in performance at the species level, compressive sensing outperforms  $l_1$ -regularized regression in 97% of the cases and outperforms random forest regression in a majority of the cases (67%).

Thus, compressive sensing provides a way to use the sparse basis for our simulated consumer–resource communities, along with limited observations of species abundances in different subcommunities, to predict species abundances out-of-sample.

**Compressive Sensing Predicts Species Abundances in Experimental Data.** Using simulated communities, we demonstrate that the challenge posed by an exponential number of subcommunities to the problem of ecological predictions is massively reduced with compressive sensing. We now test the performance of CS as a predictive tool in a diverse range of real microbiomes, in the context of unknown community dynamics which likely depart from our idealized consumer–resource model, in addition to the difficulties posed by both stochasticity and incomplete data. We draw data from four published studies of microbial communities spanning in vitro and in vivo conditions, pool sizes from 5 to 16, and microbes from environments including the human gut, soil, and fruit flies (28, 29, 63, 64).

In the first study, Gould et al. assembled all combinations of 5 species in vivo in the gut of *Drosophila melanogaster*. The remaining three studies considered larger species pool sizes and hence assembled only a fraction of all possible species combinations in vitro. Sanchez-Gorostiaga et al. assembled 53 combinations of 6 starch-degrading soil microbes. Friedman et al. assembled 101 combinations of 8 soil microbes. From Clark et al., we consider 187 combinations assembled of 16 gut microbes. Fig. 4 demonstrates the performance of compressive sensing,  $l_1$ -regularized regression, and random forests on predicting unseen communities using a k-fold cross-validation procedure (*Methods*). We used a k-fold cross-validation procedure due to the small number of data points and report the prediction score, at the optimal value of the hyperparameter, on the stacked validation sets. Further, keeping the bias-variance trade-off (65) in mind, we set  $k = 3$  for most experimental communities. A complete description is given in *Methods*. Fig. 4B–E show the performance of compressive sensing, alongside the alternative methods, on the individual species. As in simulations, we find that some species are easier to predict than others. Compressive sensing outperforms regularized regression in the majority of the cases (32 out of 35) and also does better than random forest regression in 29 of 35 cases. At the community level, as shown in Fig. 4, compressive sensing outperforms both methods when comparing the mean prediction score on a dataset. This difference in prediction is statistically significant for the two largest datasets. This is tabulated in Table 1, where we list the  $P$ -values of a one-sided permutation test. We also note that species which are harder to predict using compressive sensing remain hard to predict when using the other methods. In summary, our compressive sensing approach enables more accurate predictions of microbial community abundances.

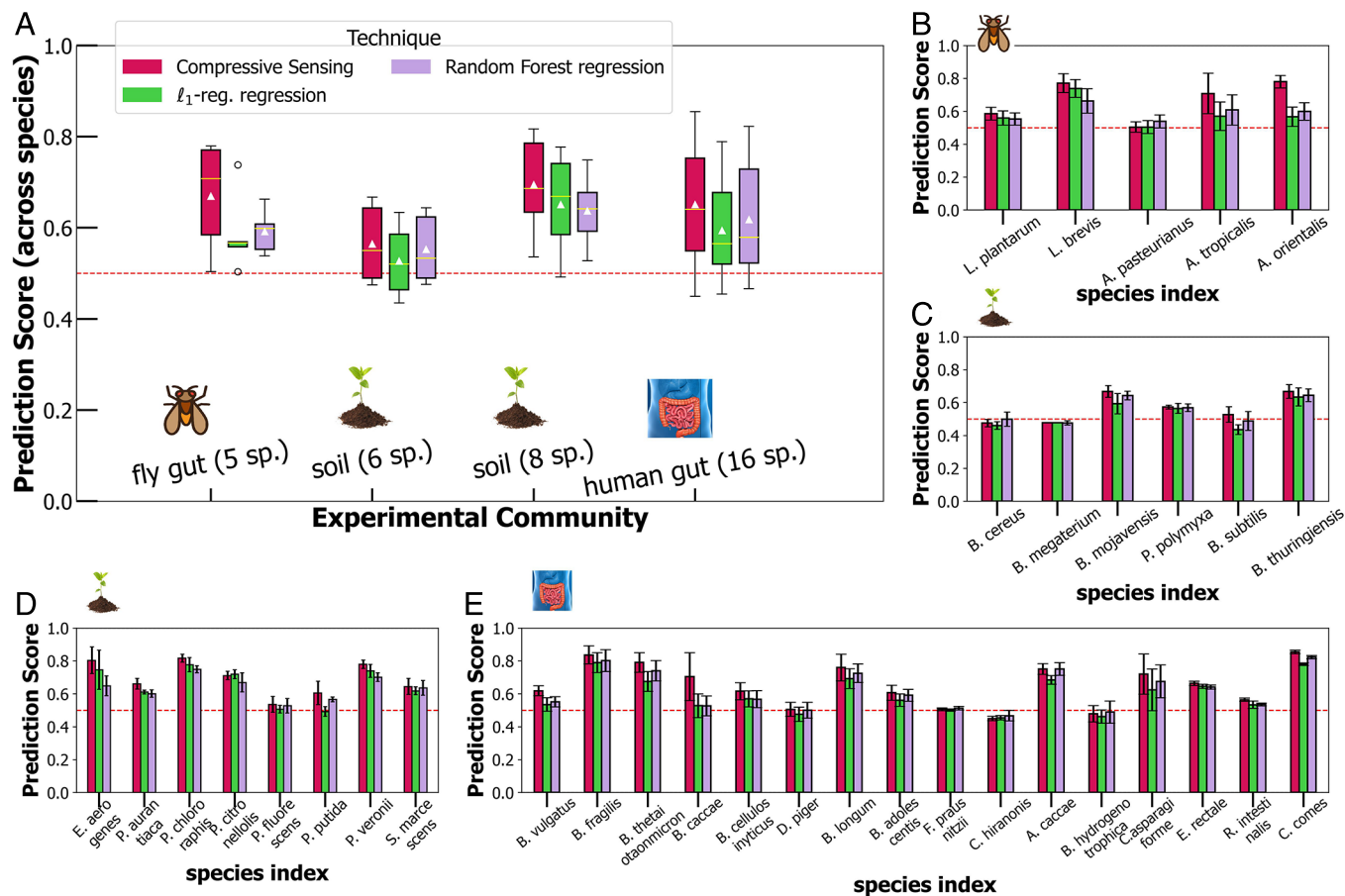
**Table 1. Performance  $P$ -values: One-sided permutation test for various datasets, comparing compressive sensing with a weighted Walsh–Hadamard basis with random forest regression and  $l_1$ -regularized regression**

Dataset	$P$ -value (CS and $l_1$ )	$P$ -value (CS and RF)
In silico (16 species)	$5 \times 10^{-5}$	$10^{-4}$
Fruit fly gut (5 species)	0.06	0.06
Starch soil (6 species)	0.1	0.015
Friedman soil (8 species)	$8 \times 10^{-3}$	$4 \times 10^{-3}$
Clark human gut (16 species)	$1.52 \times 10^{-5}$	$9 \times 10^{-3}$
Clark human gut (SCFAs)	0.0022	0.0002

**Discussion**

A central goal of ecology is to understand emergent properties of microbiomes. In particular, we want to be able to understand and quantitatively predict outcomes of community assembly. Furthermore, community composition may be predictive of community function, and systematically optimizing microbial community function is a critical goal in microbiome engineering. However, an exhaustive search through all possible ways to combine microbial taxa is impractical, and fitting mechanistic models can still require large amounts of data, while also imposing assumptions about community dynamics that may fail to hold for real microbiomes. Here, we establish the method of compressive sensing as a model-agnostic, predictive tool, applicable in situations where experimenters have access to limited data. The success of our approach relies on the assumption that species abundances are sparse in a transformed basis, which we demonstrate explicitly for simulated consumer–resource models, and test in experimental data. While this kind of sparsity appears in many areas of signal processing and has also been applied in an evolutionary genetics context (52–54, 69), the different and complex dynamics underlying microbiome communities means that this sparsity did not a priori need to hold in microbiomes. The fact that it does may open up broad avenues for assembling microbial communities and optimizing their function.

We use this method in silico, in vivo, and in vitro to predict final steady-state abundances of all species. We also show that our method outperforms  $l_1$ -regularization and a widely used machine learning method, random forest regression. We also find that the sparsity of microbial landscapes is not limited to their late-time, steady-state abundances but can be extended to community functions. In Fig. 5, we show that this method can also robustly predict community functions: For the 16-species dataset studied in ref. 64, compressive sensing can be used to accurately predict the amount of butyrate, succinate, lactate, and acetate produced by combinations of these species. By leveraging the interpretability of Walsh–Hadamard coefficients and the roughness of the community-function landscape, we also found that butyrate production, in particular, could be associated with two key species, *Desulfovibrio piger* (DP) and *Anaerostipes caccae* (AC). The absence of DP was found to increase butyrate production, while the absence of AC decreased it. This inference matches what has been found in the study in ref. 64: Hydrogen sulfide production by DP inhibits butyrate production. This reinforces a key assumption that emergent properties of microbiomes, perhaps even incredibly complex community-level functions, may be thought of as effectively arising from only a few degrees of freedom. On the side of microbiome engineering, this result is useful in practice—using sparsely sampled



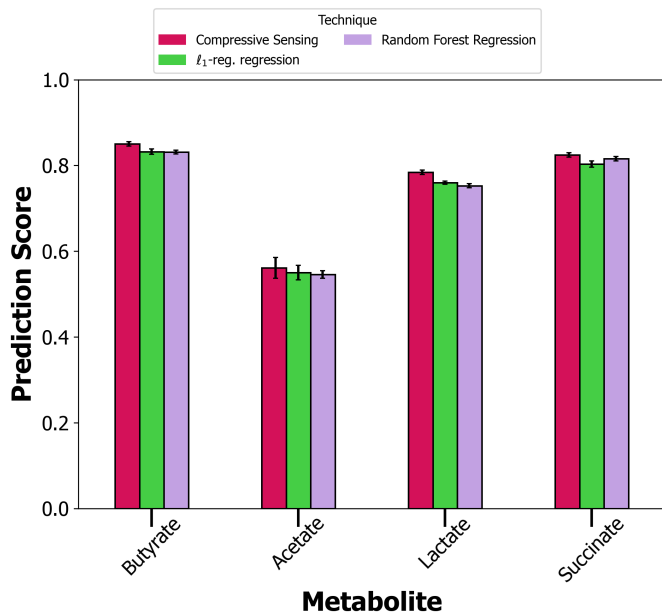
**Fig. 4.** Performance of compressive sensing on four experimental data. We tested the performance of compressive sensing in four real microbiomes, ranging from 5 to 16 species, and in both in vitro and in vivo contexts (*Methods*). (A) At a community level, compressive sensing outperforms both random forest regression and  $l_1$ -regularized regression. (B) For the in vivo gut fly dataset from refs. 28 and 66, compressive sensing does better in predicting a majority (4 of 5) of the species. Permutation test  $P$ -values are 0.09 and 0.06 for comparison against random forest regression and  $l_1$ -regularized regression. (C) For the in vitro 6-species soil community (29, 67), we find that compressive sensing outperforms random forest marginally (permutation test  $P$ -value: 0.1) but does better than  $l_1$ -regularized regression ( $P$ -value: 0.015). This dataset, overall, is harder to predict. (D) Compressive sensing outperforms both random forest regression and  $l_1$ -regularized regression ( $P$ -values:  $4 \times 10^{-3}$  and  $8 \times 10^{-3}$ ) in the 8-species in vitro soil community assembled by Friedman et al. (63). At the species level, compressive sensing improves predictions in 7 of 8 cases when compared to  $l_1$ -regularized regression and in all cases when compared to random forests. (E) Using a 16 species community assembled by Clark et al. (64, 68), we find that compressive sensing outperforms other methods reaching statistically significant values in a permutation test ( $P$ -value:  $9 \times 10^{-3}$  for comparison with random forest regression, and  $1.52 \times 10^{-5}$  for  $l_1$ -regularized regression). At the species level, compressive sensing improves predictions in 15 of 16 cases when compared to  $l_1$ -regularized regression and in 12 cases when compared to random forests. We report all results as averages of 5 runs. The error bars in panels (B–E) indicate SD from the mean across 5 runs. Species names on the x-axis are compiled from the data as provided by authors of the studies.

structure-function landscapes, we may be able to recover entire landscapes and hence reliably look for optimal functions and corresponding communities without taking into account intricate mechanistic models. This formalism does not require time-series data, making it easier to work with data obtained from high-throughput laboratory techniques. Further, it performs well with relative abundances, which is in line with how abundances are widely observed and calculated in this field, though in principle, the approach will also work well on absolute abundances (*SI Appendix, Figs. S8 and S9*).

Recently, there has been work on applying deep-learning models to predict community structures (32, 34), providing a natural point of comparison with our approach. These models can predict well, but their performance can be obscured by large numbers of hyperparameters and a difficulty in interpretability. On the contrary, our method is naturally endowed with interpretation in terms of the sparseness of higher-order landscape interactions, and the compressive sensing algorithm requires only one tunable hyperparameter. In cases where

compressive sensing predicts species abundances accurately, it is reasonable to conclude that higher-order interactions are sparse, and that low-order interactions dominate. Our approach therefore provides a kind of middle-ground between mechanistic models, which may be challenging to parametrize well, and purely statistical models, which are hard to interpret. The identification of sparsity renders the nature of the landscape in real communities both more predictable and better interpretable.

While our approach works well in the simulated and experimental data we applied it to, this method hinges on the assumption that there is a unique steady-state set of community abundances, rather than the potential for multiple equilibria, or more complex late-time dynamics (35–38, 63). There is therefore scope to consider further generalizations in cases where experimental data suggest that these outcomes are possible. However, for cases where experiment suggests a unique map from initial composition to late-time abundances, our approach paves a way to reliably engineer microbial communities.



**Fig. 5.** Compressive sensing is able to directly predict community functions. Along with abundances of species, Clark et al. (64, 68) reported concentrations of four organic acid fermentation products: butyrate, acetate, lactate, and succinate. Of these, butyrate is reported to be particularly important to human health. For these important community functions of metabolite production, we find that compressive sensing does well with  $PS > 0.5$  and also outperforms the methods of  $l_1$ -regularized regression and random forest regression. Permutation  $P$ -values for difference in means between compressive sensing and  $l_1$ -regularized regression and that between compressive sensing and random forest regression indicate statistical significance ( $P$ -values: 0.0022 and 0.0002). We report averages over 3 runs. Cross-validation  $k$  fold was chosen to be 10.

## Methods

**Representation of Abundances in the Walsh–Hadamard Basis.** We consider a combinatorially complete dataset of microbial abundances at steady state. Starting with a set of  $S$  species, there are  $2^S - 1$  of combining them, each assemblage characterised by the species initially present. We consider that the underlying population dynamics of these subcommunities, whatever they are, give rise to a unique steady state for each subcommunity, i.e., we do not consider cases of multistability. Borrowing from the language of genetics and epistasis, we consider a vector of abundances at steady state for each of the  $S$  species. For concreteness, in this section, we consider a set of three species. The set of all the subcommunities, each with a different starting composition is enumerated as  $\vec{\sigma} = \{[000], [001], [010], [011], [100], [101], [110], [111]\}$ . We include the ecologically trivial case of all species being initially absent. The element  $[011]$ , for example, denotes the subcommunity with species 2 and 3 present at the start. For each of these unique conditions, we consider the steady-state abundance (say, cell counts) of each species. This allows for a mapping to be written down explicitly, between the decimal-ordered binary elements and the  $S$  steady-state abundance vectors.

$$\begin{bmatrix} 000 \\ 001 \\ 010 \\ 011 \\ 100 \\ 101 \\ 110 \\ 111 \end{bmatrix} \mapsto \begin{bmatrix} N_{000} \\ N_{001} \\ N_{010} \\ N_{011} \\ N_{100} \\ N_{101} \\ N_{110} \\ N_{111} \end{bmatrix}_0, \begin{bmatrix} N_{000} \\ N_{001} \\ N_{010} \\ N_{011} \\ N_{100} \\ N_{101} \\ N_{110} \\ N_{111} \end{bmatrix}_1, \begin{bmatrix} N_{000} \\ N_{001} \\ N_{010} \\ N_{011} \\ N_{100} \\ N_{101} \\ N_{110} \\ N_{111} \end{bmatrix}_2 \dots \quad [4]$$

Here, the subscripts on the vectors denote the species indices. For each species, we can further build  $2^{S-1}$ -long vectors, since every species is initially present in only the half the total possible  $2^S$  subcommunities. We can define a Walsh–Hadamard transform on this vector of steady-state abundances of a species  $i$ . In particular, we work with a weighted Walsh–Hadamard transform. The weighted

transform is implemented by the matrix:  $VH$  where the matrices  $V$  and  $H$  are generated by the recursion relations:

$$H_{n+1} = \begin{bmatrix} H_n & H_n \\ H_n & -H_n \end{bmatrix}, \quad [5]$$

with  $H_0 = 1$

$$V_{n+1} = \begin{bmatrix} 0.5V_n & 0 \\ 0 & -V_n \end{bmatrix}, \quad [6]$$

with  $V_0 = 1$ .

$V$  is a diagonal weighting matrix that takes into account the order of interactions, to account for averaging over different numbers of terms as a function of the order of interactions (56). In our case, it serves to provide a way to bias inference from limited data toward lower-order interactions (SI Appendix, section 5 and Fig. S7).

**Compressive Sensing and Sparse Recovery Algorithms.** Implementation of compressive sensing involves using optimization algorithms that find a sparse representation of the data using small subsets of observations. After a sparse representation has been found by the algorithm, the rest of the unseen data are computed by taking the inverse transform of the bestguess sparse representation. This inverse transform is the inverse of the matrix,  $\Omega$ , that is expected to sparsify the data. To infer the sparse representations from limited data, we employed the basis pursuit denoising (BPDN) algorithm (55) which is an optimization problem posed as:  $\arg\min_b \left( \frac{1}{2} \|Db - a\|_2^2 + \lambda \|b\|_1 \right)$ . This is

the same as LASSO (Least Absolute Shrinkage and Selection Operator) (53). The idea is to find the  $b$  vector that has the smallest  $l_1$ -norm while keeping the error from observed data as low as possible. To find the sparsest representation, the problem should actually minimize the  $l_0$ -norm instead of  $l_1$ , but this problem is nonconvex and combinatorial (53), and  $l_1$ -norm minimization is a convex approximation to this. This algorithm has only one hyperparameter,  $\lambda$ , which can be tuned according to the expected sparsity of the representations. In real datasets, this value needs to be chosen carefully. We report all results in the main text based on the optimal hyperparameter. We used the SPORCO package (70) in Python to implement BPDN using an alternating direction method of multipliers (ADMM) algorithm (71).

**Problem set-up.** Given a  $S$ -species combinatorially incomplete dataset with  $n$  different subcommunities sampled, we want to predict the abundances of species outside of the sampled experiments. We focus on a species,  $i$ , for which we have  $n$  sampled abundances, and we make predictions for the remaining  $2^{S-1} - n$  subcommunities. We assume, for most species in the pool, that the abundances,  $\vec{a}(\vec{\sigma})$ , are sparse in the  $\vec{\omega}$  basis, i.e.,  $\vec{\beta}(\vec{\omega}) = \Omega \vec{a}(\vec{\sigma})$  is a sparse vector, with only a few significant coefficients.

**Algorithm implementation.** We implement the BPDN/LASSO algorithm with  $D = \Omega_{n \times M}^{-1}$ ; here,  $D$  is the partial  $\Omega^{-1}$  matrix with rows chosen such that the row indices correspond to the decimal representation of the sampled  $n$  subcommunities, and  $M = 2^{S-1}$ . The algorithm then returns the optimal  $\vec{\beta}_{\text{bestguess}}(\vec{\omega})$  which is a solution to the program:

$$\arg\min_{\vec{\beta}} \left( \frac{1}{2} \|\Omega_{n \times M}^{-1} \vec{\beta}(\vec{\omega})_M - \vec{a}_n^i(\vec{\sigma})\|_2^2 + \lambda \|\vec{\beta}(\vec{\omega})_M\|_1 \right). \quad [7]$$

Using Eq. 3 of the main text, we recover the complete  $2^{S-1}$ -long abundance vector for species  $i$ . For the in silico dataset, we compute the prediction score on the out-of-sample data, whose ground truth is known, using Eq. 2:

$$PS = \frac{1}{2 - R^2(\vec{a}_{\text{true}}(\vec{\sigma}), \vec{a}_{\text{bestguess}}(\vec{\sigma}))_{\text{(out-of-sample)}}}. \quad [8]$$

This procedure is repeated for all species in the pool, each with its own set of sampled subcommunities. For experimental datasets, with only partial abundance vectors available, we use  $k$ -fold CV, by dividing the available dataset for each species into  $k$ -folds, and solving the program in Eq. 7 for sampled abundances in the training folds, and computing the prediction score on the stacked abundances corresponding to the data-points in the validation folds.

**$l_1$ -regularized regression.** For a species  $i$  with relative abundances  $\vec{a}(\vec{\sigma})$ , we consider a representation

$$\vec{g} = G_n \vec{a}(\vec{\sigma}).$$

Here,  $G$  is a matrix defined recursively as

$$G_{n+1} = \begin{bmatrix} G_n & 0 \\ -G_n & G_n \end{bmatrix}, \quad [9]$$

with  $G_0 = 1$ . We note that, as elucidated in refs. 52 and 56, this transformation can be thought of as looking at the abundance landscape of a species as a local approximation around a single subcommunity—the one with all species absent. This is akin to a Taylor expansion on the landscape, while the class of Walsh–Hadamard transforms corresponds to a Fourier transform (72) on the landscape and is an approximation over the background of all subcommunities. We implement a similar BPDN algorithm as for compressive sensing with Walsh–Hadamard transform, where the algorithm tries to learn the  $g$  coefficients.

**Bench-marking with Random Forest Regression.** Tree-based ensemble learners, like the random forest regressor and xgboost, are popular choices of supervised learning algorithms, especially when the predictors are not sure to be linearly related to the target variables. We implemented the random forest regressor as available in scikit-learn in Python (73). While there are arguments for using a multioutput regression instead of several single-output ones (74), we worked with fitting a random forest model to each species individually. For data control and an apples-to-apples comparison, we used the same data splits and number of folds as we did in compressive sensing. With random forest regression, there are multiple hyperparameters to consider. To reduce the number of hyperparameters to tune, given the small size of some of the experimental datasets, we used the number of estimators (trees) to be the scikit-learn default (100) and the number of features for selection at every split to be 1/3 of the number of species in each case.

**Predictions on a 5-species Gut Community.** Gould et al. (28) assembled all possible communities from a 5-species pool of bacteria that are known to colonize fly guts in both wild and laboratory conditions. All the 32 subcommunities (termed treatments) were assembled in germ-free flies, each with 48 replicate experiments. They reported colony-forming unit (CFU) counts for each replicate and each treatment. For each replicate, we first computed the relative abundance of each species in that community and then took the average of these relative abundances of each condition (i.e., each subcommunity type). Complete data for this study are available in ref. 66.

**Predictions on a 6-species Starch-degrading Community.** Sanchez-Gorostiaga et al. (29) studied assemblages by considering a pool of 6 amylolytic soil bacterial species. Out of the 63 possible communities, they reported abundance data (in colony-forming units) for 53 communities. We dropped two communities that were inconsistent with labeling and worked with 51 data points. For each community and each replicate, we calculated the total biomass of the community and divided the abundance of each species to find the relative abundance of each species in the community. If multiple replicate experiments were reported, we first computed the relative abundance of a species in each replicate and then averaged this across replicates to find the mean relative abundance. Data for this study are available in ref. 67.

**Predictions on an 8-species Soil Microbiome Community.** Friedman et al. (63) studied a community of 8 heterotrophic soil-dwelling bacteria, reporting their optical densities (ODs) after cross-checking actual cell counts using agar plating. They considered subcommunities which included all the 8 monocultures, all pairwise combinations, all three-species communities, all 8 leave-one-out communities, and the subcommunity with all the species present. They combined species in different ratios (i.e., other than 1:1) and sampled the populations at various times, obtaining time-series data. This difference in starting abundances however did not affect the steady state in most conditions and replicates. However, for 2 subcommunities, they found that the steady-state abundances within the replicates had a much larger variation than what could be accounted for by experimental noise, and one subcommunity displayed bistability. We discarded these experimental conditions, and for all other data

points proceeded to take the average of the final time (where we assume steady state has been attained) abundances across the replicates. Data for this study are available in ref. 75.

**Predictions on a Subset of a 25-species Synthetic Human Gut Community.** Clark et al. (64) studied a community of 25 species that consisted of species spanning all major phyla in the human gut microbiome and also are representative of the major metabolic functions in the gut. Of the  $2^{25} - 1$  ecological communities possible, they sampled  $\sim 600$ . For these subcommunities, they reported read counts and computed the relative abundances for each species using total read counts for each subcommunity. They also found absolute abundances by multiplying the relative abundance with the OD<sub>600</sub> measurement for each sample. As outlined in Methods section of ref. 64, we excluded subcommunities that were flagged as contaminated. Data for this study are available in ref. 68. In such a data-limited case, the performance of any algorithm will be difficult to test. Therefore, we considered a subset of 16 species with the other 9 always being absent. Since there are many possible subsets of 16 species in the background of any consistently absent 9 species, we selected a 16-species pool such that the number of experimental data-points available to us was maximized. The distribution of available data-points with different number of species consistently absent is shown in *SI Appendix, Fig. S10*. With this, we had  $\sim 0.2\%$  of this 16-species landscape available to us. We considered relative abundances in the data, after averaging over the number of replicate experiments for each subcommunity. There were differing number of data points corresponding to each species, unlike the previous datasets. We therefore allowed  $k$  to vary in the cross-validation scheme.  $k$  was set to 3, 5, or 7 depending on the size of the dataset for each species. The authors also reported concentrations of 4 organic acid fermentation products: butyrate, acetate, lactate, and succinate. For these community functions of metabolite production, we used a 10-fold cross-validation approach.

**Microbial Consumer–Resource Models.** We considered a pool of 16 species, with a single externally supplied resource in a chemostat. Communities with cross-feeding that are supplied with a single resource externally have been studied in recent experiments (76–80). In our simulations, the supplied resource is primarily consumed by only a few species. However, in the presence of cross-feeding, there are 20 metabolites present in the community. Different species consume different resources at different rates; this information is encoded in the consumer matrix,  $C$ , with elements  $c_{i\alpha}$ , which characterizes a kind of pairwise interaction, but here between consumer and resource, rather than directly between pairs of species. Such pairwise consumer–resource interactions can give rise to both pairwise interspecific interactions and higher-order interspecific interactions (81). We choose a consumer matrix whose elements are sampled from a binary-gamma distribution (44). We chose a binary-gamma distribution because this allows for positive uptake rates but also allows us to control the sparsity of the consumer matrix. By considering all possible combinations of initial presence-absence of the species, we generated steady-state abundances of all the species, by numerically solving the ODEs for the consumer–resource model (equations given in *SI Appendix, section 3*). Simulations reached steady state when the root mean square of the logarithmic growth rates of the species fell below a threshold, i.e.,  $\text{RMS} \left( \frac{1}{N_i} \frac{dN_i}{dt} \right) < 10^{-2}$ . Further, we verified that the extinct species could not have survived in the community by simulating a reinvasion attempt of the steady-state community.

Thus, for each species, we have  $2^{16}$  abundance data points. Since each species is present in only half the combinations, the effective total data for each species is  $2^{15}$ . Further, we considered the relative abundance of a species in each combination. We calculated the relative abundance by summing the biomass of each species in the subcommunity and dividing this by the total biomass of that subcommunity. Details of parameters used in the simulations are given in *SI Appendix, section 3*.

**Data, Materials, and Software Availability.** The code to reproduce all the analyses in this manuscript is available on GitHub, [https://github.com/sayra-ecoevol/compressed\\_landscapes\\_microbiome](https://github.com/sayra-ecoevol/compressed_landscapes_microbiome). The 5-species fruit fly dataset (28) is available in ref. 66, <https://doi.org/10.5061/dryad.2sr6316>. The 6-species starch-degrading community (29) dataset is avail-

able in ref. 67, <https://github.com/djbajic/structure-function-bacilli>, while the 8-species dataset (63) is available in ref. 75, <https://doi.org/10.5281/zenodo.8176044>. The 25-species dataset with microbes relevant to the human gut (64) is available in ref. 68, <https://github.com/RyanLincolnClark/DesignSyntheticGutMicrobiomeAssemblyFunction>.

**ACKNOWLEDGMENTS.** We would like to thank members of the O'Dwyer group for feedback and comments. We acknowledge funding support from Simons

Foundation Grant #376199 to J.P.O. We acknowledge D. Loudermilk, who licensed the image used in Fig. 1 (CC BY-SA 4.0), and two anonymous reviewers for helpful comments.

Author affiliations: <sup>a</sup>Department of Physics, University of Illinois, Urbana-Champaign, Urbana, IL 61801; <sup>b</sup>Center for Artificial Intelligence and Modeling, Carl R. Woese Institute for Genomic Biology, University of Illinois at Urbana-Champaign, Urbana, IL 61801; <sup>c</sup>Broad Institute of Massachusetts Institute of Technology and Harvard, Cambridge, MA 0214; and <sup>d</sup>Department of Plant Biology, University of Illinois, Urbana-Champaign, Urbana, IL 61801

- P. G. Falkowski, T. Fenchel, E. F. Delong, The microbial engines that drive Earth's biogeochemical cycles. *Science* **320**, 1034–1039 (2008).
- H. W. Paerl, J. L. Pinckney, A mini-review of microbial consortia: Their roles in aquatic production and biogeochemical cycling. *Microbial. Ecol.* **31**, 225–247 (1996).
- H. J. Flint, K. P. Scott, S. H. Duncan, P. Louis, E. Forano, Microbial degradation of complex carbohydrates in the gut. *Gut Microbes* **3**, 289–306 (2012).
- K. Brenner, L. You, F. H. Arnold, Engineering microbial consortia: A new frontier in synthetic biology. *Trends Biotechnol.* **26**, 483–489 (2008).
- B. Olle, Medicines from microbiota. *Nat. Biotechnol.* **31**, 309–315 (2013).
- T. Tanoue *et al.*, A defined commensal consortium elicits CD8 T cells and anti-cancer immunity. *Nature* **565**, 600–605 (2019).
- J. J. Minty *et al.*, Design and characterization of synthetic fungal–bacterial consortia for direct production of isobutanol from cellulosic biomass. *Proc. Natl. Acad. Sci. U.S.A.* **110**, 14592–14597 (2013).
- S. R. Lindemann *et al.*, Engineering microbial consortia for controllable outputs. *ISME J.* **10**, 2077–2084 (2016).
- C. G. Buffie *et al.*, Precision microbiome reconstitution restores bile acid mediated resistance to *Clostridium difficile*. *Nature* **517**, 205–208 (2015).
- J. Hu *et al.*, Design and composition of synthetic fungal–bacterial microbial consortia that improve lignocellulosic enzyme activity. *Biores. Technol.* **227**, 247–255 (2017).
- F. Senne de Oliveira Lino, D. Bajic, J. C. C. Vila, A. Sánchez, M. O. A. Sommer, Complex yeast–bacteria interactions affect the yield of industrial ethanol fermentation. *Nat. Commun.* **12**, 1498 (2021).
- J. K. Weng, X. Li, N. D. Bonawitz, C. Chapple, Emerging strategies of lignin engineering and degradation for cellulosic biofuel production. *Curr. Opin. Biotechnol.* **19**, 166–172 (2008).
- P. Piccardi, B. Vessman, S. Mitri, Toxicity drives facilitation between 4 bacterial species. *Proc. Natl. Acad. Sci. U.S.A.* **116**, 15979–15984 (2019).
- W. Swenson, D. S. Wilson, R. Elias, Artificial ecosystem selection. *Proc. Natl. Acad. Sci. U.S.A.* **97**, 9110–9114 (2000).
- A. B. George, K. S. Korolev, Ecological landscapes guide the assembly of optimal microbial communities. *PLoS Comput. Biol.* **19**, e1010570 (2023).
- R. M. May, Will a large complex system be stable? *Nature* **238**, 413–414 (1972).
- D. S. Maynard, Z. R. Miller, S. Allesina, Predicting coexistence in experimental ecological communities. *Nat. Ecol. Evol.* **4**, 91–100 (2020).
- A. Skwara, P. Lemos-Costa, Z. R. Miller, S. Allesina, Modelling ecological communities when composition is manipulated experimentally. *Methods Ecol. Evol.* **14**, 696–707 (2023).
- M. Barbier, J. F. Arnoldi, G. Bunin, M. Loreau, Generic assembly patterns in complex ecological communities. *Proc. Natl. Acad. Sci. U.S.A.* **115**, 2156–2161 (2018).
- S. Allesina, S. Tang, Stability criteria for complex ecosystems. *Nature* **483**, 205–208 (2012).
- G. Bunin, Ecological communities with Lotka–Volterra dynamics. *Phys. Rev. E* **95**, 042414 (2017).
- J. Grilli *et al.*, Feasibility and coexistence of large ecological communities. *Nat. Commun.* **8**, 14389 (2017).
- G. Barabás, M. J. Michalska-Smith, S. Allesina, The effect of intra- and interspecific competition on coexistence in multispecies communities. *Am. Nat.* **188**, E1–E12 (2016).
- G. Bunin, Interaction patterns and diversity in assembled ecological communities. *arXiv [Preprint]* (2016). <http://arxiv.org/abs/1607.04734> (Accessed 23 August 2023).
- C. A. Serván, J. A. Capitán, J. Grilli, K. E. Morrison, S. Allesina, Coexistence of many species in random ecosystems. *Nat. Ecol. Evol.* **2**, 1237–1242 (2018).
- C. A. Serván, S. Allesina, Tractable models of ecological assembly. *Ecol. Lett.* **24**, 1029–1037 (2021).
- A. F. Ansari, Y. B. S. Reddy, J. Raut, N. M. Dixit, An efficient and scalable top-down method for predicting structures of microbial communities. *Nat. Comput. Sci.* **1**, 619–628 (2021).
- A. L. Gould *et al.*, Microbiome interactions shape host fitness. *Proc. Natl. Acad. Sci. U.S.A.* **115**, E11951–E11960 (2018).
- A. Sanchez-Gorostiza, D. Bajic, M. L. Osborne, J. F. Poyatos, A. Sanchez, High-order interactions distort the functional landscape of microbial consortia. *PLoS Biol.* **17**, e3000550 (2019).
- E. Bairey, E. D. Kelsic, R. Kishony, High-order species interactions shape ecosystem diversity. *Nat. Commun.* **7**, 12285 (2016).
- M. A. Morin, A. J. Morrison, M. J. Harms, R. J. Dutton, Higher-order interactions shape microbial interactions as microbial community complexity increases. *Sci. Rep.* **12**, 22640 (2022).
- M. Baranwal *et al.*, Recurrent neural networks enable design of multifunctional synthetic human gut microbiome dynamics. *eLife* **11**, e73870 (2022).
- C. K. Fisher, P. Mehta, Identifying keystone species in the human gut microbiome from metagenomic timeseries using sparse linear regression. *PLoS ONE* **9**, e102451 (2014).
- S. Michel-Mata, X. W. Wang, Y. Y. Liu, M. T. Angulo, Predicting microbiome compositions from species assemblages through deep learning. *iMeta* **1**, e3 (2022).
- D. B. Amchin, A. Martínez-Calvo, S. S. Datta, Microbial mutualism generates multistable and oscillatory growth dynamics. *bioRxiv [Preprint]* (2022). <https://www.biorxiv.org/content/10.1101/2022.04.19.488807v1> (Accessed 27 April 2023).
- D. R. Amor, C. Ratzke, J. Gore, Transient invaders can induce shifts between alternative stable states of microbial communities. *Sci. Adv.* **6**, eaay8676 (2020).
- R. Debray *et al.*, Priority effects in microbiome assembly. *Nat. Rev. Microbiol.* **20**, 109–121 (2022).
- J. R. Björk *et al.*, Synchrony and idiosyncrasy in the gut microbiome of wild baboons. *Nat. Ecol. Evol.* **6**, 955–964 (2022).
- A. Hastings, C. L. Hom, S. Ellner, P. Turchin, H. C. J. Godfray, Chaos in ecology: Is mother nature a strange attractor? *Annu. Rev. Ecol. Syst.* **24**, 1–33 (1993).
- G. Bunin, Ecological communities with Lotka–Volterra dynamics. *Phys. Rev. E* **95**, 042414 (2017).
- T. Fukami, Historical contingency in community assembly: Integrating niches, species pools, and priority effects. *Annu. Rev. Ecol. Syst.* **46**, 1–23 (2015).
- J. Hu, D. R. Amor, M. Barbier, G. Bunin, J. Gore, Emergent phases of ecological diversity and dynamics mapped in microcosms. *Science* **378**, 85–89 (2022).
- A. Sanchez *et al.*, Directed evolution of microbial communities. *Annu. Rev. Biophys.* **50**, 323–341 (2021).
- C. Y. Chang *et al.*, Engineering complex communities by directed evolution. *Nat. Ecol. Evol.* **5**, 1011–1023 (2021).
- A. Sanchez *et al.*, The community-function landscape of microbial consortia. *Cell Syst.* **14**, 122–134 (2023).
- D. M. Weinreich, Y. Lan, C. S. Wylie, R. B. Heckendorn, Should evolutionary geneticists worry about higher-order epistasis? *Curr. Opin. Genet. Dev.* **23**, 700–707 (2013).
- D. M. Weinreich, Y. Lan, J. Jaffe, R. B. Heckendorn, The influence of higher-order epistasis on biological fitness landscape topography. *J. Stat. Phys.* **172**, 208–225 (2018).
- D. H. Brookes, A. Aghazadeh, J. Listgarten, On the sparsity of fitness functions and implications for learning. *Proc. Natl. Acad. Sci. U.S.A.* **119**, e2109649118 (2022).
- A. Ballal *et al.*, Sparse epistatic patterns in the evolution of terpene synthases. *Mol. Biol. Evol.* **37**, 1907–1924 (2020).
- Z. R. Sailer, M. J. Harms, Detecting high-order epistasis in nonlinear genotype–phenotype maps. *Genetics* **205**, 1079–1088 (2017).
- G. Yang *et al.*, Higher-order epistasis shapes the fitness landscape of a xenobiotic-degrading enzyme. *Nat. Chem. Biol.* **15**, 1120–1128 (2019).
- F. J. Poelwijk, M. Socolik, R. Ranganathan, Learning the pattern of epistasis linking genotype and phenotype in a protein. *Nat. Commun.* **10**, 4213 (2019).
- T. Hastie, R. Tibshirani, M. Wainwright, *Statistical Learning with Sparsity* (Routledge, Boca Raton, FL, ed. 1, 2015).
- E. Candes, M. Wakin, An introduction to compressive sampling. *IEEE Sig. Process. Mag.* **25**, 21–30 (2008).
- S. S. Chen, D. L. Donoho, M. A. Saunders, Atomic decomposition by basis pursuit. *SIAM J. Sci. Comput.* **43**, 129–159 (2006).
- F. J. Poelwijk, V. Krishna, R. Ranganathan, The context-dependence of mutations: A linkage of formalisms. *PLoS Comput. Biol.* **12**, e1004771 (2016).
- S. Yitbarek, J. Guittar, S. A. Knutic, C. B. Ogbunugafor, Deconstructing taxa x taxa x environment interactions in the microbiota: A theoretical examination. *bioRxiv [Preprint]* (2021). <https://www.biorxiv.org/content/10.1101/647156v2> (Accessed 21 February 2023).
- S. Doro, M. A. Herman, On the Fourier transform of a quantitative trait: Implications for compressive sensing. *J. Theor. Biol.* **540**, 110985 (2022).
- R. Marsland, W. Cui, P. Mehta, A minimal model for microbial biodiversity can reproduce experimentally observed ecological patterns. *Sci. Rep.* **10**, 3308 (2020).
- R. Marsland, W. Cui, J. Goldford, P. Mehta, The community simulator: A Python package for microbial ecology. *PLoS ONE* **15**, e0230430 (2020).
- J. E. Nash, J. V. Sutcliffe, River flow forecasting through conceptual models. Part I–A discussion of principles. *J. Hydrol.* **10**, 282–290 (1970).
- R. Tibshirani, Regression shrinkage and selection via the LASSO. *J. R. Stat. Soc.: Ser. B (Methodol.)* **58**, 267–288 (1996).
- J. Friedman, L. M. Higgins, J. Gore, Community structure follows simple assembly rules in microbial microcosms. *Nat. Ecol. Evol.* **1**, 0109 (2017).
- R. L. Clark *et al.*, Design of synthetic human gut microbiome assembly and butyrate production. *Nat. Commun.* **12**, 3254 (2021).
- P. Mehta *et al.*, A high-bias, low-variance introduction to Machine Learning for physicists. *Phys. Rep.* **810**, 1–124 (2019).
- A. L. Gould *et al.*, Data from: Microbiome interactions shape host fitness. Zenodo (2018). [10.5061/dryad.2sr6316](https://doi.org/10.5061/dryad.2sr6316). Accessed 21 June 2022.
- A. Sanchez-Gorostiza, D. Bajic, M. L. Osborne, J. F. Poyatos, A. Sanchez, Data from: High-order interactions distort the functional landscape of microbial consortia. GitHub Repository (2019). <https://github.com/djbajic/structure-function-bacilli>. Accessed 25 August 2022.
- R. L. Clark *et al.*, Data from: Design of synthetic human gut microbiome assembly and butyrate production. GitHub Repository (2021). <https://github.com/RyanLincolnClark/DesignSyntheticGutMicrobiomeAssemblyFunction>. Accessed 19 May 2022.
- J. C. Ye, Compressed sensing MRI: A review from signal processing perspective. *BMC Biomed. Eng.* **1**, 8 (2019).
- B. Wohlberg, “SPORCO: A Python package for standard and convolutional sparse representations” in *Proceedings of the 16th Python in Science Conference* (2017), pp. 1–8.
- S. Boyd, N. Parikh, E. Chu, *Distributed Optimization and Statistical Learning Via the Alternating Direction Method of Multipliers* (Now Publishers Inc., Hanover, MA, 2011).
- E. D. Weinberger, Fourier and Taylor series on fitness landscapes. *Biol. Cybern.* **65**, 321–330 (1991).

73. F. Pedregosa *et al.*, Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.* **12**, 2825–2830 (2011).
74. L. Schmid, A. Gerharz, A. Groll, M. Pauly, Machine learning for multi-output regression: When should a holistic multivariate approach be preferred over separate univariate ones? *arXiv [Preprint]* (2022). <http://arxiv.org/abs/2201.05340>. Accessed 21 February 2023.
75. J. Friedman, L. M. Higgins, J. Gore, Data and code from Friedman, Higgins, and Gore 2017 ("Community structure follows simple assembly rules in microbial microcosms". *Nat. Ecol. Evol.* **1**, 0109 (2017)) (2023). <https://doi.org/10.5281/zenodo.8176044>.
76. J. E. Goldford *et al.*, Emergent simplicity in microbial community assembly. *Science* **361**, 469–474 (2018).
77. M. Gralka, R. Szabo, R. Stocker, O. X. Cordero, Trophic interactions and the drivers of microbial community assembly. *Curr. Biol.* **30**, R1176–R1188 (2020).
78. S. Estrela *et al.*, Functional attractors in microbial community assembly. *Cell Syst.* **13**, 29–42.e7 (2021).
79. T. N. Enke *et al.*, Modular assembly of polysaccharide-degrading marine microbial communities. *Curr. Biol.* **29**, 1528–1535.e6 (2019).
80. M. Dal Bello, H. Lee, A. Goyal, J. Gore, Resource-diversity relationships in bacterial communities reflect the network structure of microbial metabolism. *Nat. Ecol. Evol.* **5**, 1424–1434 (2021).
81. J. P. O'Dwyer, Whence Lotka-Volterra? *Theor. Ecol.* **11**, 441–452 (2018).