

Topological analysis reveals state transitions in human gut and marine bacterial communities

William K. Chang¹, Dave VanInsberghe², and Libusha Kelly¹

¹Systems and Computational Biology Department, Albert Einstein College of Medicine

²Department of Environmental and Civil Engineering, Massachusetts Institute of Technology

March 12, 2020

Abstract

Microbiome dynamics influence the health and functioning of human physiology and the environment and are driven in part by interactions between large numbers of microbial taxa, making large-scale prediction and modeling a challenge. Here, using topological data analysis, we identify states and dynamical features relevant to macroscopic processes. We show that gut disease processes and marine geochemical events are associated with transitions between community states, defined as topological features of the data density. We find a reproducible two-state succession during recovery from cholera in the gut microbiomes of multiple patients, evidence of dynamic stability in the gut microbiome of a healthy human after experiencing diarrhea during travel, and periodic state transitions in a marine *Prochlorococcus* community driven by water column cycling. Our approach bridges small-scale fluctuations in microbiome composition and large-scale changes in phenotype without details of underlying mechanisms, and provides a novel assessment of microbiome stability and its relation to human and environmental health.

Introduction

Complex microbial ecosystems (‘microbiomes’) inhabit a diversity of environments in the biosphere, including the global ocean [47], soil [13], and the human gut [48]. Large-scale alterations in the composition of microbiomes is often associated, whether as driver or consequence, with environmental processes such as seasonal geological cycling and nutrient fluctuations [15]; physiological processes such as menstrual cycles [16]; and clinical phenotypes such as irritable bowel syndrome [2]. Analysis and prediction of the large-scale dynamics of microbiome composition is thus a pressing issue in multiple fields of study.

As with many biological systems, understanding of the dynamics of microbiomes is complicated by their high dimensionality. Numerous variables define the state of a microbiome; these include frequencies of microbial taxa and their genetic alleles, which are decoupled due to genomic plasticity and horizontal gene transfer [36, 38], and environmental conditions such as temperature, pH, and biochemical concentrations. A microbiome thus has a vast number of potential configurations in which it may, in principle, fluctuate on a short time scale. By contrast, systemic phenotypes, such as human gut infections or aquatic algal blooms, persist for much longer than bacterial generation time, and community compositions may be diverse within a phenotype [15]. Furthermore, due to the diverse biology of microbiomes across habitats, it may be desirable to have a quantitative framework that can be generalized across biological systems.

One approach to analyzing microbiome dynamics has been to infer the network of underlying pairwise interactions between taxa by calculating the inverse covariance matrix from time series data, often as a basis for modeling population dynamics using Lotka-Volterra equations [14, 28, 46]. Such approaches are useful for predicting fine-grained taxon-taxon interactions of importance, and are challenged by the compositional nature of microbiome data [44] and possible role of higher-order interactions [3]. Notably, it is impossible to fit Lotka-Volterra models to compositional data

47 without information regarding the total population size [26]. A complementary coarse-grained
48 approach is to cluster samples according to compositional similarity, and conceptualize dynamics
49 as stochastic transitions between clusters [1, 9]. Such approaches can be used to identify large-
50 scale shifts in compositional state, with the implicit assumption that each temporal sample can be
51 assigned to one of a finite number of discrete categories.

52 In our approach to microbiome dynamics, we were motivated by the concept of potential land-
53 scapes in physics. The potential landscape formalism considers a high-dimensional phase space, in
54 which coordinates represent system states, and system dynamics correspond to trajectories through
55 phase space. The dynamics are envisioned as being influenced by features of a landscape in phase
56 space, the height of which corresponds to the value of a potential energy function: for example,
57 local minima of the potential may represent stable states, and valleys probable dynamics of the
58 system. In biology, the potential landscape and related concepts have proved useful in theoretical
59 and experimental studies of ecological dynamics [6, 7, 41]; cell phenotypes in differentiating stem
60 cells [50, 52] and cancer cells [25, 30]; and states of brain activity [20].

61 In principle, potential landscapes predict an inverse relationship between the value of the po-
62 tential and the probability of observing the corresponding system state, and thus between the
63 potential in a region of phase space and the density of observations in that region. In reality,
64 certain landscape features and dynamics may lead to the persistence of transient states and the
65 illusion of stability [22, 35], and strong external perturbations may cause the dynamics to deviate
66 from those predicted by the potential landscape, in particular in biological applications. For ex-
67 ample, perturbations to the gene expression of a differentiating stem cell may cause it to lose or
68 fail to attain a differentiated phenotype [24]. While the potential landscape formalism may not
69 be directly applicable to microbiomes due to the open nature of the system and rapid turnover
70 relative to currently-practical sampling frequency, we speculated that creating a representation of
71 the density of data points in the compositional phase space of microbial ecosystems could lead to
72 useful insights for analyzing, and eventually predicting, microbiome dynamics. Specifically, we hy-
73 pothesized that local maxima of the data density could form a basis by which to infer characteristic
74 metastable states of microbiome composition, allowing the association of observations with states
75 and the representation of dynamics as metastable state transitions while retaining the continuity
76 of the underlying phase space.

77 To characterize features of the microbial phase space, we used topological data analysis (TDA),
78 specifically the Mapper algorithm [39, 45], which has recently found application in microbiome
79 research [31]. TDA is a class of methods for inferring properties of data, represented as a point
80 cloud, in high-dimensional phase-space, that seeks to be robust to factors such as scale and res-
81 olution. Briefly, Mapper represents the underlying distribution of data in a metric space as an
82 undirected graph, where each vertex comprises a non-exclusive subset of data points spanning a
83 patch of phase space. An edge is drawn between each two vertices that share at least one data
84 point (Fig. 1A), representing connectivity between patches. We complement Mapper with a novel
85 graph-theoretical analysis using k-nearest neighbor (kNN) distance to estimate the density of data
86 points over each patch of phase space represented by a vertex, determine local maxima, and define
87 metastable community states (Fig. 1B). In contrast to established methods such as hierarchical
88 clustering, our method preserves the notion of a continuous underlying density distribution, with
89 the states representing a discrete coarse-graining, and recognizes low-density regions of phase space
90 unassociated with any metastable state. In addition, it is possible for a data point to be associated
91 with more than one vertex in the Mapper graph and thus with more than one state, allowing
92 identification of samples that fall between or are in transition between metastable states.

93 We used our method to infer the density and associated topological features of the point clouds
94 for three published microbial time series data sets, two human gut microbiomes—one of stool
95 samples collected from seven cholera patients from disease through recovery [23], one from two
96 mostly healthy adult males [8]—and one of marine *Prochlorococcus* communities spanning multiple
97 depths collected from one site in the Atlantic Ocean (BATS) and one in the Pacific (HOT) [32].
98 (For details on the sampling frequency and duration for each data set analyzed, see Supporting
99 Information Table 1.) We selected these data sets in part to test our method by recapitulating
100 biology known from the original studies, and in part to discover novel features not addressed by
101 prior methods. In both human gut and marine systems, we find that significant physiological
102 and environmental events, including recovery from infection and geochemical cycling, correspond
103 to recurrent successions of state transitions. We show that these successions are an informative
104 coarse-grained view of microbiome dynamics, with implications for the assessment of ecological

105 resilience.

106 Results

107 Dynamics of human gut microbiome recovery from cholera infection

108 We found the cholera phase space to be partitioned by clinical phenotype, i.e. diarrhea or recovery
109 (Fig. 2A). Division of the phase space into states found that vertices within a state tended to consist
110 of either samples from the diarrhea phase or from the recovery phase, rather than a mixture of both
111 (Supporting Information Fig. 6). The original study [23] recognized phases of progression according
112 to equal-time divisions of the diarrhea and recovery periods, respectively, of each patient. Our
113 identification of disease substates, in contrast, is based on community composition and integrated
114 across data from all patients. We found the diarrhea region was further subdivided into two states,
115 2 and 7 (Fig. 2B). Patients C, E, and G occupied state 7 for prolonged durations immediately before
116 clinical recovery; patients A, B, and F stably occupied state 7 for approximately 20 hours, but
117 switched to other states for the last few time points before clinical recovery (Fig. 2C). In the case of
118 patient A, the final few time points were associated with state 5, which represented an intermediate
119 region of the phase space between the diarrhea- and recovery-associated neighborhoods. These
120 results suggest that state 2 constituted a universal ‘early’ diarrhea state, and state 7 a universal
121 ‘late’ diarrhea state, with distinct community compositions. The original study noted taxa which
122 consistently changed in abundance between the start and end of the diarrhea phase, for example
123 *Streptococcus* and *Fusobacterium* [23], here we show that these compositional shifts are observable
124 on the whole-community scale.

125 Generally, patients occupied state 7 for longer than they did state 2, suggesting that the stability
126 of the late state in a given patient influences disease duration. To quantify stability, we calculated a
127 temporal correlation function for each state-patient pair during the diarrhea phase (see Methods).
128 Monotonically decreasing correlation functions indicate metastability, showing that the system
129 transiently occupies a state before transitioning to a different state; slopes become more negative
130 with decreasing stability. While this analysis revealed that all patients transiently occupied state
131 2, with greatest persistence in patient C, patients A, C, and E had non-monotonic correlation
132 functions for state 7, coinciding with prolonged times to recovery compared to the rest of the cohort,
133 with patients B and F exhibiting the expected monotonic decrease (Fig. 2D). This indicated that
134 patients A, C, and E repeatedly entered and exited state 7, suggesting that prolonged diarrhea in
135 these three patients may have been additionally influenced by the instability or inaccessibility of
136 alternative, healthy states, and that (re-)assembly of the healthy microbial community constitutes
137 a non-trivial step in recovery.

138 Dynamics of two healthy adult microbiomes with transient diarrhea

139 In contrast to the cholera data set, the two healthy adult gut microbiome time series from David
140 *et al.* [8] were separated by subject (Fig. 3A). Despite being clinically healthy for most of the
141 observation period, both subjects’ microbiomes experienced perturbations: subject A traveled
142 from his residence in the United States to southeast Asia, twice experiencing traveller’s diarrhea;
143 and subject B, also based in the US, suffered an acute infection by *Salmonella*. Previous studies [8,
144 19] noted that, while the microbiome of A returned to its original state after travel, recovery from
145 *Salmonella* left the microbiome of B in an alternative state. Confirming this, we found that subject
146 A occupied the same regions of phase space before and after travel, while subject B occupied disjoint
147 regions before and after infection. We further found that the post-*Salmonella* samples of subject
148 B distributed over several connected components, showing that the gut microbiome of subject B
149 remained in flux across several distinct compositional substates even after being clinically marked
150 as having recovered (Fig 3B). Division of the phase space into states found that vertices within a
151 state tended to be dominated by samples from a single subject (Supporting Information Fig. 7).

152 The large connected components representing the pre- and post-travel healthy samples of sub-
153 ject A and the pre-*Salmonella* healthy samples of subject B were each divided into several states
154 (Supporting Information Fig. 1), suggesting that the clinical ‘healthy’ phenotype of an individual
155 is a probability over multiple compositionally distinct states. The existence of states in microbiome
156 phase space proposes a novel metric for microbiome resilience: comparing the distribution of sam-
157 ples across states between time windows. Subject A occupied states with identical probability

158 before and after travel, exhibiting resilience; in contrast, subject B post-infection did not restore
159 the pre-infection probability across states, despite some samples sharing states with pre-infection
160 healthy samples (Fig. 4A). Thus, the restoration of the microbial community to a ‘healthy’ state
161 cannot be confirmed with a single time point.

162 Temporal correlation functions further showed that subject A, as well as subject B before
163 infection, repeatedly visited the same set of states; in contrast, subject B after infection transiently
164 occupied several states without repetition (Fig. 4B). This shows that not only did the microbiome
165 of subject B enter an alternative state, or probability across states, post-infection, but that this
166 alternative state was not fully stabilized. It is possible that the pre-infection probability across
167 states was restored in subject B after the end of the observational period.

168 Recurrent seasonal dynamics of *Prochlorococcus* communities in the Pa- 169 cific and Atlantic

170 Compared to the phase spaces of human gut microbiomes, which may be discretized by individual
171 or phenotype, the *Prochlorococcus* phase space was organized by gradients of depth (Fig. 5A)
172 and temperature (Supporting Information Fig. 4), indicating that, in these environments, small
173 changes to environmental conditions result in small changes to community structure. In contrast to
174 the two human gut microbiome data sets, division of the *Prochlorococcus* phase space into states
175 found the mean depth per vertex in each state to vary continuously (Supporting Information
176 Fig. 8). The phase space possessed multiple states (Fig. 5B), with state 4 largely representing
177 shallow fractions of the water column $\leq 100\text{m}$; states 2, 3, and 6 deeper fractions; and state 1
178 intermediate depths. State 5 represented an infrequently-occupied region sampled only by the
179 140m fraction at BATS on January 27, 2004, and by the 125m fraction at HOT on January 31,
180 2008 (Fig. 5C). As such, state 5 possibly constitutes an alternative state for deep water fractions in
181 mid-winter. Communities differing in depth rarely shared compositions, and transitioned between
182 states, in many cases periodically across calendar years (Fig. 5C), showing that some communities
183 experienced abrupt periodic shifts in environmental conditions due to geochemical events.

184 Despite the graduated variation of composition with depth and temperature, the range of
185 compositional dissimilarity across the range of environmental conditions is sufficient to constrain
186 given depth fractions to a neighborhood of phase space, such that shallow- and deep-fraction
187 *Prochlorococcus* communities rarely occupy the same compositional states over time (Fig. 5C).
188 However, it is known that the BATS water column undergoes an annual late winter upwelling [32],
189 intermixing communities that otherwise inhabit different depths, and homogenizing environmental
190 conditions across depths. We predicted that mixing would drive communities at all depths at BATS
191 to converge on a common state, while no convergence would be observed at HOT. Accordingly,
192 we observed a transition to state 1 by all depths at BATS in January of each year. After June,
193 depths 1-20m and 120-200m relax toward states characteristic of shallow and deep depth fractions,
194 respectively, while state 1 persists longer in intermediate depths 40-100m. By contrast, no such
195 upwelling occurs at HOT, and the probability of a given depth fraction occupying any state remains
196 uniform over the calendar year; the distribution is especially stationary for shallow depths (Fig. 5C).
197 This periodicity was also evident in periodic correlation functions for BATS, and non-periodic for
198 HOT (Fig. 5D).

199 Robustness of phase space characterization

200 Given that the data sets analyzed here are among the largest longitudinal microbiome data sets
201 currently available, we asked whether the biological hypotheses could have been obtained from
202 sparser data sets. We focused on our finding that microbiome phase spaces are structured by
203 latent variables representing host phenotypes or environmental conditions, and examined whether
204 this structuring was robust to data rarefaction. We found that the partitioning of the phase
205 space by clinical phenotype in the case of the cholera patients, by subject in the case of the two
206 healthy adult humans, and the gradation by depth in the case of *Prochlorococcus* communities,
207 are robust to all rarefaction tests performed. In the case of cholera patients, nodes remained
208 divided into those representing mostly samples from the diarrhea phase and those representing
209 the recovery phase, with edges being more dense between nodes of the same phenotype than
210 those of different phenotypes (Supporting Information Fig. 3). In the case of the two healthy
211 adult humans, nodes were consistently dominated by samples from one subject, with edges being

212 more dense between nodes representing the same subject than those representing different subjects
213 (Supporting Information Fig. 4). For the *Prochlorococcus* data set, nodes aggregating samples
214 from similar depth fractions were more densely connected than those representing disparate depths
215 (Supporting Information Fig. 5).

216

217 Comparison with hierarchical clustering and principal component analysis 218

219 To compare our method with standard methodologies, we performed hierarchical clustering and
220 principal component analysis (PCA) on the OTU tables for each data set. We found that, while
221 PCA confirmed the global partitioning of data by diarrhea or recovery within the cholera data set,
222 partitioning by subject within the two adult gut microbiomes data set, and gradation by depth
223 within the *Prochlorococcus* data set, it failed to make evident finer-grained features such as the
224 existence of early- and late- diarrhea states in the cholera data set, and the distinction of pre- and
225 post-*Salmonella* states for subject B in the two human gut microbiomes data set. Furthermore, the
226 reduction of dimensionality to two dimensions through PCA made the visual separation between
227 hierarchical clusters unclear, particularly for the cholera and two human gut microbiome data
228 sets, and for the *Prochlorococcus* data set introduced a strong ‘horseshoe’ effect [33] (Supporting
229 Information Fig. 9).

230 Discussion

231 We identified unrecognized dynamics governing large-scale phenotypes in microbial time series data
232 by using TDA to infer the shape of data density from 16S and ITS ribosomal RNA time series data.
233 While analyses from the original studies identified bacterial taxa that were differed in abundance
234 across host phenotypic or environmental states—for example the loss of *Firmicutes* in subject B
235 post-*Salmonella* infection [8]—our method, by contrast, aims to identify transitions between global
236 compositional states defined across all taxa without reference to metadata. Our results reveal the
237 role of latent physiological and environmental variables [34], such as disease phenotype and phase of
238 geochemical cycles, in organizing microbiomes over time. We observed common dynamics across
239 instances of ecological processes in the two gut and one environmental timeseries datasets we
240 studied. Using our approach, one can thus begin to infer general mechanisms that determine
241 large scale phenotypes of clinical and environmental importance. The elements of our method—
242 the definition of a metric phase space using the square root of the Jensen-Shannon divergence,
243 the representation of the phase space using TDA, and the characterization of topological features
244 using the adapted kNN density estimator and shortest graph distance searches—are specifically
245 advantageous for analyzing high-dimensional compositional data. Relative abundances provide
246 incomplete information on a system, and a system may be compositionally stable while remaining
247 dynamic in absolute abundance [49]. Our method can be readily adapted to work with absolute
248 abundance where such data are available. Compared to representational methods such as PCA,
249 our method benefits from using all distance information; and compared to clustering techniques,
250 our method does not require specifying the number of states, such as required in k-means.

251 While subjects in both human gut data sets experienced transient infection by bacterial pathogens,
252 the large-scale dynamics differed between the two groups. We found that multiple cholera patients
253 followed a trajectory of early- to late-stage disease states. In contrast, the two healthy subjects
254 from the year-long data set experienced apparently random jumps between states during *Salmonella*
255 infection and traveler’s diarrhea, respectively, that did not result in the stabilization in a repro-
256 ducible alternate state during the course of disease. This discordance between the two human gut
257 microbiome datasets suggests that microbial infections can potentially be classified into ‘ordered’
258 and ‘disordered’ types. Ordered infections are characterized by a reproducible trajectory through
259 phase space, while disordered infections are characterized by unpredictable progression through
260 phase space. The latter case represents a version of the ‘Anna Karenina principle,’ meaning indi-
261 vidual microbiomes are more dissimilar during a particular perturbation than during health [51],
262 while the former represents an inversion of the principle. Scale is likely important in this dis-
263 tinction: independent of the deterministic or stochastic nature of the perturbation induced by
264 an infection, if its magnitude is smaller than ‘baseline’ fluctuations of the healthy microbiome,
265 variations between individuals will remain the dominant variable in organizing the phase space. If

266 the magnitude of the perturbation is larger, it may overwhelm individual variability and cause the
267 phase space to instead appear organized by phenotype. Thus, data on the variability of healthy
268 microbiomes over time between and within individuals will be crucial to characterizing the impact
269 of a given disease on the microbiome. We also note that our conclusions are influenced by sampling
270 frequency: our method cannot capture dynamics on a shorter time scale than that of sampling,
271 and systems that seem noisy on a particular time scale may have ordered dynamics on longer time
272 scales.

273 Our analysis of the David *et al.* data set shows that the microbiome of a healthy individual
274 transitions between states over time. While key dominant taxa may persist, no single large-scale
275 compositional state defines healthy physiology. However, an individual microbiome may occupy
276 states with the same probability during two separate ‘healthy’ time windows. Integrating the
277 information over time for each of the healthy periods, the physiological phenotype can be inferred
278 to be stable despite the system state being dynamic. Put differently, if one interprets states as
279 microstates of the microbiome composition, a systemic clinical or environmental phenotype could
280 then be regarded as a *macrostate*, and a resilient ‘healthy’ microbiome will remain in a stable
281 macrostate over time.

282 This notion of resilience as identical probability across states before and after a perturbation
283 can be generalized to a notion of dynamic stability, defined as stationary probability across states
284 over time. Dynamically stable microbiomes do not necessarily stabilize within a single state,
285 but revisit a given set of states with fixed probability. Our temporal correlation analysis shows
286 that dynamically stable microbiomes, such as subject A and subject B pre-infection from the
287 study in [8], are characterized by non-monotonic temporal correlation functions, indicating the
288 microbiome revisits the same states over time. In contrast, unstable microbiomes, such as subject
289 B post-infection, exhibit monotonically decaying correlation functions, indicating the microbiome
290 transiently occupies compositional states without recurrence. Dynamical instability can persist
291 after infection even in the microbiome of an individual clinically marked as having recovered from
292 infection, as in the case of subject B, revealing additional nuances to the association between
293 stability and health in human microbiomes. The ability to assess resilience from data in the absence
294 of detailed knowledge of the underlying network of microbe-microbe interactions complements
295 model-based methods that analytically solve for fixed points and linear stability [5]. Alternate
296 means of estimating stability and resilience may be possible, for example by quantifying the degree
297 to which consecutive time points are associated with the same or adjacent Mapper vertices.

298 For the two human gut microbiome data sets, we observe some of the same phenomena as the
299 original studies: for the seven cholera patients, certain taxa were differentially abundant throughout
300 the progression of disease [23]; and for subject B of the two healthy males, that the pre-*Salmonella*
301 microbiome composition was not recovered by the end of the experiment [8]. In the first case, we
302 remark that differential abundance of individual taxa does not necessarily imply the existence of
303 large-scale compositional states consistent across patients and disease phases, such as we describe
304 here. In the second case, we additionally found multiple states in the pre- and post-perturbation
305 healthy phases of both subjects, and showed that restoration of a healthy and resilient microbiome is
306 associated with the recovery not of a specific composition but of a distribution across compositional
307 states.

308 We point out several caveats regarding our method. First, though we defined the phase space
309 using the Jensen-Shannon distance, other metrics may be used, and the results of analysis using
310 different metrics for the same data should be compared in future applications. Second, due to
311 the lack of an established protocol for selecting Mapper hyperparameters, we used a heuristic
312 method to choose their values for our analyses. A more rigorous optimization method is desirable,
313 especially one developed against synthetic data from *de novo* simulations where the ‘ground truth’
314 of the parameters, and thus the shape of the density, are known *a priori*. Third, we use Mapper to
315 create a representation of the density, but question of whether it is effective to analyze microbiome
316 dynamics via the topology of the density in a given case is independent of Mapper and TDA, and
317 other methods may be used. Fourth, we assume the data accurately represent the compositions
318 of the sampled communities, when in fact challenges exist with translating sequencing data into
319 compositions [18, 17]; addressing these challenges is outside the scope of this manuscript.

320 In addition to offering a novel quantitative description of microbiome states and dynamics, we
321 hope our analysis will, in time, facilitate predictive modeling of the dynamics and forecasting of
322 major state transitions in the microbiome. As an example, our approach to identifying states from
323 microbial time series can be used to infer state transition probabilities under different conditions,

324 and thus can serve as a basis for fitting the parameters of Markov chain models [9, 12]. The
325 concept of the potential landscape that motivated our study is closely linked to the theory of
326 critical transition forecasting [6, 7, 29, 40, 42]: as perturbations destabilize a system, it ascends
327 the potential gradient and eventually reaches a tipping point from where it can rapidly enter into
328 an alternative stable state. Topological analyses, in turn, may eventually facilitate characterization
329 of the potential landscape based on past observations, and real-time estimation of its stability and
330 state transition probability. Both of these approaches allow modeling and prediction of major
331 dynamical events without detailed knowledge of underlying mechanisms, and may prove pivotal to
332 understanding complex, data-rich biological systems not limited to microbiomes, but also including,
333 for instance, gene regulatory networks and animal ecosystems.

334 **Methods**

335 **Human gut microbiome data and preprocessing**

336 The publicly available data that we re-analyzed here were generated by David *et al* [8] accessible
337 on the European Nucleotide Archive (ENA) under the accession number ERP006059, and by Hsiao
338 *et al* [23] on the NCBI Short Read Archive (SRA) under the accession number PRJEB6358. The
339 downloaded reads were trimmed with V-xtractor version 2.1 [21] (a HMM scan based method of
340 isolating variable regions from 16S rRNA sequences) to ensure the amplicon sequences could be
341 aligned across consistent fractions of the 16S rRNA variable regions. Trimmed reads were then
342 clustered into OTUs using usearch v9.2.64 [11] with a minimum cluster size of two. Representative
343 sequences from each OTU were classified using mothur v1.36.1 [43] and the RDP reference 16S
344 rRNA sequences v16 [4].

345 ***Prochlorococcus* data**

346 Data from Malstrom *et al* [32] was obtained from the Biological and Chemical Oceanography Data
347 Management Office (<https://www.bco-dmo.org>), accession number 3381.

348 **Mapper**

349 Conceptually, the Mapper algorithm accepts as input a matrix of distances or dissimilarities be-
350 tween data, and aims to represent the shape of the distribution of data points in high-dimensional
351 phase space as an undirected graph. In this graph, vertices represent neighborhoods of phase space
352 spanned by subsets of adjacent data points, and edges represent connectivity between neighbor-
353 hoods. In brief, it does this by dividing the data into overlapping subsets that are similar according
354 to the output of at least one filter function that assigns a scalar value to each data point, perform-
355 ing local clustering on each subset, and representing the result as an undirected graph, where each
356 vertex represents a local cluster of data points, and edges between vertices represent at least one
357 shared data point between clusters.

358 **Distance matrix**

359 We interpreted microbiome relative abundances to be probability distributions, and thus used the
360 square root of the Jensen-Shannon divergence as a metric [27]. However, it is important to note
361 that any other metric can be used in place of the Jensen-Shannon distance, such as the Aitchison
362 distance [37], calculated from centered [28] or isometric [44] log-transformed relative abundances.

363 **Filter functions and binning**

364 For the filter functions used by Mapper to bin data points, we performed principal coordinate
365 analysis (PCoA, also known as classical multidimensional scaling) in two dimensions on the pairwise
366 distance matrix, and used the ranked values of principal coordinates (PCo) 1 and 2 as the first and
367 second filter values for Mapper, following Rizvi *et al.* [39]. PCo ranks are an appropriate filter for
368 our purposes, as it assigns similar filter values to points that are relatively close together in the
369 original phase space. We wish to note that while PCoA leads to loss of information, the following
370 local clustering step is performed using subsets of distances from the original distance matrix, and

371 is thus not affected. The data points were then binned by overlapping intervals of the two ranked
372 principal coordinates. For hyperparameters specifying these bins and their overlaps, see Table 1.

373 Local clustering

374 The algorithm first performs hierarchical clustering from all pairwise distances between data points
375 within a bin of filter values. Then, it creates a histogram of branch lengths using a predefined
376 number of bins, and uses the first empty bin in the histogram as a cutoff value, separating the
377 hierarchical tree into single-linkage clusters. The algorithm thus finds a separation of length scales
378 within each neighborhood of phase space represented by a bin of the filter values. We used the
379 default number of histogram bins, 10, for each data set (Table 1).

380 Creating the undirected Mapper graph

381 The final output is produced by representing each local cluster of data points as a vertex, and
382 drawing an edge between each pair of vertices that share at least one data point. When plotting,
383 the size of each vertex represents the number of data points therein.

384 Selection of hyperparameters

385 The Mapper algorithm is relatively new, and there are currently no standard protocols to optimize
386 the values of the hyperparameters. For our purposes, it was important that the algorithm achieved
387 a sufficiently high resolution in partitioning data, but also adequately represented connections
388 between regions of phase space. We thus used the following heuristic to set the number of intervals
389 and percent overlap for each data set.

- 390 1. The largest vertex in the resultant Mapper graph should represent no more than $\approx 10\%$ of
391 the total number of data points in the set;
- 392 2. the number of connected components representing only one data point should be minimized.

393 We acknowledge that a heuristic determination of appropriate hyperparameter values leaves
394 much to be desired; as such, we recommend future in-depth theoretical explorations of how the
395 Mapper output depends on the choice of hyperparameters.

396 Density estimation

We estimated the inverse density for each vertex by calculating the k -nearest neighbors (kNN)
distance [10] for each constituent data point i :

$$\text{kNN}(i, k) = \frac{\sum_j^k d_{ij}}{k} \quad (1)$$

where d_{ij} is the distance between points i and j , choosing k equal to 10% of the number of samples
in each data set, rounded to the nearest integer. For a vertex V representing n points, we define
its inverse density as

$$D_{\text{inv}}(V) = \frac{\sum_{i \in V} \text{kNN}(i, k)}{n^2} \quad (2)$$

The n^2 term in the denominator compensates for the differing sizes of vertices. Finally, we invert
the inverse density to obtain the estimated density:

$$D(V) = \frac{1}{D_{\text{inv}}} \quad (3)$$

397

398 State assignment

We then defined states as topological features of the density surrounding local maxima of D . We designated each vertex with higher D than its neighbors to be a local maximum of the potential. Connected vertices tied for maximum D were each assigned to be a local maximum. To approximate a gradient, we converted the undirected Mapper graph to a directed graph, with each edge pointing from the the vertex with lower D to the one with higher D . For each non-maximum vertex, we found the graph distance d_g to each local maximum constrained by edge direction. We defined the state B_x of a maximum V_x as the set of vertices V with uniquely shortest graph distance to V_x :

$$V \in B_x \text{ if } d_g(V, V_x) < d_g(V, V_y) \quad (4)$$

399 for all $y \neq x$ and $V_y \in M$, where M is the set of all local maxima (Fig 1B). Vertices equidistant to
400 multiple maxima were defined to be unstable regions unassigned to any state. Multiple connected
401 maxima were defined as belonging to the same state. Notably, one data point may be associated
402 with multiple vertices and states, or an unstable region and at least one state: we interpreted this
403 to mean that the point is near a saddle point separating states, and as the ‘true’ coordinates of the
404 saddle point are unknown, the data point is assigned to *all* such states and/or an unstable region
405 with uniform weight.

406 Calculating the temporal correlation function

Given that a system occupied state B_x at time t , we defined the temporal correlation to be the probability that it will still (or again) occupy state B_x at time $t + \tau$:

$$f_x(t+\tau) = \begin{cases} 1 & \text{if system is associated with state } B_x \text{ at time } t + \tau \\ 0 & \text{otherwise.} \end{cases} \quad (5)$$

$$\text{corr}_x(\tau) = \langle f_x(t + \tau) \rangle \quad (6)$$

407 We calculated the correlation function for each state x visited by a subject during a characteristic
408 period and for all sampled intervals between pairs of samples of length τ , where the subject was
409 in state B_x in the sample at the start of the interval. For the cholera data set, we calculated
410 correlation functions for each state visited by each subject over the disease period. For the data
411 set of two healthy adult males, we calculated correlation functions for each state visited by each
412 subject in each healthy period, either before or after infection. For the *Prochlorococcus* data set,
413 we calculated correlation functions for each state at each depth fraction at either site. Where a
414 data point is associated with multiple states, we weigh the association with each state as $f'_x(t) =$
415 $\frac{1}{p} f_x(t)$, with p the total number of unique states associated with the system at time t , with the
416 unassigned/unstable state regarded as a single distinct state. Notably, this means $f'_x(t + \tau)$ can
417 have values of $1, \frac{1}{2}, \frac{1}{3} \dots$

418 Rarefaction test

419 We created random subsets of each data set representing 90%, 50%, and 10% of the original data
420 points, repeating 10 times for each data set and downsampling ratio. We then created Mapper
421 graphs representing the rarefied data using the same hyperparameters as for each of the full data
422 sets. We colored the vertices to indicate the same features as for the full data sets: for the cholera
423 data set, by fraction of samples belonging to the diarrhea or recovery phase; for the two healthy
424 adult gut microbiomes data set, by fraction of samples obtained from each subject; and for the
425 *Prochlorococcus* data set, by the mean depth from which samples originated. We ordered the
426 vertices by feature value and used a circularized linear layout algorithm, such that vertices with
427 similar feature values are adjacent. Finally, we used shading to display edge densities.

428 Software and data

429 The main repository for the study can be found on GitHub, at [http://github.com/kellylab/
430 microbial-landscapes](http://github.com/kellylab/microbial-landscapes).

431 An open-source implementation of Mapper in R, `TDAmapper`, was used for the main analysis
432 and can be found at <http://github.com/wkc1986/TDAmapper>. This package was forked from the
433 original implemented by Daniel Müllner which is maintained by Paul T. Pearson and can be found
434 at <https://github.com/paultpearson/TDAmapper>.

435 Funding

436 L.K. is supported in part by a Peer Reviewed Cancer Research Program Career Development
437 Award from the United States Department of Defense (CA171019).

438 Author’s contributions

439 W.K.C. designed and performed the analysis. D.V. processed and performed OTU calling on the
440 data from Hsiao *et al.*[23] and David *et al.*[8]. W.K.C., D.V., and L.K. wrote the manuscript.

441 Competing interests

442 The authors declare that there are no competing interests.

443 References

- 444 [1] J. Paul Brooks et al. “Changes in vaginal community state types reflect major shifts in the
445 microbiome”. In: *Microbial Ecology in Health and Disease* 28.1 (Jan. 1, 2017), p. 1303265.
446 ISSN: null. DOI: [10.1080/16512235.2017.1303265](https://doi.org/10.1080/16512235.2017.1303265). URL: [https://doi.org/10.1080/](https://doi.org/10.1080/16512235.2017.1303265)
447 [16512235.2017.1303265](https://doi.org/10.1080/16512235.2017.1303265) (visited on 08/05/2019).
- 448 [2] C. Casén et al. “Deviations in human gut microbiota: a novel diagnostic test for determining
449 dysbiosis in patients with IBS or IBD”. In: *Alimentary Pharmacology & Therapeutics* 42.1
450 (July 2015), pp. 71–83. ISSN: 1365-2036. DOI: [10.1111/apt.13236](https://doi.org/10.1111/apt.13236).
- 451 [3] Hasan Celiker and Jeff Gore. “Clustering in community structure across replicate ecosys-
452 tems following a long-term bacterial evolution experiment”. In: *Nature Communications* 5
453 (Aug. 8, 2014). ISSN: 2041-1723. DOI: [10.1038/ncomms5643](https://doi.org/10.1038/ncomms5643). URL: [http://www.nature.com/](http://www.nature.com/doi/finder/10.1038/ncomms5643)
454 [doifinder/10.1038/ncomms5643](http://www.nature.com/doi/finder/10.1038/ncomms5643) (visited on 12/18/2014).
- 455 [4] James R. Cole et al. “Ribosomal Database Project: data and tools for high throughput rRNA
456 analysis”. In: *Nucleic Acids Research* 42 (D1 Jan. 1, 2014). Publisher: Oxford Academic,
457 pp. D633–D642. ISSN: 0305-1048. DOI: [10.1093/nar/gkt1244](https://doi.org/10.1093/nar/gkt1244). URL: [https://academic.](https://academic.oup.com/nar/article/42/D1/D633/1063201)
458 [oup.com/nar/article/42/D1/D633/1063201](https://academic.oup.com/nar/article/42/D1/D633/1063201) (visited on 03/02/2020).
- 459 [5] Katharine Z. Coyte, Jonas Schluter, and Kevin R. Foster. “The ecology of the microbiome:
460 Networks, competition, and stability”. In: *Science* 350.6261 (Nov. 6, 2015), pp. 663–666. ISSN:
461 0036-8075, 1095-9203. DOI: [10.1126/science.aad2602](https://doi.org/10.1126/science.aad2602). URL: [http://www.sciencemag.](http://www.sciencemag.org/content/350/6261/663)
462 [org/content/350/6261/663](http://www.sciencemag.org/content/350/6261/663) (visited on 11/07/2015).
- 463 [6] L. Dai et al. “Generic Indicators for Loss of Resilience Before a Tipping Point Leading to
464 Population Collapse”. In: *Science* 336.6085 (June 1, 2012), pp. 1175–1177. ISSN: 0036-8075,
465 1095-9203. DOI: [10.1126/science.1219805](https://doi.org/10.1126/science.1219805). URL: [http://www.sciencemag.org/cgi/doi/](http://www.sciencemag.org/cgi/doi/10.1126/science.1219805)
466 [10.1126/science.1219805](http://www.sciencemag.org/cgi/doi/10.1126/science.1219805) (visited on 09/12/2014).
- 467 [7] Vasilis Dakos and Jordi Bascompte. “Critical slowing down as early warning for the onset
468 of collapse in mutualistic communities”. In: *Proceedings of the National Academy of Sciences*
469 111.49 (Dec. 9, 2014), pp. 17546–17551. ISSN: 0027-8424, 1091-6490. DOI: [10.1073/pnas.](https://doi.org/10.1073/pnas.1406326111)
470 [1406326111](https://doi.org/10.1073/pnas.1406326111). URL: [http://www.pnas.org/lookup/doi/10.1073/pnas.](http://www.pnas.org/lookup/doi/10.1073/pnas.1406326111)
471 [1406326111](http://www.pnas.org/lookup/doi/10.1073/pnas.1406326111) (visited on 11/08/2016).
- 472 [8] Lawrence A. David et al. “Host lifestyle affects human microbiota on daily timescales”. In:
473 *Genome Biology* 15 (2014), R89. ISSN: 1474-760X. DOI: [10.1186/gb-2014-15-7-r89](https://doi.org/10.1186/gb-2014-15-7-r89). URL:
474 <http://dx.doi.org/10.1186/gb-2014-15-7-r89> (visited on 08/12/2016).
- 475 [9] Daniel B. DiGiulio et al. “Temporal and spatial variation of the human microbiota during
476 pregnancy”. In: *Proceedings of the National Academy of Sciences* 112.35 (Sept. 1, 2015),
477 pp. 11060–11065. ISSN: 0027-8424, 1091-6490. DOI: [10.1073/pnas.1502875112](https://doi.org/10.1073/pnas.1502875112). URL: <https://www.pnas.org/content/112/35/11060>
478 [/www.pnas.org/content/112/35/11060](https://www.pnas.org/content/112/35/11060) (visited on 01/03/2019).
- 479 [10] Richard O. Duda, Peter E. Hart, and David G. Stork. *Pattern classification*. Google-Books-
480 ID: YoxQAAAAMAAJ. Wiley, 2001. 688 pp. ISBN: 978-0-471-05669-0.

- 481 [11] Robert C. Edgar. “Search and clustering orders of magnitude faster than BLAST”. In: *Bioinformatics* 26.19 (Oct. 1, 2010), pp. 2460–2461. ISSN: 1367-4803. DOI: [10.1093/bioinformatics/btq461](https://doi.org/10.1093/bioinformatics/btq461). URL: <https://academic.oup.com/bioinformatics/article/26/19/2460/230188/Search-and-clustering-orders-of-magnitude-faster> (visited on 05/01/2017).
- 482
483
484
- 485 [12] Mathieu Faure and Sebastian J. Schreiber. “Quasi-stationary distributions for randomly per-
486 turbed dynamical systems”. In: *The Annals of Applied Probability* 24.2 (Apr. 2014), pp. 553–
487 598. ISSN: 1050-5164. DOI: [10.1214/13-AAP923](https://doi.org/10.1214/13-AAP923). URL: <http://projecteuclid.org/euclid.aoap/1394465365> (visited on 10/31/2018).
- 488
- 489 [13] N. Fierer and R. B. Jackson. “The diversity and biogeography of soil bacterial communities”.
490 In: *Proceedings of the National Academy of Sciences* 103.3 (Jan. 17, 2006), pp. 626–631. ISSN:
491 0027-8424, 1091-6490. DOI: [10.1073/pnas.0507535103](https://doi.org/10.1073/pnas.0507535103). URL: <http://www.pnas.org/cgi/doi/10.1073/pnas.0507535103> (visited on 08/05/2019).
- 492
- 493 [14] Jonathan Friedman and Eric J. Alm. “Inferring Correlation Networks from Genomic Survey
494 Data”. In: *PLOS Computational Biology* 8.9 (Sept. 20, 2012), e1002687. ISSN: 1553-7358.
495 DOI: [10.1371/journal.pcbi.1002687](https://doi.org/10.1371/journal.pcbi.1002687). URL: <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1002687> (visited on 09/04/2018).
- 496
- 497 [15] Jed A. Fuhrman, Jacob A. Cram, and David M. Needham. “Marine microbial community
498 dynamics and their ecological interpretation”. In: *Nature Reviews Microbiology* 13.3 (Mar.
499 2015), pp. 133–146. ISSN: 1740-1526. DOI: [10.1038/nrmicro3417](https://doi.org/10.1038/nrmicro3417). URL: <http://www.nature.com/nrmicro/journal/v13/n3/abs/nrmicro3417.html> (visited on 03/03/2015).
- 500
- 501 [16] Pawel Gajer et al. “Temporal Dynamics of the Human Vaginal Microbiota”. In: *Science
502 Translational Medicine* 4.132 (May 2, 2012), 132ra52–132ra52. ISSN: 1946-6234, 1946-6242.
503 DOI: [10.1126/scitranslmed.3003605](https://doi.org/10.1126/scitranslmed.3003605). URL: <http://stm.sciencemag.org/content/4/132/132ra52> (visited on 02/04/2016).
- 504
- 505 [17] Gregory B. Gloor et al. “Microbiome Datasets Are Compositional: And This Is Not Optional”.
506 In: *Frontiers in Microbiology* 8 (2017). ISSN: 1664-302X. DOI: [10.3389/fmicb.2017.02224](https://doi.org/10.3389/fmicb.2017.02224).
507 URL: <https://www.frontiersin.org/articles/10.3389/fmicb.2017.02224/full?report=reader> (visited on 07/11/2019).
- 508
- 509 [18] Gregory Brian Gloor et al. “Compositional uncertainty should not be ignored in high-
510 throughput sequencing data analysis”. In: *Austrian Journal of Statistics* 45.4 (July 28, 2016),
511 pp. 73–87. ISSN: 1026-597X. DOI: [10.17713/ajs.v45i4.122](https://doi.org/10.17713/ajs.v45i4.122). URL: <https://www.ajs.or.at/index.php/ajs/article/view/vol45-4-5> (visited on 07/11/2019).
- 512
- 513 [19] Didier Gonze et al. “Microbial communities as dynamical systems”. In: *Current Opinion in
514 Microbiology* 44 (Aug. 1, 2018), pp. 41–49. ISSN: 1369-5274. DOI: [10.1016/j.mib.2018.07.004](https://doi.org/10.1016/j.mib.2018.07.004). URL: <http://www.sciencedirect.com/science/article/pii/S1369527418300092>
515 (visited on 07/24/2018).
- 516
- 517 [20] Shi Gu et al. “The Energy Landscape of Neurophysiological Activity Implicit in Brain Net-
518 work Structure”. In: *Scientific Reports* 8.1 (Feb. 6, 2018), p. 2507. ISSN: 2045-2322. DOI:
519 [10.1038/s41598-018-20123-8](https://doi.org/10.1038/s41598-018-20123-8). URL: <https://www.nature.com/articles/s41598-018-20123-8> (visited on 11/19/2018).
- 520
- 521 [21] Martin Hartmann et al. “V-Xtractor: an open-source, high-throughput software tool to iden-
522 tify and extract hypervariable regions of small subunit (16S/18S) ribosomal RNA gene se-
523 quences”. In: *Journal of Microbiological Methods* 83.2 (Nov. 2010), pp. 250–253. ISSN: 1872-
524 8359. DOI: [10.1016/j.mimet.2010.08.008](https://doi.org/10.1016/j.mimet.2010.08.008).
- 525
- 526 [22] Alan Hastings et al. “Transient phenomena in ecology”. In: *Science* 361.6406 (Sept. 7, 2018),
527 eaat6412. ISSN: 0036-8075, 1095-9203. DOI: [10.1126/science.aat6412](https://doi.org/10.1126/science.aat6412). URL: <http://science.sciencemag.org/content/361/6406/eaat6412> (visited on 09/11/2018).
- 528
- 529 [23] Ansel Hsiao et al. “Members of the human gut microbiota involved in recovery from *Vibrio*
530 *cholerae* infection”. In: *Nature* 515.7527 (Nov. 20, 2014), pp. 423–426. ISSN: 0028-0836. DOI:
531 [10.1038/nature13738](https://doi.org/10.1038/nature13738). URL: <http://www.nature.com/nature/journal/v515/n7527/full/nature13738.html> (visited on 02/11/2016).
- 532
- 533 [24] Sui Huang. “The molecular and mathematical basis of Waddington’s epigenetic landscape: A
534 framework for post-Darwinian biology?” In: *BioEssays* 34.2 (2012), pp. 149–157. ISSN: 1521-
535 1878. DOI: [10.1002/bies.201100031](https://doi.org/10.1002/bies.201100031). URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/bies.201100031> (visited on 01/14/2020).

- 536 [25] Sui Huang, Ingemar Ernberg, and Stuart Kauffman. “Cancer attractors: A systems view of
537 tumors from a gene network dynamics and developmental perspective”. In: *Seminars in cell*
538 *& developmental biology* 20.7 (Sept. 2009), pp. 869–876. ISSN: 1084-9521. DOI: [10.1016/j.](https://doi.org/10.1016/j.semcdb.2009.07.003)
539 [semcdb.2009.07.003](https://doi.org/10.1016/j.semcdb.2009.07.003). URL: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2754594/>.
- 540 [26] E Jarauta-Bragulat and J J Egozcue. “Approaching predator-prey Lotka-Volterra equations
541 by simplicial linear differential equations”. In: *Proceedings of the 4th International Workshop*
542 *on Compositional Data Analysis* (2011), p. 9.
- 543 [27] Omry Koren et al. “A Guide to Enterotypes across the Human Body: Meta-Analysis of Micro-
544 bial Community Structures in Human Microbiome Datasets”. In: *PLOS Computational Bi-*
545 *ology* 9.1 (Jan. 10, 2013), e1002863. ISSN: 1553-7358. DOI: [10.1371/journal.pcbi.1002863](https://doi.org/10.1371/journal.pcbi.1002863).
546 URL: [http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.](http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1002863)
547 [1002863](http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1002863) (visited on 04/27/2017).
- 548 [28] Zachary D. Kurtz et al. “Sparse and Compositionally Robust Inference of Microbial Ecological
549 Networks”. In: *PLOS Computational Biology* 11.5 (May 7, 2015), e1004226. ISSN: 1553-7358.
550 DOI: [10.1371/journal.pcbi.1004226](https://doi.org/10.1371/journal.pcbi.1004226). URL: [https://journals.plos.org/ploscompbiol/](https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1004226)
551 [article?id=10.1371/journal.pcbi.1004226](https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1004226) (visited on 04/30/2019).
- 552 [29] Ingrid A. van de Leemput et al. “Critical slowing down as early warning for the onset and
553 termination of depression”. In: *Proceedings of the National Academy of Sciences* 111.1 (Jan. 7,
554 2014), pp. 87–92. ISSN: 0027-8424, 1091-6490. DOI: [10.1073/pnas.1312114110](https://doi.org/10.1073/pnas.1312114110). URL: [http:](http://www.pnas.org/content/111/1/87)
555 [//www.pnas.org/content/111/1/87](http://www.pnas.org/content/111/1/87) (visited on 01/29/2016).
- 556 [30] Qin Li et al. “Dynamics inside the cancer cell attractor reveal cell heterogeneity, limits of
557 stability, and escape”. In: *Proceedings of the National Academy of Sciences of the United*
558 *States of America* 113.10 (Mar. 8, 2016), pp. 2672–2677. ISSN: 1091-6490. DOI: [10.1073/](https://doi.org/10.1073/pnas.1519210113)
559 [pnas.1519210113](https://doi.org/10.1073/pnas.1519210113).
- 560 [31] Tianhua Liao et al. “tmap: an integrative framework based on topological data analysis
561 for population-scale microbiome stratification and association studies”. In: *Genome Biology*
562 20.1 (Dec. 23, 2019), p. 293. ISSN: 1474-760X. DOI: [10.1186/s13059-019-1871-4](https://doi.org/10.1186/s13059-019-1871-4). URL:
563 <https://doi.org/10.1186/s13059-019-1871-4> (visited on 01/07/2020).
- 564 [32] Rex R. Malmstrom et al. “Temporal dynamics of *Prochlorococcus* ecotypes in the Atlantic
565 and Pacific oceans”. In: *The ISME Journal* 4.10 (Oct. 2010), pp. 1252–1264. ISSN: 1751-7362.
566 DOI: [10.1038/ismej.2010.60](https://doi.org/10.1038/ismej.2010.60). URL: [http://www.nature.com/ismej/journal/v4/n10/](http://www.nature.com/ismej/journal/v4/n10/full/ismej201060a.html)
567 [full/ismej201060a.html](http://www.nature.com/ismej/journal/v4/n10/full/ismej201060a.html) (visited on 07/13/2016).
- 568 [33] James T. Morton et al. “Uncovering the Horseshoe Effect in Microbial Analyses”. In: *mSys-*
569 *tems* 2.1 (Feb. 28, 2017). Publisher: American Society for Microbiology Journals Section:
570 Opinion/Hypothesis. ISSN: 2379-5077. DOI: [10.1128/mSystems.00166-16](https://doi.org/10.1128/mSystems.00166-16). URL: [https :](https://msystems.asm.org/content/2/1/e00166-16)
571 [//msystems.asm.org/content/2/1/e00166-16](https://msystems.asm.org/content/2/1/e00166-16) (visited on 03/03/2020).
- 572 [34] Lan Huong Nguyen and Susan Holmes. “Bayesian Unidimensional Scaling for visualizing
573 uncertainty in high dimensional datasets with latent ordering of observations”. In: *BMC*
574 *Bioinformatics* 18.10 (Sept. 13, 2017), p. 394. ISSN: 1471-2105. DOI: [10.1186/s12859-017-](https://doi.org/10.1186/s12859-017-1790-x)
575 [1790-x](https://doi.org/10.1186/s12859-017-1790-x). URL: <https://doi.org/10.1186/s12859-017-1790-x> (visited on 01/03/2019).
- 576 [35] Ben C. Nolting and Karen C. Abbott. “Balls, cups, and quasi-potentials: quantifying stability
577 in stochastic systems”. In: *Ecology* 97.4 (Apr. 1, 2016), pp. 850–864. ISSN: 1939-9170. DOI:
578 [10.1890/15-1047.1](https://doi.org/10.1890/15-1047.1). URL: [http://onlinelibrary.wiley.com/doi/10.1890/15-](http://onlinelibrary.wiley.com/doi/10.1890/15-1047.1/abstract)
579 [1047.1/abstract](http://onlinelibrary.wiley.com/doi/10.1890/15-1047.1/abstract) (visited on 03/07/2018).
- 580 [36] Howard Ochman, Jeffrey G. Lawrence, and Eduardo A. Groisman. “Lateral gene transfer
581 and the nature of bacterial innovation”. In: *Nature* 405.6784 (May 18, 2000), pp. 299–304.
582 ISSN: 0028-0836. DOI: [10.1038/35012500](https://doi.org/10.1038/35012500). URL: [http://www.nature.com/nature/journal/](http://www.nature.com/nature/journal/v405/n6784/full/405299a0.html)
583 [v405/n6784/full/405299a0.html](http://www.nature.com/nature/journal/v405/n6784/full/405299a0.html) (visited on 09/01/2015).
- 584 [37] Vera Pawlowsky-Glahn, Juan José Egozcue, and Raimon Tolosana-Delgado. *Modeling and*
585 *Analysis of Compositional Data*. Google-Books-ID: eG25BgAAQBAJ. John Wiley & Sons,
586 Feb. 17, 2015. 273 pp. ISBN: 978-1-119-00313-7.
- 587 [38] Martin F. Polz, Eric J. Alm, and William P. Hanage. “Horizontal Gene Transfer and the
588 Evolution of Bacterial and Archaeal Population Structure”. In: *Trends in genetics : TIG*
589 29.3 (Mar. 2013), pp. 170–175. ISSN: 0168-9525. DOI: [10.1016/j.tig.2012.12.006](https://doi.org/10.1016/j.tig.2012.12.006). URL:
590 <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3760709/> (visited on 08/14/2017).

- 591 [39] Abbas H. Rizvi et al. “Single-cell topological RNA-Seq analysis reveals insights into cellular
592 differentiation and development”. In: *Nature biotechnology* 35.6 (June 2017), pp. 551–560.
593 ISSN: 1087-0156. DOI: [10.1038/nbt.3854](https://doi.org/10.1038/nbt.3854). URL: [https://www.ncbi.nlm.nih.gov/pmc/
594 articles/PMC5569300/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5569300/) (visited on 11/15/2017).
- 595 [40] M. Scheffer et al. “Anticipating Critical Transitions”. In: *Science* 338.6105 (Oct. 19, 2012),
596 pp. 344–348. ISSN: 0036-8075, 1095-9203. DOI: [10.1126/science.1225244](https://doi.org/10.1126/science.1225244). URL: [http:
597 //www.sciencemag.org/cgi/doi/10.1126/science.1225244](http://www.sciencemag.org/cgi/doi/10.1126/science.1225244) (visited on 01/24/2015).
- 598 [41] Marten Scheffer et al. “Early-warning signals for critical transitions”. In: *Nature* 461.7260
599 (Sept. 3, 2009), pp. 53–59. ISSN: 0028-0836, 1476-4687. DOI: [10.1038/nature08227](https://doi.org/10.1038/nature08227). URL:
600 <http://www.nature.com/doi/10.1038/nature08227> (visited on 02/09/2016).
- 601 [42] Marten Scheffer et al. “Generic Indicators of Ecological Resilience: Inferring the Chance of a
602 Critical Transition”. In: *Annual Review of Ecology, Evolution, and Systematics* 46.1 (2015),
603 pp. 145–167. DOI: [10.1146/annurev-ecolsys-112414-054242](https://doi.org/10.1146/annurev-ecolsys-112414-054242). URL: [http://dx.doi.org/
604 10.1146/annurev-ecolsys-112414-054242](http://dx.doi.org/10.1146/annurev-ecolsys-112414-054242) (visited on 11/19/2015).
- 605 [43] Patrick D. Schloss et al. “Introducing mothur: Open-Source, Platform-Independent, Community-
606 Supported Software for Describing and Comparing Microbial Communities”. In: *Applied and
607 Environmental Microbiology* 75.23 (Dec. 1, 2009), pp. 7537–7541. ISSN: 0099-2240, 1098-5336.
608 DOI: [10.1128/AEM.01541-09](https://doi.org/10.1128/AEM.01541-09). URL: <https://aem.asm.org/content/75/23/7537> (visited
609 on 08/08/2019).
- 610 [44] Justin D. Silverman et al. “A phylogenetic transform enhances analysis of compositional
611 microbiota data”. In: *eLife* 6 (Feb. 15, 2017), e21887. ISSN: 2050-084X. DOI: [10.7554/eLife.
612 21887](https://doi.org/10.7554/eLife.21887). URL: <https://elifesciences.org/articles/21887> (visited on 07/12/2017).
- 613 [45] Gurjeet Singh, Facundo Mémoli, and Gunnar Carlsson. “Topological Methods for the Analysis
614 of High Dimensional Data Sets and 3D Object Recognition”. In: *Eurographics Symposium on
615 Point-Based Graphics* (2007), p. 11.
- 616 [46] Richard R. Stein et al. “Ecological Modeling from Time-Series Inference: Insight into Dy-
617 namics and Stability of Intestinal Microbiota”. In: *PLoS Comput Biol* 9.12 (Dec. 12, 2013),
618 e1003388. DOI: [10.1371/journal.pcbi.1003388](https://doi.org/10.1371/journal.pcbi.1003388). URL: [http://dx.doi.org/10.1371/
619 journal.pcbi.1003388](http://dx.doi.org/10.1371/journal.pcbi.1003388) (visited on 12/18/2014).
- 620 [47] Curtis A. Suttle. “Marine viruses — major players in the global ecosystem”. In: *Nature Re-
621 views Microbiology* 5.10 (Oct. 2007), pp. 801–812. ISSN: 1740-1526. DOI: [10.1038/nrmicro1750](https://doi.org/10.1038/nrmicro1750).
622 URL: [http://www.nature.com/nrmicro/journal/v5/n10/full/nrmicro1750.html#B2
623](http://www.nature.com/nrmicro/journal/v5/n10/full/nrmicro1750.html#B2) (visited on 03/23/2017).
- 624 [48] Peter J. Turnbaugh et al. “The human microbiome project”. In: *Nature* 449.7164 (Oct. 18,
625 2007), pp. 804–810. ISSN: 1476-4687. DOI: [10.1038/nature06244](https://doi.org/10.1038/nature06244).
- 626 [49] Doris Vandeputte et al. “Quantitative microbiome profiling links gut community variation
627 to microbial load”. In: *Nature* 551.7681 (Nov. 2017), pp. 507–511. ISSN: 1476-4687. DOI:
628 [10.1038/nature24460](https://doi.org/10.1038/nature24460). URL: <https://www.nature.com/articles/nature24460> (visited
629 on 01/03/2020).
- 630 [50] C. H. Waddington. *The Strategy Of The Genes*. 1957. URL: [http://archive.org/details/
631 in.ernet.dli.2015.547782](http://archive.org/details/in.ernet.dli.2015.547782) (visited on 08/05/2019).
- 632 [51] Jesse R. Zaneveld, Ryan McMinds, and Rebecca Vega Thurber. “Stress and stability: applying
633 the Anna Karenina principle to animal microbiomes”. In: *Nature Microbiology* 2.9 (Sept.
634 2017), p. 17121. ISSN: 2058-5276. DOI: [10.1038/nmicrobiol.2017.121](https://doi.org/10.1038/nmicrobiol.2017.121). URL: [https://www.
635 nature.com/articles/nmicrobiol2017121](https://www.nature.com/articles/nmicrobiol2017121) (visited on 11/27/2018).
- 636 [52] J. X. Zhou et al. “Quasi-potential landscape in complex multi-stable systems”. In: *Journal of
637 The Royal Society Interface* 9.77 (Dec. 7, 2012), pp. 3539–3553. ISSN: 1742-5689, 1742-5662.
638 DOI: [10.1098/rsif.2012.0434](https://doi.org/10.1098/rsif.2012.0434). URL: [http://rsif.royalsocietypublishing.org/cgi/
639 doi/10.1098/rsif.2012.0434](http://rsif.royalsocietypublishing.org/cgi/doi/10.1098/rsif.2012.0434) (visited on 10/11/2018).

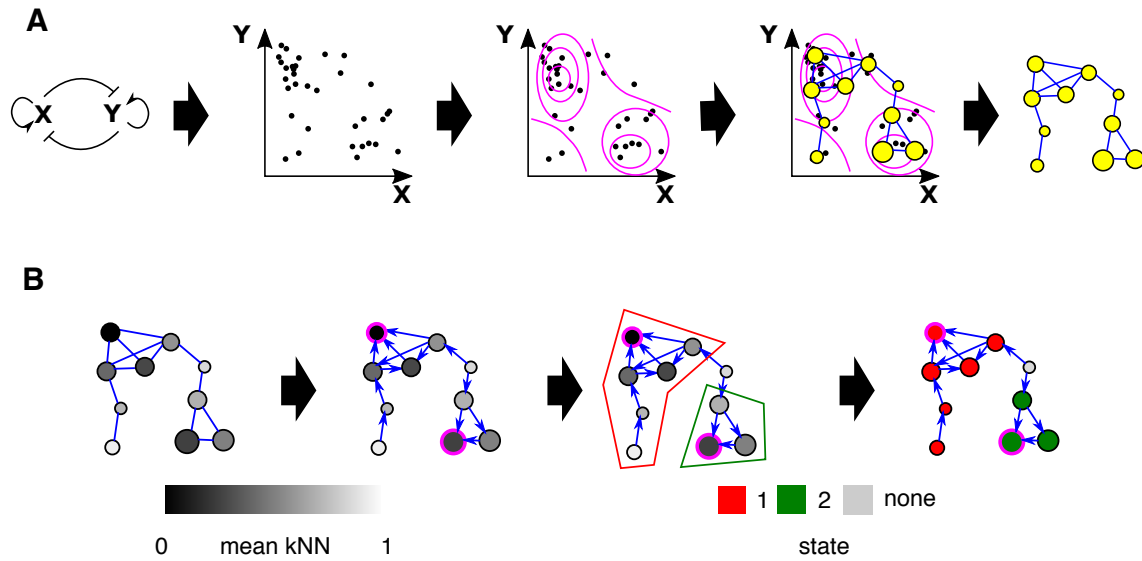


Figure 1: Using Mapper to characterize the microbial phase space. **A.** Cartoon of use of the Mapper algorithm to infer the probability density of a toy ecosystem. The mutually antagonistic interaction between species X and Y leads to denser sampling of the phase space where either X or Y is abundant and the other is rare than in other regions; configurations in which X and Y are similar in abundance are unstable, as small uncertainties in numerical advantage will eventually lead to the dominance of one species over the other. Mapper infers a ‘skeleton’ of density from the data represented as a point cloud. This representation preserves major features of the density such as the two densely-sampled clusters separated by a sparsely-sampled region. Size of vertices indicates number of data points aggregated in each vertex. **B.** Identification of local maxima and metastable states in the Mapper graph shown in A. Data density for each vertex is estimated by the inverse of the mean kNN distance (see Methods) for samples associated with that vertex. Shading indicates mean kNN distance over all data points included in a vertex. The graph is converted to a directed graph, with each edge pointing in the direction of increasing estimated density. A local maximum, highlighted in pink, is defined as a vertex that has higher density than all its neighbors. Finally, the state associated with a local maximum is defined as the set of vertices that have uniquely shortest directed graph distance to that maximum. Non-maxima vertices with equal graph distances to multiple local maxima are unassociated with any state (grey).

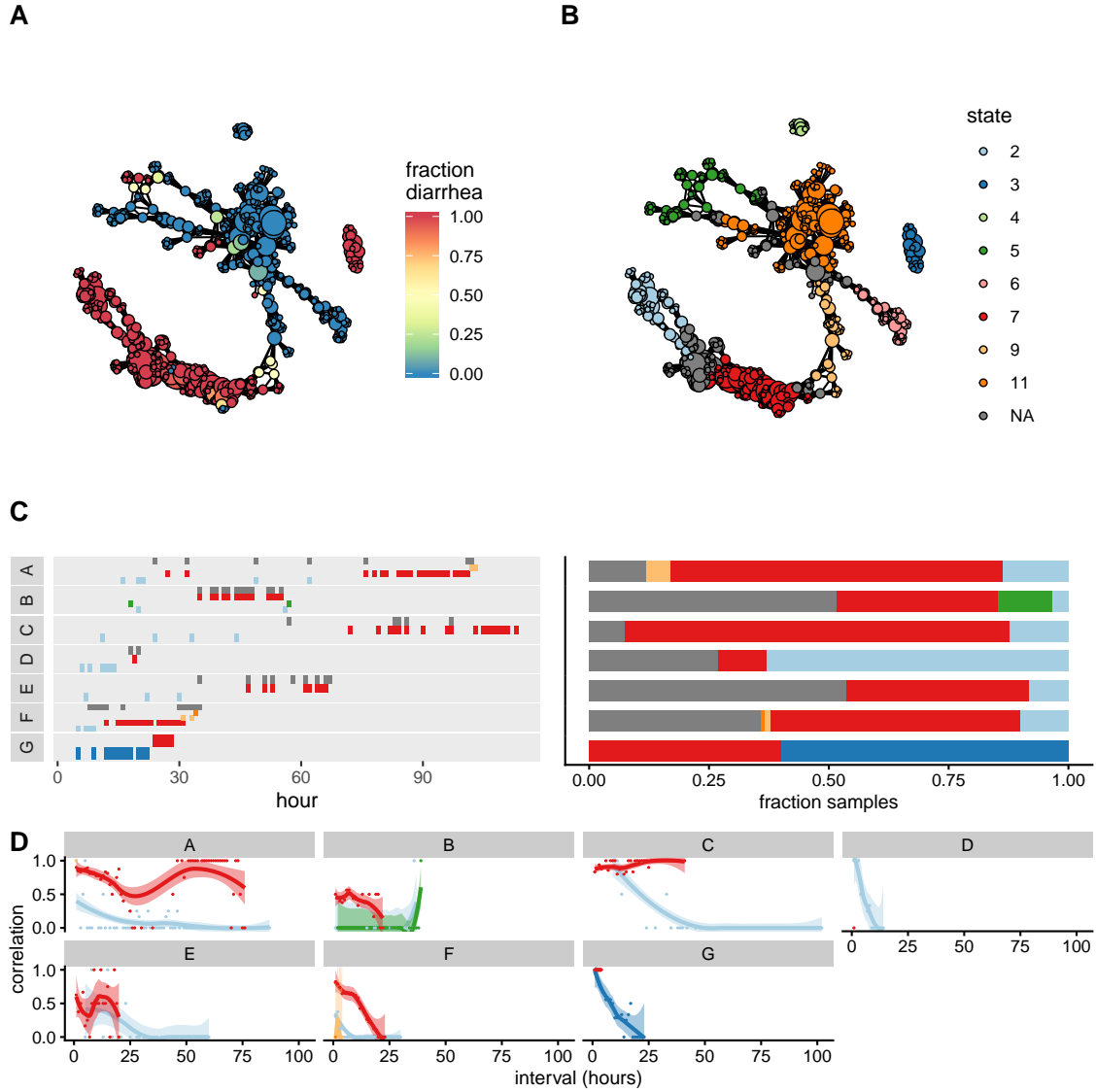


Figure 2: The phase space of the cholera gut microbiome. A. Mapper representation of the combined cholera data reveals disease- and healthy-associated neighborhoods of the phase space. Color: fraction of samples in each vertex associated with diarrhea. Connected components of the Mapper graph representing only one sample are not shown. Disjoint regions of phase space are represented as separate connected components. B. Partitioning of the phase space into metastable states. Vertices unassigned to any state are colored in grey. C. Left: progression of subject compositions during the diarrhea phase by state, showing persistence of states over time. Y axis and color indicate state index, with color indexing as in B. Where a sample was associated with multiple states, all were included. Right: frequency of samples associated with each states during the diarrhea phase for each subject with colors as in B. D. Temporal correlation function for the diarrhea phase of each subject. Dots: raw values of f'_x for pairs of samples (see Methods). Lines: smoothed empirical mean of f'_x . Ribbons: standard error of the mean. Values outside the range of $0 \leq y \leq 1$ omitted.

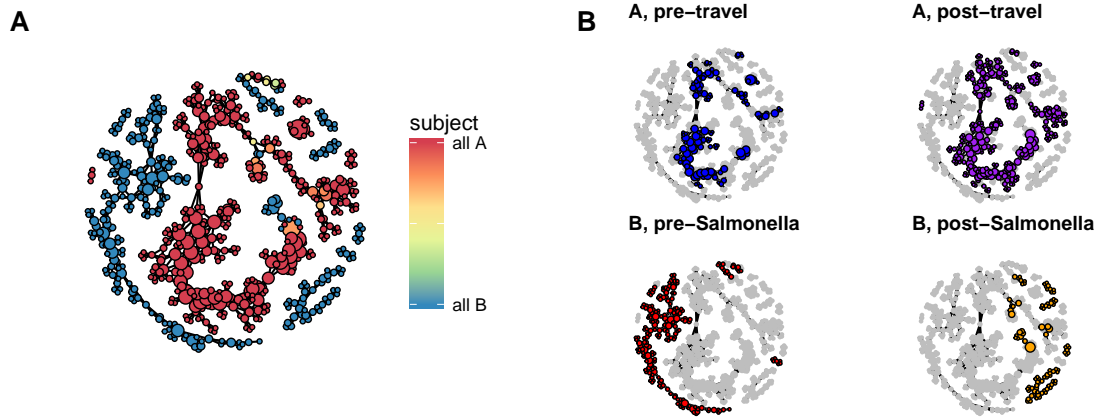


Figure 3: The phase space of two healthy adult male gut microbiomes. A. Mapper representation of the combined daily time series of two healthy adult human gut microbiomes. Connected components of the Mapper graph representing only one sample are not shown. B. Regions of phase space occupied by each subject before after perturbation.

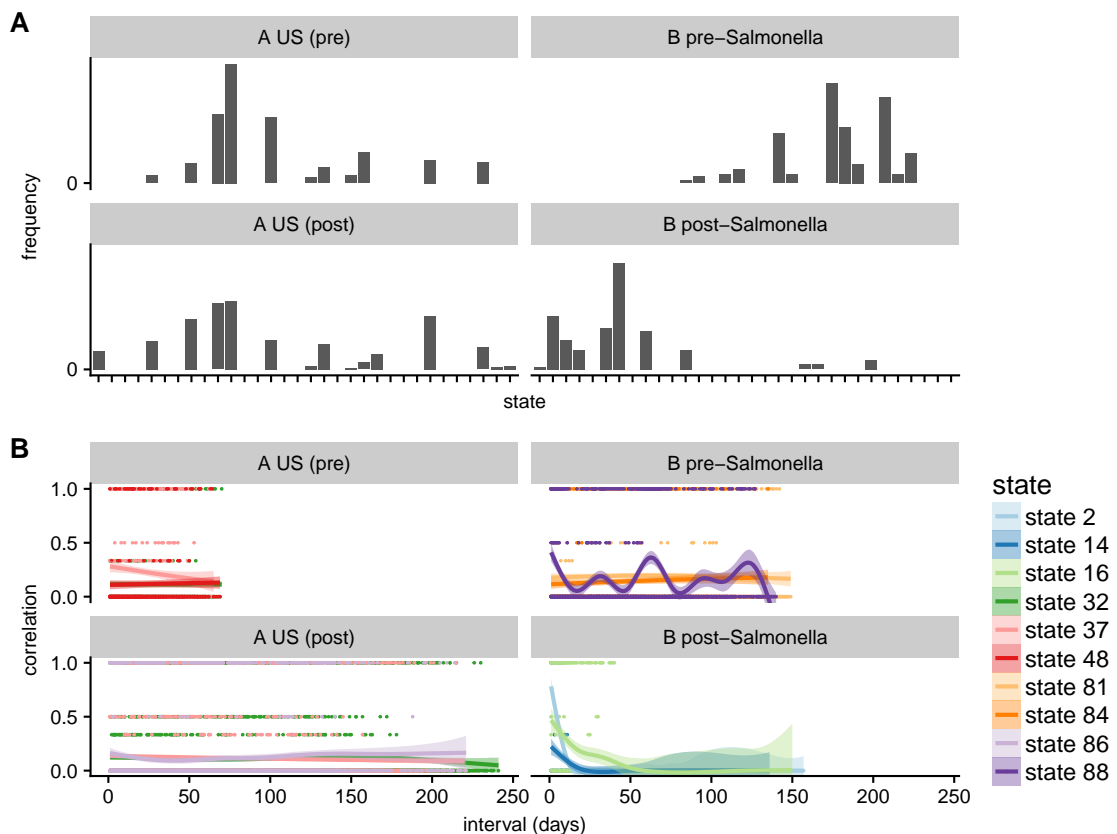


Figure 4: States and dynamics of two healthy adult male gut microbiomes. A. Frequency of states for healthy periods before and after perturbation. X axis: state index. Y axis: frequency of samples. B. Temporal correlation functions for the three most probable states during each event in the 'healthy' phases of each subject. Dots: raw values of f'_x for pairs of samples. Lines: smoothed empirical mean of f'_x . Ribbons: standard error of the mean.

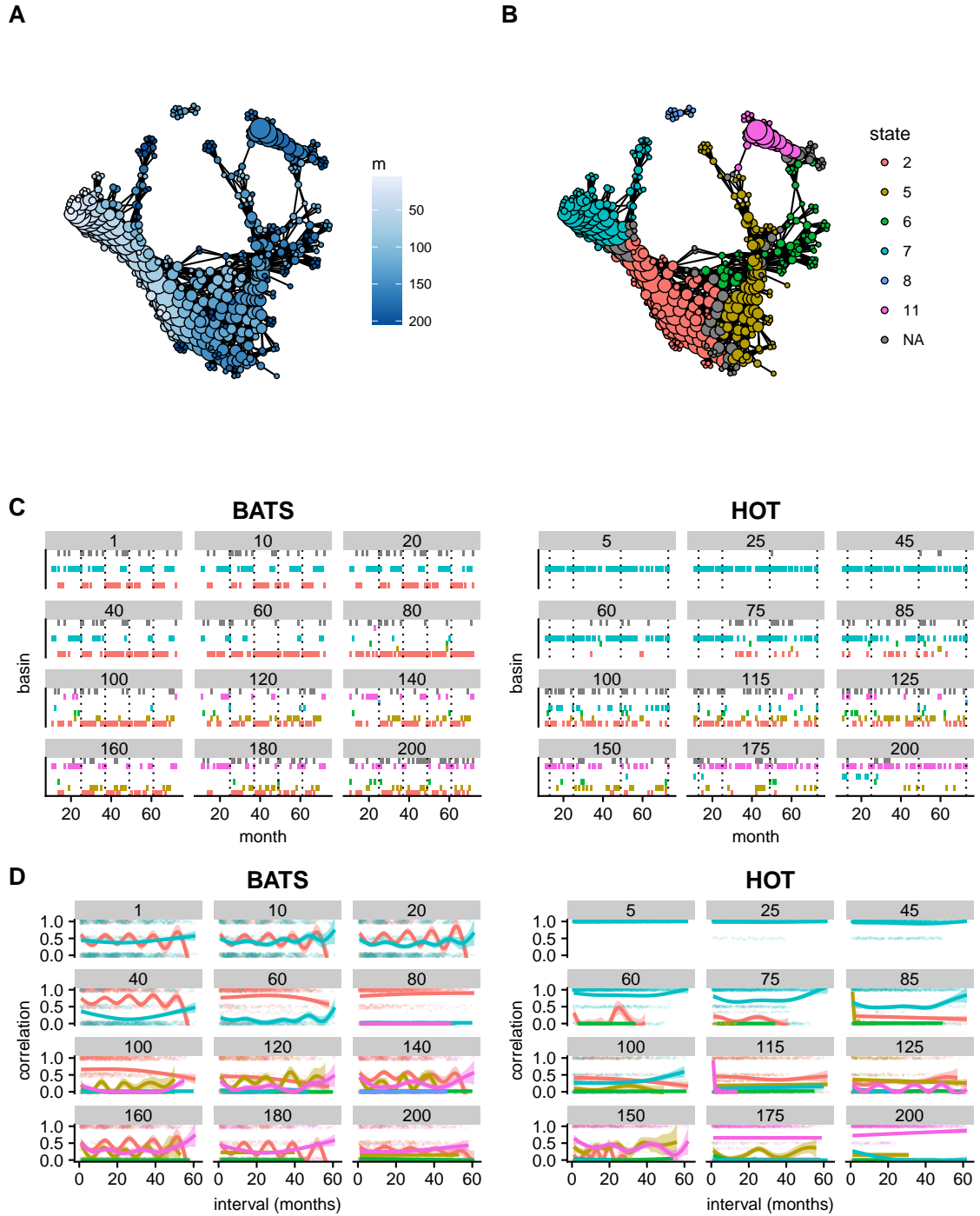


Figure 5: The combined phase space of two *Prochlorococcus* communities inhabiting the Atlantic and Pacific Oceans, respectively. Connected components of the Mapper graph representing only one sample are not shown. A. Vertices colored by mean depth in meters of represented samples. B. Partitioning of the phase space into states. C. Successions of states for each site-depth fraction combination. Dotted lines indicate samples during January. Colors indicate states as in B. D. Temporal correlation functions for each state per site-depth fraction combination. Dots: raw values of f_x for pairs of samples. Lines: smoothed empirical mean of f'_x . Ribbons: standard error of the mean.

641 Tables

Data set	# intervals for (rank(PCo1), rank(PCo2))	% overlap	# bins
Cholera	(15, 15)	70	10
Two healthy adult males	(30, 30)	50	10
<i>Prochlorococcus</i>	(20, 20)	60	10

Table 1: Hyperparameters used to generate the Mapper representation of each data set.

642 Additional Files

643 Supporting information

644

- 645 • Supporting methods describing PCA and hierarchical clustering.
- 646 • Supporting table showing sampling frequency and duration for each of the data sets analyzed.
- 647 • Supporting figure showing the states of the two human gut microbiomes data set.
- 648 • Supporting figure showing the temperature gradients across the *Prochlorococcus* phase space.
- 649 • Supporting figures showing the results of the data rarefaction test.
- 650 • Supporting figures showing the mean physiological or environmental properties per state for
- 651 each data set.
- 652 • Supporting figure showing the results of PCA and hierarchical clustering.

653

654 Supporting data

- 655 • Taxonomy tables showing the mean composition of each state for each data set.