

1 **Controlling complex microbial communities: a network-** 2 **based approach**

3 Marco Tulio Angulo^{1*}, Claude H. Moog² & Yang-Yu Liu^{3,4*}

4 ¹*CONACyT - Institute of Mathematics, Universidad Nacional Autónoma de México 76230, México.*

5 ²*Laboratoire des Sciences du Numérique de Nantes, UMR CNRS 6004, Nantes 44321, France.*

6 ³*Channing Division of Network Medicine, Brigham and Women's Hospital and Harvard Medical School,*
7 *Boston, Massachusetts 02115, USA.*

8 ⁴*Center for Cancer Systems Biology, Dana-Farber Cancer Institute, Boston, Massachusetts 02115, USA.*

9 **Microbes comprise nearly half of all biomass on Earth. Almost every habitat on Earth is**
10 **teeming with microbes, from hydrothermal vents to the human gastrointestinal tract. Those**
11 **microbes form complex communities and play critical roles in maintaining the integrity of**
12 **their environment or the well-being of their hosts. Controlling microbial communities can**
13 **help us restore natural ecosystems and maintain healthy human microbiota. Yet, our abil-**
14 **ity to precisely manipulate microbial communities has been fundamentally impeded by the**
15 **lack of a systematic framework to control them. Here we fill this gap by developing a con-**
16 **trol framework based on the new notion of structural accessibility. This framework allows**
17 **identifying minimal sets of “driver species” through which we can achieve feasible control**
18 **of the entire microbial community. We numerically validate our control framework on large**
19 **microbial communities, and then we demonstrate its application for controlling the gut mi-**
20 **crobiota of gnotobiotic mice infected with *Clostridium difficile* and the core microbiota of the**
21 **sea sponge *Ircinia oros*.**

1 INTRODUCTION

2 Microorganisms form complex communities that play critical roles in maintaining the well-being
3 of their hosts or the integrity of their environment¹⁻⁴. This deep relationship can have severe
4 consequences to the host or the environment when a microbial community is disrupted. In humans,
5 for example, a disruption to the gut microbiota—the aggregate of microorganisms residing in our
6 intestine—has been associated to gastrointestinal diseases such as irritable bowel syndrome, and
7 *Clostridium difficile* Infection (CDI)^{5,6}. A variety of non-gastrointestinal disorders as divergent as
8 autism, obesity, and cavernous cerebral malformations have also been associated with disrupted
9 gut microbiota^{5,7}. For agriculture crops, a disruption to rhizosphere microbiota can reduce their
10 disease resistance and hence affect the overall crop yield^{8,9}. In the oceans, a disruption to their
11 microbiota can impact global climate by altering carbon sequestration rates^{3,4,10}. Driving these
12 microbial communities back to their healthy states has the potential to bring novel solutions to
13 prevent and treat complex human diseases, enhance sustainable agriculture, and regulate global
14 warming^{11,12}. For example, inoculation of soil microbes can restore terrestrial ecosystems¹³, and
15 Fecal Microbiota Transplantation (FMT) is so far the most successful therapy in treating patients
16 with recurrent CDI by restoring disrupted gut microbiota¹⁴. Despite the success of these two
17 empirical strategies, a broad application of microbial-manipulation strategies will only be possible
18 if we can efficiently and systematically control large complex microbial communities¹⁵.

19 There are two main challenges to efficiently control a large complex microbial community.
20 First and foremost, an efficient control method should only manipulate the minimal necessary

1 number of species in the community. However, we still lack a method to systematically identify
2 minimal sets of those “driver species” whose control can help us drive the whole community to
3 desired states. Here, we use the term “species” in the general context of ecology, i.e., as a set
4 of organisms adapted to a particular set of resources in the environment. It doesn’t necessarily
5 represent the lowest major taxonomic rank. In fact, one could think of organizing microbes by
6 strains, genera, or operational taxonomical units as well. Second, even if those driver species were
7 known, calculating the control strategy that should be applied to them for driving the community
8 towards the desired state remains somewhat tricky (e.g., it is difficult to calculate how much the
9 abundance of those drive species needs to be increased or decreased). The difficulty in solving
10 this second challenge is not only due to our insufficient knowledge of microbial dynamics and
11 interactions, but also because of the inherently complex dynamics they often display.

12 To efficiently and systematically control large complex communities, here we develop a
13 framework showing that the above two challenges can be addressed by focusing on the ecolog-
14 ical network underlying the microbial community. We first introduce the new notion of “structural
15 accessibility” and derive its graph-theoretical characterization. This theoretical result enables us
16 to efficiently identify minimal sets of driver species of any microbial community purely from the
17 topology of its underlying ecological network, even if some microbial interactions are missing
18 and its population dynamics is unknown. Structural accessibility is a generalization of the notion
19 of structural controllability¹⁶ —which only applies to systems with linear dynamics— to systems
20 with nonlinear dynamics. Linear structural controllability is receiving increasing attention from the
21 viewpoint of Network Science¹⁷. Once the driver species are identified, we systematically design

1 feedback control strategies to drive a microbial community towards the desired state, even if the
2 microbial dynamics is not precisely known. We numerically validated our control framework in
3 large microbial communities, analyzing its performance for different parameters of the community
4 we aim to control (e.g., the connectivity of its underlying ecological network), and with respect to
5 errors in the ecological network used to identify the driver species. Finally, we demonstrate our
6 framework by controlling the core microbiota of the sea sponge *Ircinia oros*, and restoring the gut
7 microbiota of gnotobiotic mice infected by *Clostridium difficile*. Our results provide a rational and
8 systematic framework to control microbial communities and other complex ecosystems based only
9 on knowing their underlying ecological networks.

10 **PROBLEM STATEMENT**

11 In our modeling framework, we focus on exploring the impact that manipulating a subset of species
12 has on the abundances of other species. We thus consider a microbial community whose *state* at
13 time t can be determined from the abundance profile $x(t) \in \mathbb{R}^N$ of its N species, where the i -th
14 entry $x_i(t)$ of $x(t)$ represents the abundance of the i -th species at time t . Let us assume that the
15 state evolves according to some general population dynamics

$$\dot{x}(t) = f(x(t)), \quad (1)$$

16 where the function $f : \mathbb{R}^N \rightarrow \mathbb{R}^N$ models the intrinsic growth and inter/intra-species interactions
17 of the community (see Supplementary Note 1 for details). For most microbial communities the
18 function f is unknown and difficult to infer due to the manifold of interaction mechanisms between
19 microbes, such as cross-feeding and modulation by the host immune system¹⁸. Thus we assume

1 that $f(x)$ is some unknown *meromorphic* function (i.e., each entry $f_i(x)$ is the quotient of analytic
2 functions of x). This is a very mild assumption that is satisfied by most population dynamics
3 models¹⁹.

4 Instead of knowing the population dynamics of the microbial community, we assume we
5 know its underlying *ecological network* $\mathcal{G} = (X, E)$. This network is defined as a directed graph
6 where nodes $X = \{x_1, \dots, x_N\}$ represent species and edges $(x_j \rightarrow x_i) \in E$ denote that the j -
7 th species has a direct ecological impact (e.g., direct promotion or inhibition) on the i -th species
8 (Fig.1a). Mapping these ecological networks requires performing mono-culture and co-culture
9 experiments^{20,21}, using time-resolved abundance data and system identification techniques^{22,23}, or
10 using steady-state abundance data via a recently developed inference method²⁴. The accuracy of all
11 these methods strongly depends on how informative is the available data²⁵. Note that these ecolog-
12 ical networks are different from correlation or co-occurrence based networks because correlation
13 doesn't imply causation²⁶. Correlation-based networks can be readily constructed from abundance
14 profiles of different samples^{20,27} and, under certain specific conditions²⁸, they could be a proxy of
15 the underlying ecological network.

16 Controlling a microbial community consists in driving its state from an initial value $x_0 =$
17 $x(0) \in \mathbb{R}^N$ at time $t = 0$ (e.g., a “diseased” state) towards a desired value $x_d \in \mathbb{R}^N$ (e.g., a
18 “healthier” state, Fig. 1b). We consider that the community will not naturally evolve to the desired
19 state. To drive the microbial community, we consider a set of M *control inputs* $u(t) \in \mathbb{R}^M$
20 that directly affect certain species that we call *actuated species* (Fig. 1a). These control inputs

1 encode a combination of M control actions that are simultaneously applied to the community at
2 time t . There are four types of control actions that we consider. If $u_j(t) < 0$, the j -th control
3 action at time t is either a *bacteriostatic agent* or a *bactericide*, which decreases the abundance
4 of the species it actuates by inhibiting their reproduction or directly killing them, respectively²⁹.
5 If $u_j(t) > 0$, the j -th control action at time t is either a *prebiotic*³⁰ or a *transplantation*, which
6 stimulate the growth or engrafts a consortium of the species it actuates, respectively. For the
7 human gut microbiota, probiotics administration³¹ and FMTs¹⁴ are examples of transplantations.
8 We introduce the *controlled ecological network* of the community $\mathcal{G}^c = (X \cup U, E \cup B)$ to specify
9 which species are actuated by each control input. Here, the set $U = \{u_1, \dots, u_M\}$ is the set of
10 control input nodes, and the edge $(u_j \rightarrow x_i) \in B$ denotes that the the j -th control input actuates
11 the i -th species (Fig. 1a).

12 Given a controlled ecological network describing the interactions between species and which
13 species are actuated by the control inputs, we next introduce two control schemes describing how
14 the control inputs will affect the species. The first control scheme models a combination of prebi-
15 otics (if $u_j(t) > 0$) and bacteriostatic agents (if $u_j(t) < 0$) as *continuous* control inputs modifying
16 the growth of the actuated species (Fig. 1c):

$$\dot{x}(t) = f(x(t)) + g(x(t))u(t), \quad t \in \mathbb{R}. \quad (2)$$

17 The second control scheme considers a combination of transplantations (if $u_j(t) > 0$) and bacteri-
18 cides (if $u_j(t) < 0$) applied at discrete *intervention instants* $\mathbb{T} = \{t_1, t_2, \dots\}$, rendering *impulsive*

1 control inputs that instantaneously modify the abundance of the actuated species (Fig. 1d):

$$\dot{x}(t) = f(x(t)) \text{ if } t \notin \mathbb{T}, \quad x(t^+) = x(t) + g(x(t))u(t) \text{ if } t \in \mathbb{T}. \quad (3)$$

2 In the above equation, the symbol $x(t^+)$ denotes the state “right after time t ”, so a control input
3 $u(t) \neq 0$ at $t \in \mathbb{T}$ makes $x(t)$ “jump” at that time instant. Thus, control actions are classified as
4 impulsive if they instantaneously modify the abundance of some species, and continuous otherwise
5 (see Supplementary Note 1.2 for details).

6 Both control schemes are characterized by the pair of functions $\{f, g\}$, describing the con-
7 trolled population dynamics of the microbial community. As we have seen, the function $f : \mathbb{R}^N \rightarrow$
8 \mathbb{R}^N models the intrinsic growth and inter/intra-species interactions. The function $g : \mathbb{R}^N \rightarrow \mathbb{R}^{N \times M}$
9 models the direct susceptibility of the species to the control actions. The i -th species is actuated
10 by the j -th control input if the (i, j) -th entry of $g(x)$ satisfies $g_{ij}(x) \neq 0$. As in the uncon-
11 trolled community of Eq. (1), the function $g(x)$ is typically unknown because the mechanisms of
12 susceptibility to the control actions can be uncertain. Thus we assume that g is some unknown
13 meromorphic function such that $g_{ij} \neq 0$ iff $(u_j \rightarrow x_i) \in \mathcal{B}$.

14 Notice that when all species are directly controlled (i.e., each species is actuated by an inde-
15 pendent control input so $M = N$ and $g(x)$ is full rank), the state of the whole microbial commu-
16 nity can obviously be fully controlled. Fortunately, as we next show, controlling all the species in
17 a community is far from being necessary. Indeed, several species can be *indirectly* controlled by
18 the same control input when this signal is adequately propagated through the ecological network
19 underlying the community. Thus, our first goal is to identify minimal sets of species that we need

1 to actuate in order to drive the entire community. We call those species *driver species*. We will
2 also study if the impulsive control scheme can be as effective as the continuous control scheme for
3 controlling microbial communities. Indeed, the former is more feasible than the latter, especially
4 for human-associated microbial communities. Finally, we will design the control inputs that should
5 be applied to the identified driver species to drive the whole community towards the desired state.

6 IDENTIFYING DRIVER SPECIES

7 **Driver species are characterized by the absence of autonomous elements.** To understand when
8 a set of actuated species is a set of driver species, consider a three-species community with the clas-
9 sical Generalized Lotka-Volterra (GLV) population dynamics (Fig. 2a). This toy community has
10 one control input actuating the third species x_3 . Actuating this species alone creates an *autonomous*
11 *element* —namely, a constraint between some species abundances that the control input cannot
12 break, confining the state of the community to a low-dimensional manifold. More precisely, our
13 mathematical formalism reveals that $\xi = x_1x_2$ is an autonomous element for this microbial com-
14 munity (Example 2 in Supplementary Note 2). Indeed, differentiating ξ with respect to time yields
15 $\dot{\xi} = x_1x_2(1 - x_3) + x_1x_2(-1 + x_3) \equiv 0$, which implies that the state of the community is con-
16 strained to the low-dimensional manifold $\{x \in \mathbb{R}^3 | x_1x_2 = x_1(0)x_2(0)\}$ for all control inputs (Fig.
17 2a right). Intuitively, an autonomous element exists because the control input cannot change the
18 abundance of species x_1 without changing the abundance of species x_2 in a predefined way (i.e.,
19 $x_2 = x_1(0)x_2(0)/x_1$). It is thus impossible to drive the whole community in its three-dimensional
20 state space, implying that x_3 cannot be a driver species for this community. Introducing a second

1 control input actuating species x_1 eliminates this autonomous element by helping the system to
2 jump out of the low-dimensional manifold. Hence, the community can be driven in any direction
3 within its three-dimensional state space (Fig. 2b). This indicates that $\{x_1, x_3\}$ is a minimal set of
4 driver species for this community. Actually, by using these two driver species we can steer the
5 community to any desired state with positive abundances (Example 6 in Supplementary Note 5).

6 In the general case of N species and M control inputs, we define a set of actuated species
7 as a set of driver species if the corresponding controlled population dynamics $\{f, g\}$ of the micro-
8 bial community lacks autonomous elements. For a given pair $\{f, g\}$, the absence of autonomous
9 elements can be mathematically deduced using a formalism based on differential one-forms (Sup-
10 plementary Note 2). Indeed, for the continuous control scheme of Eq. (2), the conditions for the
11 absence of autonomous elements are well understood because they define when a system is *accessi-*
12 *ble*³². As a cornerstone concept in nonlinear control theory, accessibility has been instrumental for
13 developing technological advances such as robotics. Since it is more natural to control microbial
14 communities with impulsive control actions, in this paper we extended the study of autonomous
15 elements to the impulsive control systems of Eq. (3). For this, we first introduced a mathematical
16 definition of autonomous elements for impulsive control systems (Definition 3 in Supplementary
17 Note 2). Using this definition, we characterized necessary and sufficient conditions for the absence
18 of autonomous elements in a given controlled population dynamics (Theorem 2 in Supplementary
19 Note 2).

20 To our surprise, we found that the conditions for the absence of autonomous elements for the

1 continuous and the impulsive control schemes are identical (Remark 2 in Supplementary Note 2).
2 This result suggests that, for controlling microbial communities, transplantsations and bactericides
3 (impulsive control actions) can be as effective as prebiotics and bacteriostatic agents (continuous
4 control actions). Since impulsive control actions could be simpler to implement for many mi-
5 crobial communities such as the human gut microbiota, this result assures us to further develop
6 microbiome-based therapies in the form of probiotic cocktails and FMTs.

7 **Structural accessibility characterizes the generic absence of autonomous elements.** For com-
8 plex microbial communities such as the human gut microbiota, it is very difficult to choose an
9 adequate pair $\{f, g\}$ to model its controlled population dynamics. As the autonomous elements
10 depend on such a pair, this might suggest that it is impossible to predict their presence and thus to
11 identify the driver species of complex microbial communities. We now show that this seemingly
12 unavoidable limitation can be solved by focusing on the topology of the controlled ecological
13 network of the community.

14 Define the graph $\mathcal{G}_{f,g} = (X \cup U, E_{f,g} \cup B_{f,g})$ associated with a meromorphic function pair
15 $\{f, g\}$ as follows. First, the edge $(x_j \rightarrow x_i) \in E_{f,g}$ exists if x_j appears in the right-hand side of \dot{x}_i
16 or $x_i(t^+)$ in Eqs. (2) or (3), respectively. Second, the edge $(u_j \rightarrow x_i) \in B_{f,g}$ exists if $g_{ij} \neq 0$. In
17 this definition, the interaction $x_j \rightarrow x_i$ can originate in the uncontrolled population dynamics (i.e.,
18 $f_i(x)$ depends on x_j) or, in a more general case, also in the controlled dynamics (i.e., the i -th row
19 of $g(x)$ depends on x_j). Using this definition and given a controlled ecological network \mathcal{G}^c , we can
20 describe the class \mathcal{D} of all possible controlled population dynamics that the controlled microbial

1 community can have. Mathematically, we describe the class \mathfrak{D} as containing all *base models*
 2 $\{f^*, g^*\}$ such that $\mathcal{G}_{f^*, g^*} = \mathcal{G}^c$, together with all *deformations* $\{f, g\}$ of each of those base models.
 3 The base models characterize the simplest controlled population dynamics that the community can
 4 have. We have chosen them as controlled GLV models with constant susceptibilities:

$$f_i^*(x) = r_i x_i + \sum_{j=1}^N a_{ij} x_i x_j, \quad g_{ij}^*(x) = b_{ij}, \quad (4)$$

5 for $i = 1, \dots, N$. The base models are parametrized by $A = (a_{ij}) \in \mathbb{R}^{N \times N}$, $r = (r_i) \in \mathbb{R}^N$,
 6 and $B = (b_{ij}) \in \mathbb{R}^{N \times M}$, representing the interaction matrix, the intrinsic growth rate vector,
 7 and the susceptibility matrix of the community, respectively. Thus, the base models in \mathfrak{D} are all
 8 controlled GLV models such that their graph matches \mathcal{G}^c . As a classical population dynamics
 9 model, the GLV model has been applied to microbial communities in lakes, soils, and human
 10 bodies^{14, 15, 20, 33–39}. Notice that in a microbial community, any species that gets extinct cannot
 11 “resurrect” by itself without some external influence such as a transplantation or migration. Eq.
 12 (4) is the simplest population dynamics that satisfies this condition in the following sense: it is
 13 obtained by considering population dynamics of the form $f_i(x) = x_i F_i(x)$, and then choosing the
 14 functions $F_i(x)$ to be simple affine functions.

15 Next, we say that a meromorphic pair $\{f, g\}$ is a deformation of a base model $\{f^*, g^*\}$
 16 if it satisfies the following three conditions: (i) it has the same graph as the base model (i.e.,
 17 $\mathcal{G}_{f, g} = \mathcal{G}_{f^*, g^*}$); (ii) there exists a finite set of parameters $\theta \in \mathbb{R}^C$ such that $\{f(x), g(x)\} =$
 18 $\{\tilde{f}(x; \theta), \tilde{g}(x; \theta)\}$; and (iii) the identity $\{\tilde{f}(x; 0), \tilde{g}(x; 0)\} = \{f^*(x), g^*(x)\}$ holds. The minimal
 19 integer $C \geq 0$ for which these conditions are satisfied is called the *size* of the deformation, quanti-
 20 fying the cardinality of the parameter set θ that is needed to obtain the deformation from the base

1 model. A rather general class of controlled population dynamics can be described by deformations
 2 of the base model of Eq. (4), such as

$$f_i(x; \theta) = \theta_{i,1} + x_i (-r_i - \theta_{i,2}x_i) (\theta_{i,3}x_i - 1) + \sum_{j=1}^N a_{ij} \frac{x_i x_j}{1 + \theta_{ij,4} + \theta_{ij,5}x_i + \theta_{ij,6}x_i x_j + \theta_{ij,7}x_j}, \quad (5)$$

3 for $i = 1, \dots, N$. In Eq. (5), the parameters $\theta_{i,1}$ are migration rates from/to neighboring habi-
 4 tats, $\theta_{i,2}^{-1}$ are the carrying capacities of the environment, $\theta_{i,3}^{-1}$ are the Allee constants, and the rest
 5 $\{\theta_{ij,k}\}_{k=4}^7$ characterize the saturation of the functional responses⁴⁰. Note that $\theta_{i,1} > 0$ can also
 6 model species like *C. difficile* that sporulate into “inactive” forms and then recover. Note also
 7 that “higher-order” interactions can be described as deformations. For example, if species x_i is
 8 directly affected by species x_j and x_k , then a deformation can include the third-order interaction
 9 $\theta_{ij}x_i x_j x_k$. Similarly, deformations allow cases when the susceptibility of the i -th species to j -
 10 th control input is mediated by the abundance of other species. For example, the deformation
 11 $g_{ij}(x; \theta) = b_{ij} + \theta_{ijk}x_k$ models a case when the i -th species is actuated by the j -th control input
 12 but its effect is mediated by the abundance of the k -th species.

13 We call the class \mathfrak{D} *structurally accessible* if almost all of its base models and almost all
 14 of their deformations lack autonomous elements. This means that, except for a zero-measure
 15 set of “singularities”, all the controlled population dynamics that the community may take have
 16 to lack autonomous elements. The conditions under which \mathfrak{D} is structurally accessible are fully
 17 characterized using our mathematical formalism (Supplementary Note 3), and they depend only
 18 on the underlying controlled ecological network \mathcal{G}^c . We first proved that, generically, increasing
 19 the size of a deformation cannot create autonomous elements (see Proposition 1 in Supplementary
 20 Note 3, and Fig. 2c for an illustration). This result reduces the search for autonomous elements to

1 the deformations in \mathcal{D} with minimal size $C = 0$. That is, to all base models whose graph matches
2 \mathcal{G}^c . Finally, we proved that \mathcal{D} is structurally accessible if and only if \mathcal{G}^c satisfies the following two
3 conditions: (i) each species is the end-node of a path that starts at a control input node; and (ii)
4 there is a disjoint union of cycles (excluding self-loops) and paths that cover all species nodes (see
5 Theorem 3 of Supplementary Note 3). If these two graph conditions are satisfied, we also call \mathcal{G}^c
6 structurally accessible.

7 The notion of structural accessibility introduced above is a nonlinear counterpart of the no-
8 tion of structural controllability for linear systems¹⁶. For linear systems we have $\{f(x), g(x)\} =$
9 $\{Ax, B\}$, and the absence of autonomous elements is equivalent to their controllability³² —the
10 intrinsic ability to drive the system between two arbitrary states, which can be verified by the cel-
11 ebrated Kalman’s rank condition: $\text{rank}(B, AB, A^2B, \dots, A^{N-1}B) = N$. Condition (i) above is
12 necessary for both structural accessibility and linear structural controllability, requiring that the
13 network contains paths that spread the influence of the control inputs to all species. However, for
14 linear structural controllability, condition (ii) is sufficient but not necessary. More precisely, for
15 linear structural controllability, the required disjoint union of cycles that cover the species nodes
16 can also include self-loops due to intrinsic nodal dynamics (see Remark 4 in Supplementary Note
17 3).

18 **Identifying minimal sets of driver species in microbial communities.** The above result provides
19 a complete graph-characterization of driver species: a set of actuated species is a set of driver
20 species (for all but a zero-measure set of controlled population dynamics that the community may

1 have) if and only if its corresponding \mathcal{G}^c is structurally accessible. We used this characterization
2 to build an algorithm that identifies a minimal set of driver species from the ecological network
3 of the community. More precisely, we mapped the satisfaction of the graph conditions (i) and (ii)
4 into solving a maximum matching problem over the graph \mathcal{G} without self-loops (Proposition 3 in
5 Supplementary Note 4). This result provides a polynomial time algorithm to identify one minimal
6 set of driver species, making it feasible for large networks (Remark 5 in Supplementary Note 4).

7 Note that once \mathcal{G}^c is structurally accessible this network cannot lose its structural accessibil-
8 ity when new edges are added to it. This observation implies that a set of driver species remains
9 valid even if new edges (e.g., new inter/intra-species interactions) are added to the ecological net-
10 work of the community. Therefore, it is possible to find the driver species of a microbial community
11 using an “incomplete” ecological network that only includes some of the ecological interactions
12 (e.g., high-confidence interactions).

13 **DRIVING THE DRIVER SPECIES**

14 Next we turn to the question of calculating the control signal $u(t)$ that needs to be applied to a
15 set of driver species to drive the whole community towards the desired state. We will show that
16 impulsive control actions can make this calculation easier.

17 **Calculating optimal control strategies for microbial communities with known population dy-**
18 **namics.** To calculate the impulsive control inputs $\{u(t_k), t_k \in \mathbb{T}\}$ needed to drive the micro-
19 bial community to the desired state x_d we adopt a *model predictive control* (MPC) approach⁴¹.

1 First, based on the current state of the community $x(t_k)$ at the intervention instant $t_k \in \mathbb{T}$, we
 2 use knowledge of its controlled population dynamics to predict the sequence of states $\hat{X}_{k,L} =$
 3 $\{\hat{x}(t_{k+1}), \dots, \hat{x}(t_{k+L+1})\}$ that the community will take in response to a sequence of L impulsive
 4 control inputs $U_{k,L} = \{u(t_k), \dots, u(t_{k+L-1})\}$. The *prediction horizon* $L > 0$ quantifies how far
 5 into the future we predict. We then choose $u(t_k) = u_1^*(t_k)$, where $u_1^*(t_k)$ is the first element of the
 6 optimal control input sequence $U_{k,L}^*$ calculated by solving the following optimization problem:

$$U_{k,L}^* = \arg \min_{U_{k,L} \in \mathbb{R}^{M \times L}} J_{x_d}(\hat{X}_{k,L}, U_{k,L}) \text{ subject to } U_{k,L} \in \Omega. \quad (6)$$

7 Here $\Omega \subseteq \mathbb{R}^{M \times L}$ is a set that specifies constraints in the control inputs we can use, and J_{x_d} is some
 8 cost function penalizing deviations of the predicted trajectory $\hat{X}_{k,L}$ from the desired state x_d . For
 9 example, the simplest cost function $J_{x_d}(\hat{X}_{k,L}, U_{k,L}) = \|\hat{x}(t_{k+L+1}) - x_d\|$ penalizes deviations of
 10 the predicted final state from the desired state. Penalizing the deviations of intermediate states can
 11 provide a smoother transition to the desired state.

12 To choose the prediction horizon L in Eq. (6), we proved that it is possible to distinguish
 13 between two cases (Theorem 4 in Supplementary Note 5). The first case is when the commu-
 14 nity can be driven to the desired state using a finite number L of impulsive control actions. This
 15 number can be calculated from its controlled population dynamics. The second case is when the
 16 community can only be asymptotically driven to the desired state as time goes by, meaning that a
 17 “sufficiently large” $L \gg N$ should be used. This second case could be circumvented by increasing
 18 the number of driver species (Remark 8 in Supplementary Note 5). Note that by recalculating $U_{k,L}^*$
 19 at each intervention instant $t_k \in \mathbb{T}$ using the actual state of the community, the MPC method cre-
 20 ates a feedback loop that enhances its robustness against prediction errors due to uncertainty in the

1 dynamics⁴¹. For $L = 1$ the proposed MPC methodology is similar to the network control method
2 of Ref. (40). Eq. (6) is a finite-dimensional optimization problem that can be solved using several
3 algorithms such as “DIRECT”⁴². By contrast, for continuous control actions, the analogous opti-
4 mization problem is defined over the infinite-dimensional space of all M -dimensional continuous
5 functions. Solving such optimization problem is apparently more difficult, significantly limiting
6 our ability to calculate optimal continuous control actions.

7 We studied the performance of the above MPC strategy in the three-species microbial com-
8 munity with a solo driver species of Fig. 1. Given the dynamics of this community (see caption
9 in Fig. 1), we find that $L = 3$ impulsive control inputs are sufficient to drive the whole commu-
10 nity (Example 4 in Supplementary Note 5). To calculate the optimal control inputs we selected
11 $J_{x_d}(\hat{X}_{k,L}, U_{k,L}) = \|\hat{x}(t_{k,L}) - x_d\|_2$ in Eq. (6). Solving the optimization problem using DIRECT
12 yields the MPC strategy $u^*(t_1) = -0.8815$, $u^*(t_2) = 2.0089$ and $u^*(t_3) = -10^{-4}$ (pink in Fig.
13 3a). We use this example to compare the performance of applying two other control strategies
14 to drive this community. The first strategy uses a transplantation to restore the abundance of the
15 driver species (i.e., increase its abundance to its desired value), expecting that such control action
16 will drive the rest of the community to the desired state (purple in Fig. 3a). This control strat-
17 egy is reminiscent of a probiotic administration that restores the “healthy” abundance of the driver
18 species. The second control strategy ignores the driver species of this community, using two con-
19 trol inputs (instead of one) to set the abundance of the non-driver species to their desired values
20 (blue in Fig. 3a).

1 Among the above control strategies, only the MPC applied to the driver species succeeds
2 (Fig. 3b). Actually, this strategy succeeds in a somewhat unconventional way: despite the driver
3 species is more abundant in the desired state than in the initial state, the first control action de-
4 creases further its abundance. This first control action makes the non-driver species reach their
5 desired abundances and, once that happens, the abundance of the driver species is finally increased
6 to its desired value (pink in Fig. 3b). The second control strategy succeeds in driving species x_2
7 and x_3 , but it fails to drive x_1 to the desired abundance because it approaches the desired state from
8 an unstable direction (purple in Fig. 3b). Finally, not actuating the driver species results in the
9 worst strategy, failing to drive a single species to the desired state (blue in Fig. 3b). This example
10 demonstrates the importance of actuating the driver species.

11 **Calculating control strategies for microbial communities with uncertain population dynam-**
12 **ics or a large number of species.** In general, solving the non-convex optimization problem of
13 Eq. (6) is challenging as the number of species or prediction horizon increase. Also, a prerequisite
14 for solving this optimization problem is a reasonable knowledge of the controlled population dy-
15 namics of the community, which may not available. To circumvent these two drawbacks, next we
16 leverage the network underlying the controlled microbial community.

17 Consider that it is possible to obtain a weighted adjacency matrix $\hat{A} \in \mathbb{R}^{N \times N}$ from the
18 ecological network \mathcal{G} of the community, providing a proxy for its interaction matrix. Without
19 additional knowledge of the susceptibility matrix of the community, we assume it is possible to
20 increase or decrease as desired the abundance of each driver species. Under this assumption, we

1 define $\hat{B} \in \{0, 1\}^{N \times M}$ as a proxy for the susceptibility matrix, with $b_{ij} = 1$ if the j -th control
 2 input actuates the i -th driver species. Next, by rewriting the controlled population dynamics of the
 3 community as $\{f(x), g(x)\} = \{\hat{A}x + w_x, \hat{B} + w_u\}$, we use the pair $\{\hat{A}x, \hat{B}\}$ to provide a linear
 4 prediction for the response of the community to the control inputs. Here, the nonlinear functions
 5 $(w_x, w_u) = (f - \hat{A}x, g - \hat{B})$ represent perturbations whose magnitude depend on how well the
 6 linear pair $\{\hat{A}x, \hat{B}\}$ approximates the true dynamics $\{f(x), g(x)\}$ of the community. Using this
 7 linear pair for predicting the response of the community to impulsive control actions, we design a
 8 *linear MPC* by solving the optimization problem of Eq. (6) with the quadratic cost function

$$J_{x_d}(\hat{X}_{k,L}, U_{k,L}) = \sum_{i=k}^L [\hat{x}(t_i) - x_d]^\top Q [\hat{x}(t_i) - x_d] + u(t_i)^\top R u(t_i).$$

9 In the above equation, the positive definite matrices $Q = Q^\top \in \mathbb{R}^{N \times N}$ and $R = R^\top \in \mathbb{R}^{M \times M}$
 10 are design parameters. The matrix Q penalizes the deviations of the predicted trajectory from the
 11 desired state, and R quantifies the “cost” of using the control inputs. Under this scenario (i.e., a
 12 linear prediction model and quadratic cost), the solution to the optimization problem of Eq. (6)
 13 can be obtained in closed form⁴³ even if $L \rightarrow \infty$. This result enabled us to obtain the explicit form
 14 $u(t_k) = Kx(t_k)$ for the linear MPC at time $t_k \in \mathbb{T}$, where $K \in \mathbb{R}^{M \times N}$ is computed by solving a
 15 Riccati algebraic equation (Supplementary Note 6). Since the Riccati equation can be efficiently
 16 solved for large N , the linear MPC can be calculated for large microbial communities. The above
 17 linear MPC has several other advantages: it requires minimal knowledge of the controlled popula-
 18 tion dynamics of the community (i.e., the weighted adjacency matrix of its underlying ecological
 19 network); it is robust to the perturbations (w_x, w_u) and other uncertainties (Remark 12 in Sup-
 20 plementary Note 6); and it also allows calculating the control signals for the continuous control

1 scheme (Remark 10 in Supplementary Note 6).

2 We used the above linear MPC for controlling the three-species community of Fig. 1, assum-
3 ing its dynamics is unknown. Based on the ecological network of this community and its popula-
4 tion dynamics (see Fig. 1 and its caption), we choose $\hat{A} = (-0.5, 0, -0.1; 0, -5, 1; 0, 0, -1)$ as a
5 proxy for its interaction matrix. Note that \hat{A} is a rather rough approximation of the linearization of
6 the population dynamics at the desired state given by $(-0.37, 0, -0.05; 0, -5.31, 0.52; 0, 0, -1)$.
7 Since $\{x_3\}$ is a solo driver species for this community, we use $\hat{B} = (0; 0; 1)$. Choosing $Q =$
8 $\text{diag}(20, 1, 10)$, we compared the performance of three different linear MPCs obtained by using
9 the values $R = 10^{-4}, 10^{-3}, 10^{-2}$ (Fig. 3c). The performance of the linear MPC strongly depends
10 on the selection of these parameters. For $R = 10^{-4}$, despite not using knowledge of the population
11 dynamics, the performance of the linear MPC (pink in Fig. 3d) is very similar to the performance
12 of the MPC that uses full knowledge of the nonlinear population dynamics (pink in Fig. 3b). The
13 success of the linear MPC in driving a community with nonlinear population dynamics illustrates
14 the robustness of the MPC strategy, since the controller succeeds despite having non-zero pertur-
15 bations (w_x, w_u) . As R increases, the performance of the linear MPC deteriorates, first using more
16 interventions to reach the desired state (green in Fig. 3d), and finally failing to drive the system to
17 the desired state (blue in Fig. 3d). Indeed, since $R > 0$ quantifies the “cost” of using control inputs,
18 increasing R reduces the magnitude of the control inputs, to the point they are not large enough to
19 drive the system towards the desired state. We emphasize that, in general, the performance of the
20 linear MPC also depends on the chosen (\hat{A}, \hat{B}) and the desired state (Remark 11 in Supplementary
21 Note 6).

1 **Numerical validation of the control framework on large microbial communities.** To system-
2 atically validate our control framework, we considered communities of $N = 100$ species having
3 random Erdős-Rényi ecological networks with a prescribed connectivity $c \in [0, 1]$, see Fig. 4a.
4 The network edge-weights are chosen from a normal distribution with zero mean and standard
5 deviation $\sigma \geq 0$, where σ characterizes the typical interspecies interaction strength. Negative self-
6 loops with weights -1 were added to each species to ensure stability, representing intraspecies
7 interactions. We use this ecological network to identify the driver species of the community, and
8 its corresponding weighted adjacency matrix as the interaction matrix to construct the linear MPC.
9 The parameters $Q = 20 \times 10^4 I_{N \times N}$, $R = 0.15 I_{M \times M}$ of the linear MPC were fixed for all commu-
10 nities, and the intervention time instants $t_k \in \mathbb{T}$ were chosen such that $t_{k+1} - t_k = 0.1$. Next, we
11 used Eq. (5) to numerically simulate the population dynamics of these communities. For this, we
12 set the weighted adjacency matrix of the ecological network we built as the interaction matrix A in
13 Eq. (5). We choose $\theta_{i,j} = 0$ for $j = 1, \dots, 6$, and $\theta_{i,j,7}$ uniformly at random from $[0, \theta_{\max}]$, where
14 θ_{\max} is a parameter. Last, we choose the intrinsic growth rates r_i to ensure all generated random
15 communities share the desired state $x_d \in \mathbb{R}^N$ as an equilibrium point. Note all the constructed
16 communities have nonlinear population dynamics, and their linearization at the desired state is not
17 equal to the interaction matrix used for the linear MPC (see Supplementary Note 8 for details of
18 this construction).

19 To quantify the success of our control framework on a particular community, we generate 300
20 initial species abundances that are uniformly distributed at a distance $d > 0$ from the desired state
21 (distance is measured using the Euclidean norm). Then, the *success rate* of our control framework

1 at distance d is defined as the proportion of those initial conditions that are driven to the desired
2 state only when the linear MPC is applied to a minimal set of driver species of the community
3 (Fig. 4b-d). Namely, the success rate discards all initial conditions that naturally evolve to the
4 desired state. Finally, we calculated the *mean success rate* by averaging the success rate over 100
5 randomly constructed ecological networks (see items 7 and 8 of Supplementary Note 8 for details).

6 The mean success rate of our control framework changes with the distance to the desired
7 state, being close to 1 for small distances regardless of the parameters of the microbial community
8 (Fig. 4e-f). This result agrees well with the theoretical prediction that success is guaranteed pro-
9 vided that the distance to the desired state is small enough. We next investigated how the success
10 rate changes with the distance d for different interspecies interaction strengths, and for different
11 connectivities of the ecological network underlying the community. The success rate decreases
12 as the interspecies interaction strength increase, especially for large distances (Fig. 4e). Since
13 increasing the interspecies interaction strength damages the stability of the population dynamics⁴⁴,
14 this result suggests that microbial communities become “harder” to control as they lose stability.
15 The success rate of our control framework is also higher in microbial communities whose ecolog-
16 ical networks have lower connectivity (Fig. 4f). Note that, in general, the size of a minimal set of
17 driver species decreases as the network connectivity increases. Therefore, this observations sug-
18 gest that the success rate may increase as the number of driver species increases. Indeed, regardless
19 of the distance to the desired state, we find that our control framework attains a success rate > 0.8
20 provided that we drive at least 6 of the 100 species (Fig. 4g). This last result also suggests that
21 the success rate of our control framework can be enhanced by directly controlling a few additional

1 species.

2 Finally, we investigated the robustness of our control framework to errors in the ecological
3 network used for both identifying the driver species, and for calculating the linear MPC. Note that,
4 despite structural accessibility is insensitive to missing interactions in the ecological network, the
5 calculated linear MPC is not. Additionally, structural accessibility can be lost if some ecological
6 interactions do not really exist in the ecological network. To introduce errors in the ecological
7 network, we randomly rewire each of its edges with probability $p \in [0, 1]$. This rewiring probability
8 determines the percentage of error introduced to the ecological network (e.g., $p = 0.05$ corresponds
9 to a 5% error). Our control framework is robust to these errors, in the sense that the success rate
10 deteriorates but remains larger than zero despite large errors (Fig. 4h). However, just a 5% error
11 decreases the success rate in about 30%. This result illustrates that our framework is feasible for
12 controlling large microbial communities provided we have an accurate map of their ecological
13 networks.

14 APPLICATION

15 Mapping the ecological network of a microbial community allow us to identify its driver species.
16 We identified a minimal set of driver species in the gut microbiota of germ-free mice that are
17 pre-colonized with a mixture of human commensal bacterial type strains and then infected with
18 *Clostridium difficile* spores²². We identified a minimal set of five driver species in this 14-species
19 community: *R. obeum* (x_1), *R. mirabilis* (x_{12}), *B. ovalus* (x_2), *C. ramnosum* (x_6) and *A. muciniphila*

1 (x_{10}), see Fig. 5a. We also used the ecological network underlying the core microbiota of the sea
2 sponge *Ircinia oros*²³, finding ten driver species in this twenty-species community (Fig. 5b).

3 We studied by simulation the efficacy of the identified driver species and the linear MPC
4 method for these two microbial communities, assuming that their dynamics are uncertain (see
5 Supplementary Note 7 for details of the dynamics used for the simulation). For the mice gut
6 microbiota, our framework succeeds in driving the whole community from an initial state where
7 *Clostridium difficile* is overabundant, towards a desired state with a better balance of species (Figs.
8 5c and 5d). Similar results were obtained for controlling the core microbiota of *Ircinia oros*, using
9 the ten identified driver species to drive the twenty species constituting this microbial community
10 (Figs. 5e and 5f). The success of our control framework shows again that the linear MPC method is
11 robust enough to drive microbial communities despite the presence of the perturbations (w_x, w_u) .

12 **DISCUSSION**

13 An influential method to understand and manage complex ecosystems has been identifying species
14 with a “big impact” on the entire ecosystem, leading to notions such as keystone^{45,46} or core⁴⁷
15 species. In general, the keystone or core species of an ecosystem are not necessarily its driver
16 species. For example, the driver species of an ecosystem do not depend on their abundance, while
17 the definition of keystone species does depend on the abundance —namely, species whose removal
18 cause a disproportionate deleterious effect relative to their abundance⁴⁵.

19 It was suggested that notion of *controllability* —the ability to drive a system between any two

1 states— could help predicting the success of ecosystem management strategies⁴⁸. For microbial
2 communities and many other biological systems, it is inadequate to use the notion of controllability
3 because there are states that those systems cannot reach by their nature (e.g., those states corre-
4 sponding to negative abundances). Additionally, since dynamic models for microbial communities
5 and other complex ecosystems are nonlinear, uncertain, and often very difficult to infer, it is impos-
6 sible to even test if those systems are controllable or not. The notion of structural accessibility at
7 the basis of our framework overcomes these two limitations, generalizing the control-theoretic no-
8 tion of accessibility³² to systems with uncertain dynamics and impulsive control inputs. As result,
9 our framework allows efficiently controlling microbial communities only knowing their underlying
10 ecological networks. We note that our framework can be used to identify minimal sets of “driver
11 variables” for biological systems beyond microbial communities when their underlying networks
12 are known. For this, we just need to choose the adequate base model⁴⁹ for each class of system. For
13 example, we identified a single “driver protein” in the repressilator⁵⁰ —a synthetic three-gene reg-
14 ulatory network that generates sustained oscillations— allowing us to eliminate those oscillations
15 (Supplementary Note 8 and Fig. S2).

16 In this paper, we used a maximum matching based algorithm to identify a minimum set of
17 driver species from the ecological network of a given microbial community. In principle, there
18 could be multiple maximum matchings associated with the same network, rendering potentially
19 different minimum sets of driver species. Note that those minimum driver species sets share the
20 same cardinality. We claim that a minimum set of driver species is optimal only in the sense that
21 its cardinality is minimal. If the cost of choosing any species as a driver species is known, one can

1 develop a combinatorial optimization scheme to further pick up the best driver species set. But we
2 feel this is beyond the scope of the current work and hence leave it for future work.

3 Rather counterintuitively, our mathematical formalism shows that increasing the complexity
4 of the community's population dynamics (measured by the size of the deformation) can only reduce
5 the number of necessary driver species. In practice, however, increasing the complexity of the
6 dynamics could render the design of the control strategies more difficult. Note that, in general,
7 it can be expected that the design of control strategies becomes more difficult as the number of
8 used driver species decreases (see Remark 9 in Supplementary Note 5). Additionally, we note that
9 despite the minimal number of driver species decreases as the ecological network becomes denser,
10 this condition is only sufficient. Indeed, the minimal number of driver species of a microbial
11 community should be mainly determined by the degree distribution of the ecological network, since
12 the maximum matching size of a directed network is largely determined by its degree distribution⁵¹.

13 For large communities with uncertain controlled population dynamics, we calculated the
14 control actions using a linear prediction model with an infinite horizon. More sophisticated control
15 algorithms, such as those based on reinforcement learning⁵² (RL), could provide better perfor-
16 mance. Note that RL algorithms typically require specifying a-priori the “driver variables” they
17 can actuate⁵³. Our characterization of minimal sets of driver species should help to efficiently ap-
18 ply RL methods for controlling microbial communities and other biological systems. In practice,
19 the performance of the control algorithms can also be improved by using more detailed models that
20 incorporate the dynamics of the susceptibility of species to the control actions (e.g., the pharma-

1 cokinetics of prebiotics). In such case, different control actions could be modeled by different pairs
2 $\{f, g\}$ in Eqs. (2) or (3), making the conditions for the absence of autonomous elements different
3 for continuous and impulsive control actions. We note that altering the ecological network of a mi-
4 crobial community or obtaining a “simplified” network, in the spirit of Refs.⁵⁴ and ⁵⁵, respectively,
5 could be an alternative and complementary approach to controlling microbial communities (e.g.,
6 to reduce the number of necessary driver species).

7 Note also that in our deterministic framework we don’t consider the effects of stochasticity
8 due to, e.g., immigration in microbial communities. From a theoretical viewpoint, incorporating
9 stochastic effects into the model will turn Eqs. (2) and (3) into controlled stochastic differential
10 equations, which are the material of a different scientific area. To the best of our knowledge, the
11 characterization of the accessibility properties of those class of equations remains an open prob-
12 lem and their analysis become intractable in practice. Indeed, the very notion of an autonomous
13 element—the basis for the concept of accessibility— would need to be reformulated. We consider
14 this is beyond the scope of the current work and call for research activities of the control theory
15 community in this area.

16 In conclusion, by identifying driver species, our framework shows that an accurate map of
17 the ecological network underlying a microbial community opens the door for an efficient and sys-
18 tematic control. The driver species can be identified despite missing interactions in the ecological
19 network, but our methods to calculate the adequate control actions can be sensitive to them. The
20 design of controllers that are robust to missing interactions will be a necessary step for controlling

1 real microbial communities. To fully harvest the potential benefits of controlling microbial com-
2 munities a stronger synergy between microbiology, ecology, and control theory will be necessary.

3 **Correspondence** Correspondence and requests for materials should be addressed to M.T.A. (email: man-
4 gulo@im.unam.mx) or Y.-Y.L. (email: yyl@channing.harvard.edu).

5 **Data availability** All the experimental datasets analyzed in this study are publicly available.

6 **Code availability** A Julia implementation of the algorithm for identifying a minimal set of driver species,
7 as well as all other functions necessary to reproduce the results of the paper, is provided at the GitHub
8 repository: <https://github.com/mtangulo/DriverSpecies>.

9 **BIBLIOGRAPHY**

- 10 1. Pepper, J. W. & Rosenfeld, S. The emerging medical ecology of the human gut microbiome.
12 *Trends in ecology & evolution* **27**, 381–384 (2012).
- 13 2. DeLeon-Rodriguez, N. *et al.* Microbiome of the upper troposphere: Species composition
14 and prevalence, effects of tropical storms, and atmospheric implications. *Proceedings of the*
15 *National Academy of Sciences* **110**, 2575–2580 (2013).
- 16 3. Buchan, A., LeClerc, G. R., Gulvik, C. A. & González, J. M. Master recyclers: features and
17 functions of bacteria associated with phytoplankton blooms. *Nature Reviews Microbiology* **12**,
18 686–698 (2014).

- 1 4. Hultman, J. *et al.* Multi-omics of permafrost, active layer and thermokarst bog soil micro-
2 biomes. *Nature* **521**, 208–212 (2015).
- 3 5. Karczewski, J., Poniedziałek, B., Adamski, Z. & Rzymiski, P. The effects of the microbiota on
4 the host immune system. *Autoimmunity* **47**, 494–504 (2014).
- 5 6. Cox, L. M. & Blaser, M. J. Antibiotics in early life and obesity. *Nature Reviews Endocrinology*
6 **11**, 182–190 (2015).
- 7 7. Tang, A. T. *et al.* Endothelial tlr4 and the microbiome drive cerebral cavernous malformations.
8 *Nature* **545**, 305–310 (2017).
- 9 8. East, R. Microbiome: Soil science comes to life. *Nature* **501**, S18–S19 (2013).
- 10 9. Mueller, U. G. & Sachs, J. L. Engineering microbiomes to improve plant and animal health.
11 *Trends in microbiology* **23**, 606–617 (2015).
- 12 10. Guidi, L. *et al.* Plankton networks driving carbon export in the oligotrophic ocean. *Nature*
13 **532**, 465–470 (2016).
- 14 11. Alivisatos, A. P. *et al.* A unified initiative to harness earth’s microbiomes. *Science* **350**, 507–
15 508 (2015).
- 16 12. N. Dubilier, M. M.-N. & Zhao, L. Create a global microbiome effort. *Nature* (2015).
- 17 13. Wubs, E. J., van der Putten, W. H., Bosch, M. & Bezemer, T. M. Soil inoculation steers
18 restoration of terrestrial ecosystems. *Nature plants* **2**, 16107 (2016).

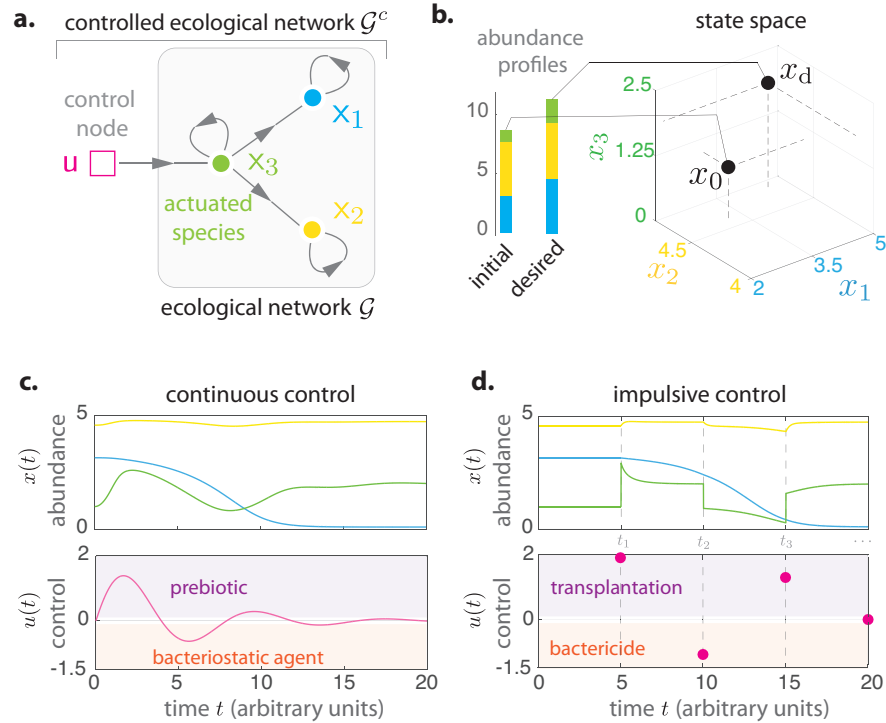
- 1 14. Buffie, C. G. *et al.* Precision microbiome reconstitution restores bile acid mediated resistance
2 to *Clostridium difficile*. *Nature* **517**, 205–208 (2015).
- 3 15. Gibson, T. E., Bashan, A., Cao, H.-T., Weiss, S. T. & Liu, Y.-Y. On the origins and control of
4 community types in the human microbiome. *PLoS Comput Biol* **12**, e1004688 (2016).
- 5 16. Lin, C. T. Structural controllability. *Automatic Control, IEEE Transactions on* **19**, 201–208
6 (1974).
- 7 17. Liu, Y.-Y. & Barabási, A.-L. Control principles of complex systems. *Reviews of Modern*
8 *Physics* **88**, 035006 (2016).
- 9 18. Phelan, V. V., Liu, W.-T., Pogliano, K. & Dorrestein, P. C. Microbial metabolic exchange
10 [mdash] the chemotype-to-phenotype link. *Nature chemical biology* **8**, 26–35 (2012).
- 11 19. Turchin, P. *Complex population dynamics: a theoretical/empirical synthesis*, vol. 35 (Prince-
12 ton University Press, 2003).
- 13 20. Faust, K. & Raes, J. Microbial interactions: from networks to models. *Nature Reviews Micro-*
14 *biology* **10**, 538–550 (2012).
- 15 21. Friedman, J., Higgins, L. M. & Gore, J. Community structure follows simple assembly rules
16 in microbial microcosms. *Nature Ecology & Evolution* 0109 (2017).
- 17 22. Bucci, V. *et al.* Mdsine: Microbial dynamical systems inference engine for microbiome time-
18 series analyses. *Genome biology* **17**, 121 (2016).

- 1 23. Thomas, T. *et al.* Diversity, structure and convergent evolution of the global sponge micro-
2 biome. *Nature Communications* **7** (2016).
- 3 24. Xiao, Y. *et al.* Mapping the ecological networks of microbial communities. *Nature communi-*
4 *cations* **8**, 2042 (2017).
- 5 25. Angulo, M. T., Moreno, J. A., Lippner, G., Barabási, A.-L. & Liu, Y.-Y. Fundamental limita-
6 tions of network reconstruction from temporal data. *Journal of the Royal Society Interface* **14**,
7 20160966 (2017).
- 8 26. Sugihara, G. *et al.* Detecting causality in complex ecosystems. *science* 1227079 (2012).
- 9 27. Friedman, J. & Alm, E. J. Inferring correlation networks from genomic survey data. *PLoS*
10 *computational biology* **8**, e1002687 (2012).
- 11 28. Berry, D. & Widder, S. Deciphering microbial interactions and detecting keystone species
12 with co-occurrence networks. *Frontiers in microbiology* **5**, 219 (2014).
- 13 29. Waksman, S. A. What is an antibiotic or an antibiotic substance? *Mycologia* **39**, 565–569
14 (1947).
- 15 30. Oremland, R. S. & Capone, D. G. Use of “specific” inhibitors in biogeochemistry and micro-
16 bial ecology. In *Advances in microbial ecology*, 285–383 (Springer, 1988).
- 17 31. Schrezenmeir, J. & de Vrese, M. Probiotics, prebiotics, and synbiotics—approaching a defi-
18 nition. *The American journal of clinical nutrition* **73**, 361s–364s (2001).

- 1 32. Conte, G., Moog, C. H. & Perdon, A. M. *Algebraic methods for nonlinear control systems*
2 (Springer Science & Business Media, 2007).
- 3 33. Moore, J. C., de Ruiter, P. C., Hunt, H. W., Coleman, D. C. & Freckman, D. W. Microcosms
4 and soil ecology: critical linkages between fields studies and modelling food webs. *Ecology*
5 **77**, 694–705 (1996).
- 6 34. Mounier, J. *et al.* Microbial interactions within a cheese microbial community. *Applied and*
7 *environmental microbiology* **74**, 172–181 (2008).
- 8 35. Stein, R. R. *et al.* Ecological modeling from time-series inference: insight into dynamics and
9 stability of intestinal microbiota. *PLoS Comput Biol* **9**, e1003388 (2013).
- 10 36. Gerber, G. K. The dynamic microbiome. *FEBS letters* **588**, 4131–4139 (2014).
- 11 37. Coyte, K. Z., Schluter, J. & Foster, K. R. The ecology of the microbiome: Networks, compe-
12 tition, and stability. *Science* **350**, 663–666 (2015).
- 13 38. Bashan, A. *et al.* Universality of human microbial dynamics. *Nature* **534**, 259–262 (2016).
- 14 39. Dam, P., Fonseca, L. L., Konstantinidis, K. T. & Voit, E. O. Dynamic models of the complex
15 microbial metapopulation of lake mendota. *NPJ Systems Biology and Applications* **2**, 16007
16 (2016).
- 17 40. Jost, C. & Ellner, S. P. Testing for predator dependence in predator-prey dynamics: a non-
18 parametric approach. *Proceedings of the Royal Society of London B: Biological Sciences* **267**,
19 1611–1620 (2000).

- 1 41. Camacho, E. F. & Alba, C. B. *Model predictive control* (Springer Science & Business Media,
2 2013).
- 3 42. Jones, D. R., Perttunen, C. D. & Stuckman, B. E. Lipschitzian optimization without the
4 lipschitz constant. *Journal of Optimization Theory and Applications* **79**, 157–181 (1993).
- 5 43. Aström, K. J. & Murray, R. M. *Feedback systems: an introduction for scientists and engineers*
6 (Princeton university press, 2010).
- 7 44. May, R. M. *Stability and complexity in model ecosystems*, vol. 6 (Princeton university press,
8 2001).
- 9 45. Power, M. E. *et al.* Challenges in the quest for keystones: identifying keystone species is dif-
10 ficult—but essential to understanding how loss of species will affect ecosystems. *BioScience*
11 **46**, 609–620 (1996).
- 12 46. Ortiz, M. *et al.* Quantifying keystone species complexes: ecosystem-based conservation man-
13 agement in the king george island (antarctic peninsula). *Ecological Indicators* **81**, 453–460
14 (2017).
- 15 47. Jain, S. & Krishna, S. Crashes, recoveries, and “core shifts” in a model of evolving networks.
16 *Physical Review E* **65**, 026103 (2002).
- 17 48. Loehle, C. Control theory and the management of ecosystems. *Journal of applied ecology* **43**,
18 957–966 (2006).
- 19 49. Barzel, B. & Barabási, A.-L. Universality in network dynamics. *Nature physics* **9**, 673 (2013).

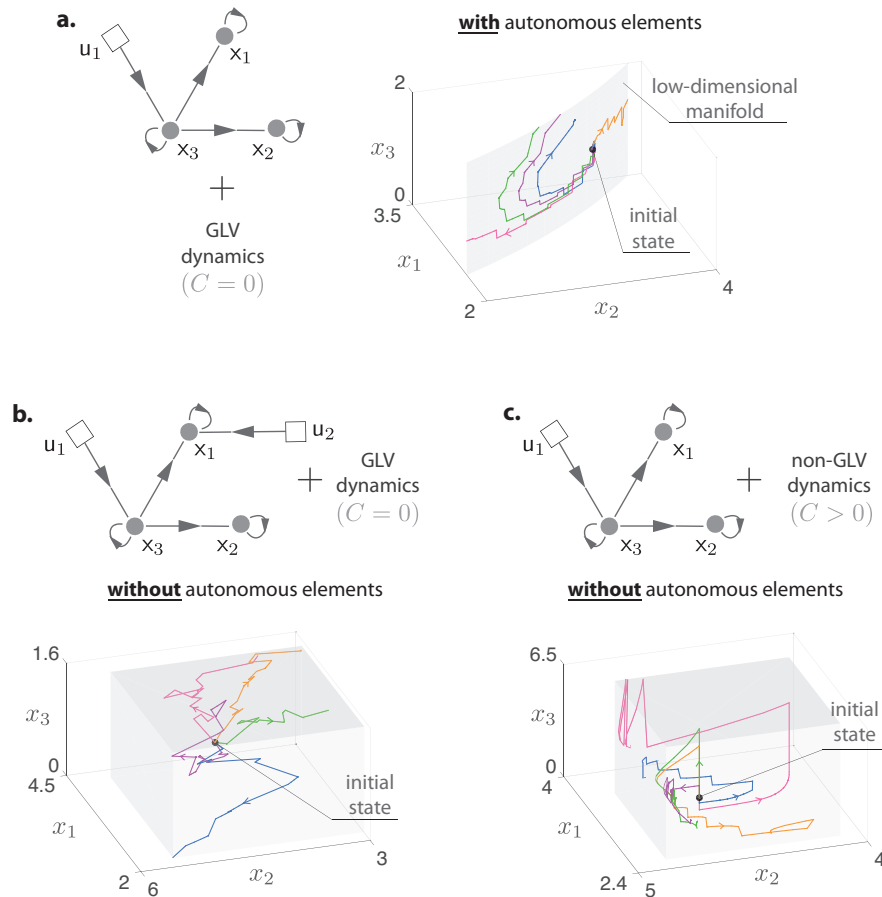
- 1 50. Buse, O., Pérez, R. & Kuznetsov, A. Dynamical properties of the repressilator model. *Physical*
2 *Review E* **81**, 066206 (2010).
- 3 51. Liu, Y.-Y., Slotine, J.-J. & Barabási, A.-L. Controllability of complex networks. *Nature* **473**,
4 167–173 (2011).
- 5 52. Sutton, R. S. & Barto, A. G. *Introduction to reinforcement learning*, vol. 135 (MIT Press
6 Cambridge, 1998).
- 7 53. Mnih, V. *et al.* Human-level control through deep reinforcement learning. *Nature* **518**, 529–
8 533 (2015).
- 9 54. Campbell, C. & Albert, R. Stabilization of perturbed boolean network attractors through
10 compensatory interactions. *BMC systems biology* **8**, 53 (2014).
- 11 55. Zañudo, J. G. & Albert, R. An effective network reduction approach to find the dynamical
12 repertoire of discrete dynamic networks. *Chaos: An Interdisciplinary Journal of Nonlinear*
13 *Science* **23**, 025111 (2013).



1

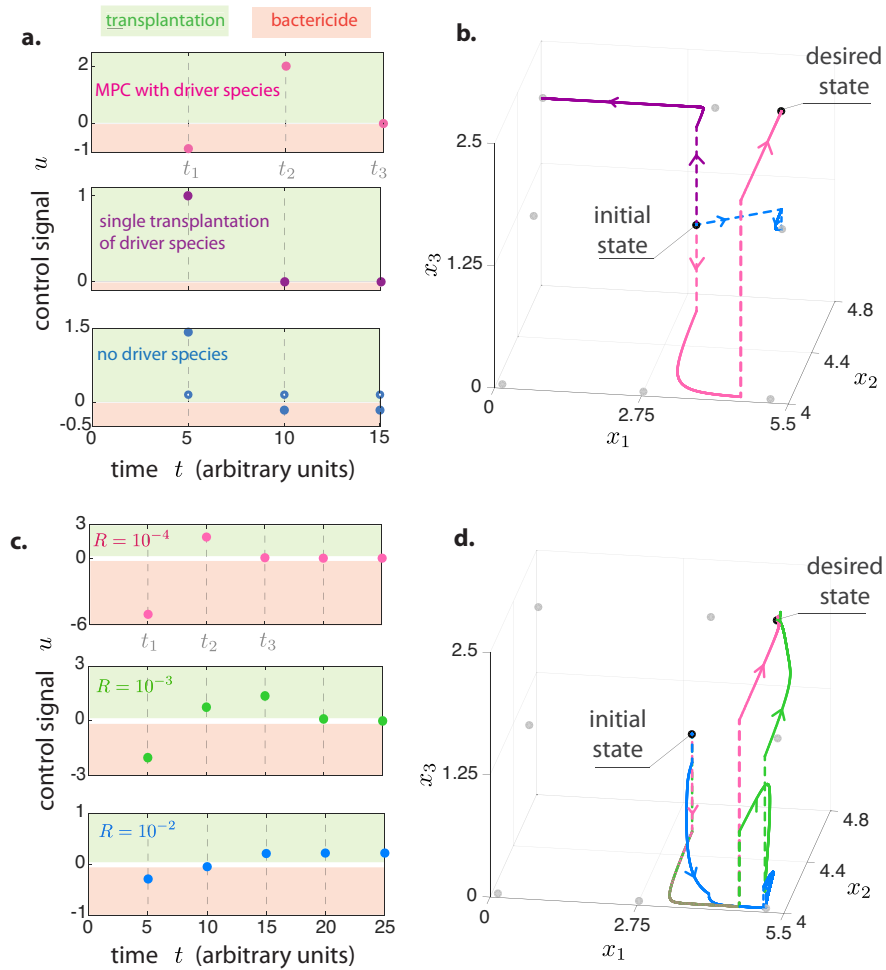
Figure 1 Controlling a microbial community. **a.** Ecological network \mathcal{G} for a toy microbial community of $N = 3$ species (green, yellow, blue). The controlled ecological network \mathcal{G}^c contains $M = 1$ control input actuating the third species. **b.** Initial and desired abundance profiles (bars). Controlling the community consists in driving its state from the initial state x_0 to the desired state x_d , represented by two points in the state space of the community. **c.** In the continuous control scheme, the control inputs $u(t)$ are continuous signals modifying the growth of the actuated species. The controlled population dynamics of this community is given by $\dot{x}_1 = 0.1 + x_1(1 - x_1/5)(x_1/3 - 1) - (0.1x_1x_3)/(1 + x_3)$, $\dot{x}_2 = 0.1 + x_2(1 - x_2/4)(x_2 - 1) + (x_2x_3)/(1 + x_3)$, $\dot{x}_3 = x_3(1 - x_3/2)(x_3 - 1) + u$. In the absence of control, this community has two equilibria $x_0 = (3.14, 4.58, 1)^T$ and $x_d = (4.57, 4.73, 2)^T$, chosen as the initial and desired states, respectively. **d.** In the impulsive control scheme, the control inputs $u(t)$ are impulses applied at the intervention instants $\mathbb{T} = \{t_1, t_2, \dots\}$, instantaneously changing the abundance of the actuated species. The controlled population dynamics is the same as in panel c, except that $\dot{x}_3 = x_3(1 - x_3/K_3)(x_3/C_3 - 1)$ and $x_3(t^+) = x_3(t) + u(t)$ if $t \in \mathbb{T} = \{5, 10, 15\}$. Under this controlled population dynamics, our mathematical formalism identifies x_3 as the solo driver species needed to drive this microbial community (Example 1 in Supplementary Note 2).

2



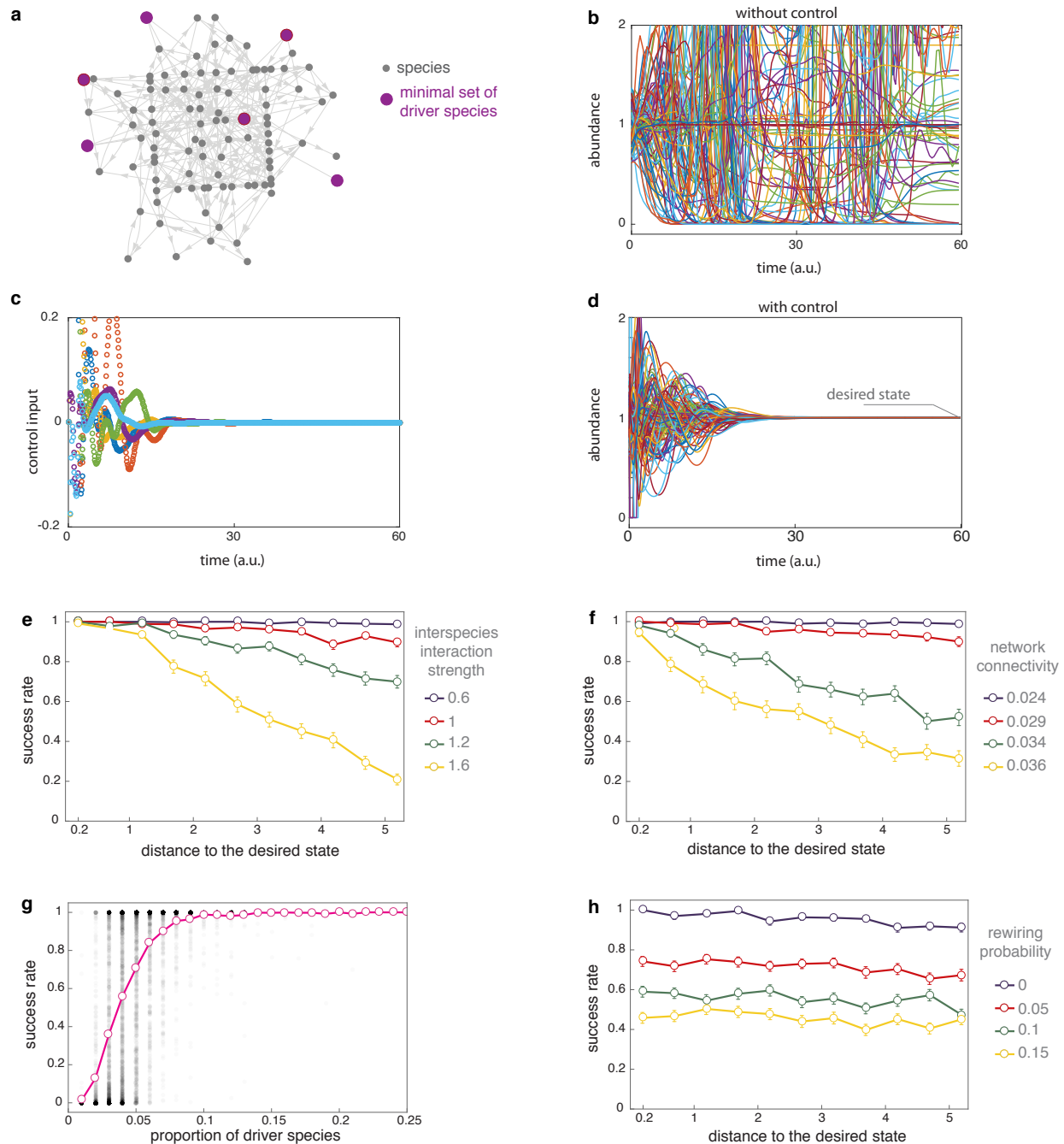
1

2 **Figure 2** Autonomous elements constrain the state of microbial communities, characterizing their
3 driver species. **a.** A three-species community with GLV dynamics $\dot{x}_1 = x_1(-1+x_3)$, $\dot{x}_2 = x_2(1-x_3)$,
4 $\dot{x}_3 = x_3(-0.5 + 1.5x_3)$. For actuating x_3 , we consider the impulsive control scheme with $x_3(t^+) =$
5 $x_3(t) + u_1(t)$ for $t \in \mathbb{T}$. With this controlled population dynamics, our mathematical formalism
6 reveals the autonomous element x_1x_2 that constraints the state of this microbial community to
7 the low-dimensional manifold $\{x \in \mathbb{R}^3 | x_1x_2 = x_1(0)x_2(0)\}$ (gray) for all control inputs. Five state
8 trajectories (in colors) with random control inputs illustrate this fact. Hence, $\{x_3\}$ alone cannot be
9 a set of driver species for this controlled population dynamics. **b.** Including a second control input
10 $u_2(t)$ actuating x_1 (i.e., $x_1(t^+) = x_1(t) + u_2(t)$ for $t \in \mathbb{T}$) eliminates the autonomous element, since
11 the state of the microbial community (colors) can explore a three-dimensional space (gray). Hence
12 $\{x_1, x_3\}$ is a minimal set of driver species for this community with GLV dynamics. **c.** We proved
13 that, generically, increasing the complexity of the controlled population dynamics cannot create
14 autonomous elements. In this example, increasing the deformation size C from the GLV in panel a
15 (with $C = 0$) to the controlled population dynamics in Fig.1 (with $C > 0$) eliminates the autonomous
16 element that was present by actuating x_3 alone (Example 1 in Supplementary Note 2). Therefore,
17 increasing the complexity of the population dynamics makes $\{x_3\}$ a solo driver species.



1

2 **Figure 3 Success and failure of different control strategies.** **a.** Three control strategies for driving
3 the microbial community of Fig. 1a toward the desired state. First, MPC applied to the identified
4 driver species $\{x_3\}$ (pink dots). The second control strategy increases the abundance of the driver
5 species to match its value at the desired state $x_3(t_1) = x_{3,d}$ (purple dots). The third control
6 strategy does not actuate the driver species, but actuates the other two species $\{x_1, x_2\}$ by setting
7 their abundance to their desired values (i.e., $x_1(t_k) = x_{1,d}$ and $x_2(t_k) = x_{2,d}$, solid and hollow blue
8 dots, respectively). **b.** The response of the microbial community to these three control strategies.
9 Here and in panel d, the “jumps” produced by the control inputs are depicted by dashed lines. The
10 equilibria of the population dynamics are shown as gray dots. Only the first strategy applying MPC
11 to the driver species succeeds in driving the community to x_d . **c.** Control strategies obtained by
12 using the linear MPC with parameters $Q = \text{diag}(20, 1, 10)$ and different values for R : 10^{-4} (pink),
13 10^{-3} (green), 10^{-2} (blue). **d.** Trajectories of the controlled community using the linear MPC control
14 strategies described in panel c. Colors correspond to the different values of R .

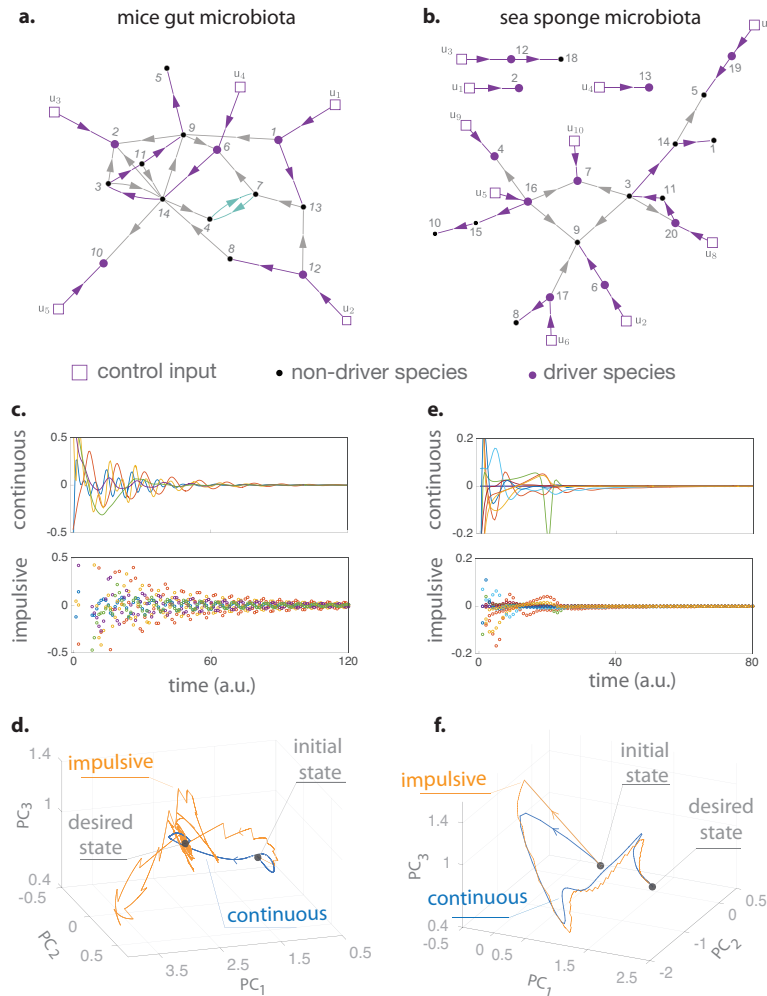


1

2

3 **Figure 4 Numerical validation of the control framework on large microbial communities.** **a.** Ex-
 4 ample of the ecological network of a random microbial community with $N = 100$ species with
 5 connectivity $c = 0.03$. We used our framework to identify a minimal set of $M = 6$ driver species.
 6 The desired state is chosen as $x_d = (1, \dots, 1)^T$. **b.** We randomly set the initial abundance x_0 of
 7 species at a distance $d = 0.4$ from the desired state x_d . Without control, the state of the microbial
 8 community does not reach the desired state x_d . **c. and d.** For the same community and initial
 9 abundance as in panel b, we apply the control input generated by the linear MPC (panel c) to

1 the six driver species we identified. This control strategy successfully drives the state of the com-
2 munity to the desired state (panel d). **e., f. and h.** Mean success rate of our control framework
3 as a function of the distance d of the initial state from the desired state. Error bars denote the
4 standard error of the mean. Here, the simulation parameters are: $c = 0.025, \theta_{\max} = 0.05$ for panel
5 e, $\sigma = 0.8, \theta_{\max} = 0.05$ for panel f, and $c = 0.025, \sigma = 0.8, \theta_{\max} = 0.05$ for panel h. **g.** Success
6 rate of our control framework for different proportions of driver species M/N . Black dots show the
7 success rate of 7700 random communities plotted as a function of the proportion of driver species.
8 Pink shows the mean success rate as a function of the proportion of driver species.



1

2 **Figure 5 Controlling host-associated microbial communities.** **a.** Inferred ecological network of
3 the gut microbiota of germ-free mice pre-colonized with a mixture of human commensal bacterial
4 type strains and then infected with *C. difficile* (species 7). **b.** Inferred ecological network of the
5 core microbiota of the sea sponge *Ircinia oros*. In both networks, self-loops are omitted to improve
6 readability. A minimal set of driver species is shown, providing a disjoint union of paths (purple)
7 and cycles (green) covering all species nodes. Refer to Table 1 in Supplementary Note 7 for the
8 species name. The controlled population dynamics of both microbial communities were simulated
9 using the cGLV equations (see Supplementary Note 7 for details). The intrinsic growth rates
10 were adjusted such that the community has an initial “diseased” equilibrium state x_0 in which one
11 species (*C. difficile* for the mice gut microbiota) is overabundant compared to the rest of species.
12 We chose the desired state x_d as another equilibrium with a more balanced abundance profile. **c, e.**
13 Control actions obtained using the linear MPC for the impulsive and continuous control schemes.
14 **d, f.** Projection of the high-dimensional abundance profiles (states of the microbial communities)
15 into their first three principal components (PCs). See Supplementary Fig.S1 for the temporal
16 response of each species. The calculated control strategies applied to the driver species succeed
17 in driving the community to the desired state, using either continuous or impulsive control.