## Supplementary Information

### 1. Inference of ECNs from longitudinal data

We consider that abundance of $O$ bacterial operational taxonomic units (OTUs) are measured over a period of $T$ days in $S$ subjects. We model the read counts $n_{os}(t)$ of OTUs "$o$" on any given day $t$ in subject $s$ as a multinomial distribution:

$$p(\{n_{os}(t)\}) = \frac{N_s(t)!}{\prod_o n_{os}(t)!} \prod_o q_{os}(t)^{n_{os}(t)} \tag{S1}$$

where $N_s(t) = \sum_o n_{os}(t)$ is the total read count on a given day and $q_{os}(t)$ are the underlying propensities for individual OTUs. We model these propensities using the exponential Gibbs-Boltzmann distribution which allows us to capture large variations in OTU abundances[33]

$$q_{os}(t) = \frac{1}{\Omega_{st}} \exp\left(-\sum_{k=1}^{K} z_k(t)\theta_{kos}\right) \tag{S2}$$

where $z_k(t)$ are time-specific latents that are shared by all OTUs and subjects, and $\theta_{kos}$ are OTU- and subject-specific loadings that are shared across all time points. The number $K$ of latents/loadings is chosen such that $K \ll O, T$ thereby achieving a lower dimensional description of the time series data. We obtain the $z$s and the $\theta$s using the maximum likelihood approach.

To that end, we write down the log-likelihood of the data:

$$L = const. + \sum_{t,o,s} n_{os}(t) \log q_{os}(t). \tag{S3}$$

The constant term of the likelihood does not depend on the parameters and can thus be omitted in likelihood maximization. Simplifying using Eq. S1 and S2, we have

$$L = -\sum_{t,o,s,k} N_s(t)x_{os}(t)\, z_k(t)\theta_{kos} - \sum_{t,s} \log \Omega_{st} \tag{S4}$$

Here $x_{os}(t) = n_{os}(t)/N_s(t)$ is the relative abundance of OTU $o$ at time $t$. We obtain the gradients

$$\frac{\partial L}{\partial z_k(t)} = -\sum_{o,s} N_s(t)\big(x_{os}(t) - q_{os}(t)\big)\theta_{kos} \text{ and} \tag{S5}$$

$$\frac{\partial L}{\partial \theta_{kos}} = -\sum_{t} N_s(t)z_k(t)\big(x_{os}(t) - q_{os}(t)\big) \tag{S6}$$

We use gradient ascent algorithm to find the latents and the loadings that maximize the likelihood.

For a given $K$, using the microbiome data $x_{os}(t)$ and starting from random initialization, we first simultaneously infer the latents $z_k(t)$ and the features $\Theta_{kos}$. We observe that the $T \times K$ matrix

583    $z$ of latents can be multiplied by an invertible matrix $B$ ($z \rightarrow zB$) and the corresponding matrix
584    $K \times O \times S$ matrix of features can be multiplied by the inverse $B^{-1}$ ($\Theta \rightarrow B^{-1}\Theta$) and the
585    abundance predictions from the model do not change. Therefore, we use the Gram-Schmidt
586    procedure to orthogonalize the matrix of latents such that $z \rightarrow z'$ where $z'^T z' = I_K$ is an identity
587    matrix. For any matrix of latents $z$, the matrix multiplier $B$ that leads to the orthonormal
588    transformation can be found by solving the equation $B^T(z^T z)B = I_K$. Once $B$ is identified, we
589    also transform the $\Theta$ matrix ($\Theta \rightarrow \Theta' = B^{-1}\Theta$). At the end of this procedure, we end up with
590    orthonormal latents $z'$ and corresponding features $\Theta'$ that correspond to the same abundances
591    as $z$ and $\Theta$. For the sake of simplicity of notation, we drop the primes.
592
593    Next, we model the dynamics of the orthonormal latents using a linear dynamical system:

$$z_k(t) = \sum_{k'} A_{kk'} z_{k'}(t) + u_k + \eta_k(t) \tag{S7}$$

595    where we assume that $A_{kk'} = A_{k'k}$, and $\eta_k(t)$ are Gaussian distributed uncorrelated noise
596    vectors: $\langle \eta_k(t_1) n_{k'}(t_2) \rangle = \delta_{12}\delta_{kk'}$ where $\delta_{ab}$ is the Kronecker delta function. Our task is to find
597    the interaction matrix $A$ and the vector $u$ that fits this model. We achieve this using squared error
598    minimization. We write

$$E(A, u) = \sum_t \left( z_k(t) - z_{k,pred}(t) \right)^2 \tag{S7}$$

600    where $z_k(t)$ is the inferred latent and $z_{k,pred}(t)$ is the corresponding prediction using $z_k(t-1)$
601    and Eq. S7. We restrict the summation only over time points $t$ such that measurements are
602    available for time points $t$ and $t - 1$. The squared error is minimized using a simulated annealing
603    approach. Once the matrix $A$ is identified, we transform the orthonormal latents $z_k(t)$ into
604    ecological normal modes $y_k(t)$ as described in the manuscript.
605
606    The scripts for obtaining ECNs $y$ and corresponding loadings $\Phi$ from read count data can be
607    found at: https://github.com/mayar-shahin/EMBED.
608

### 2. Generating *in silico* data
610    **Out of Phase Sinusoids.** We generated 40 OTU abundances for 30 time points. The un-
611    normalized abundances of 20 OTUs followed sinusoidal oscillation: $X(t) = A_1(B_1 \sin(0.5t) + 1)$
612    and the un-normalized abundances of the other 20 OTUs followed phase-shifted oscillation with
613    the same frequency $X(t) = A_1(B_1 \cos(0.5t) + 1)$. $A$s and $B$s are uniform random numbers
614    between 0 and 1. We normalize the generated abundances to produce the underlying probability
615    distribution of the data. We used multinomial sampling with a sequencing depth of 25000 to
616    generate relative OTU abundances on each day (**SI Fig. 1**, panels **A** and **B**).
617

618    **Exponentials and Sinusoids.** We generated 40 OTU abundances for 30 time points. 20 OTUs

619    followed an exponential decay $X(t) = 10A_1 \exp(-0.1t)$, 10 OTUs oscillated according to

620    $X(t) = A_2(B_2 \sin(0.5t) + 1)$ and 10 OTUs oscillated with double the frequency $X(t) =$

621    $A_3(B_3 \sin(t) + 1)$. As above, $A$s and $B$s are uniform random numbers between 0 and 1. We

622    sampled the abundances using the multinomial distribution as above (**SI Fig. 1**, panels **C** and **D**).

623

624    **Sum of Sinusoids.** We generated 40 OTU abundances for 30 time points. 20 OTUs followed a

625    single high frequency oscillation $X(t) = A_1(B_1 \cos(1.5t) + 1)$. The remaining OTU abundances

626    were generated by the addition of two different sinusoids: $X(t) = A_2(B_2 \sin(0.5t) +$

627    $B_3 \sin(t) + 1)$. As above, $A$s and $B$s are uniform random numbers between 0 and 1. We sampled

628    the abundances using the multinomial distribution as above (**SI Fig. 1**, panels **E** and **F**).

629

630        **3. Obtaining the microbiome time series from sequencing data**

631    **Murine gut microbiome response to oscillating diet.** We downloaded the microbiome

632    abundance time series data on mice fed an alternating diet of high fat high sugar chow (HFHS)

633    and low-fat plant polysaccharide chow (LFPP) from Carmody et al.[25] as described previously[14].

634    Each mouse that was subjected to an oscillatory diet was treated separately. Based on our

635    previous work on technical noise in 16s measurements, we only analyzed OTUs with mean

636    abundances > 0.1%[13] averaged across all time points and mice. On every day, the abundances of

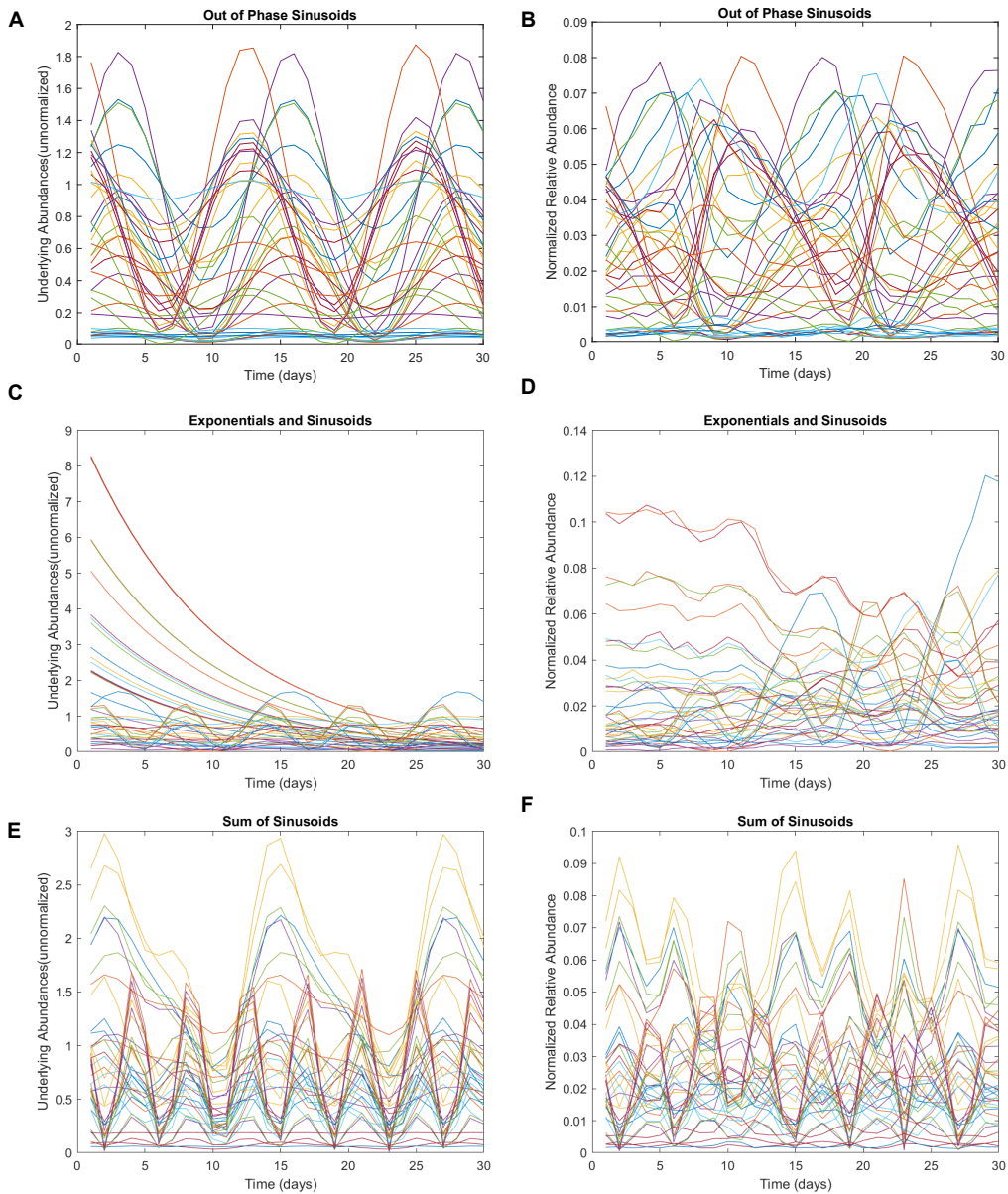637    the rest of the OTUs were lumped together in a single meta-species.

638

639    **Murine gut microbiome response to antibiotics.** We downloaded microbiome abundance data

640    from Ng et al.[10]. We focused on the data where mice were administered the antibiotic

641    ciprofloxacin. Out of the 10 cages in which the mice were housed, we omitted data from cages

642    2, 4, 5, and 8 where many time points were missing. As above, we analyzed OTUs with mean

643    abundance > 0.1% and combined the rest of the OTUs in a meta-species.

644

645        **4. Performing CLR and SSVD on 4 data sets**

646    We downloaded 4 publicly available data sets from 4 different studies[10–12,25]. Each data set

647    comprised microbiome abundance tables for multiple subjects, see SI Table 1. We use the

648    package released by Martino et al.[32] ([https://github.com/biocore/gemelli](https://github.com/biocore/gemelli)) to perform Robust

649    Centered-Log Ratio transform (CLR) on the abundances followed by sparse singular value

650    decomposition. To test how well the dimensionality reduced version capture the data, we

651    calculate an approximate reconstruction of the abundance time series using the first $K$ singular

652    values. As suggested by Martino et al. [31], the resulting approximation was re-exponentiated and

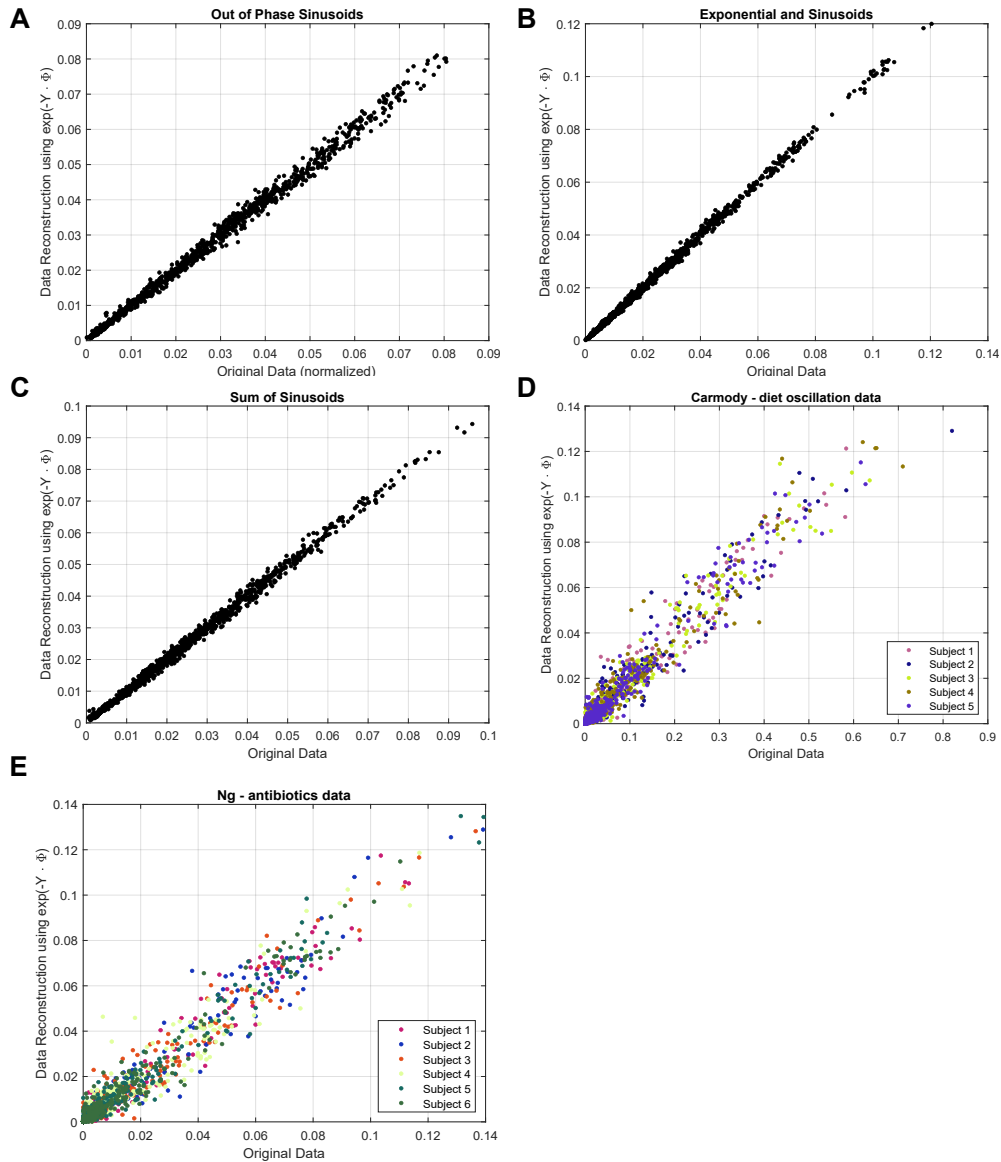653    then normalized. We Then calculated the KL-divergence with the true abundances.

654

655

## Supplementary Figures



**Supplementary Figure 1.** Collective abundance variation of bacterial species in *in silico* data sets representing out of phase oscillations (A and B), exponential decays and oscillations (C and D), and sum of sinusoids (E and F).

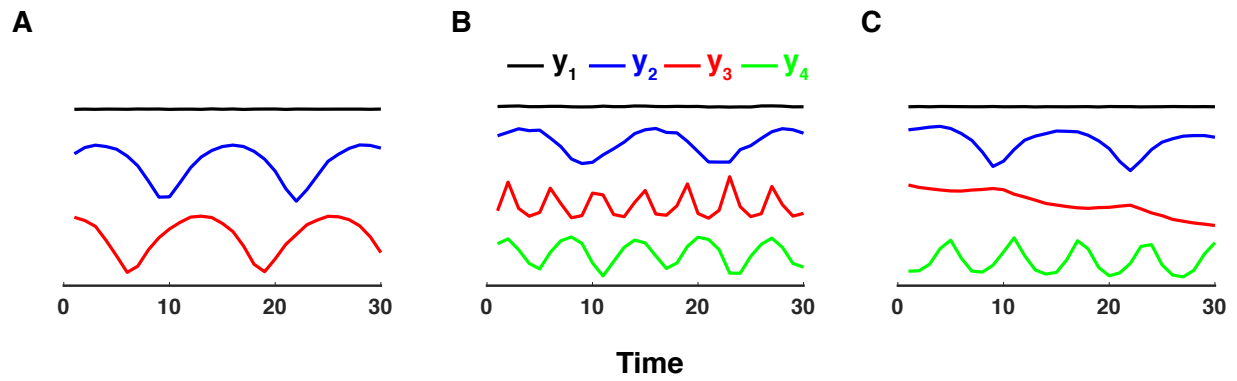**A** — Out of Phase Sinusoids

**B** — Exponential and Sinusoids

**C** — Sum of Sinusoids

**D** — Carmody - diet oscillation data

**E** — Ng - antibiotics data

667
**Supplementary Figure 2.** EMBED-based reconstruction of time series (y axis) compared to the microbiome
time series data (x-axis) for the three *in silico* data sets (A, B, and C) and the two experimental data sets
(D and E) considered in this study.

671

672

673

674

675

676

677

678

679

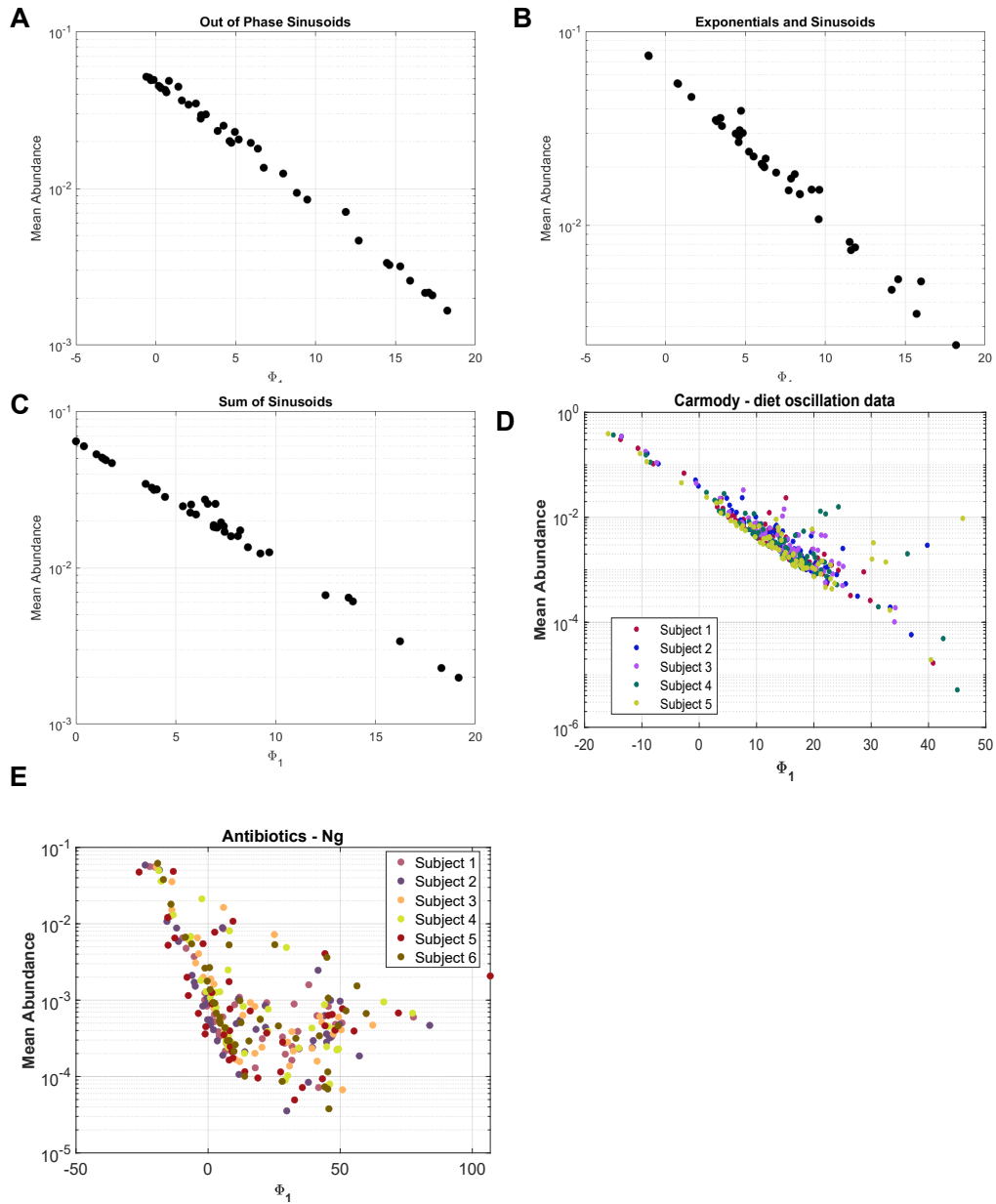**A**     **B**     **C**



681

682    **Supplementary Figure 3.** EMBED-based inference of ecological normal modes for the three in silico data

683    sets. A: out of frequency oscillations, B: sum of sinusoids, and C: exponential decay and sinusoids
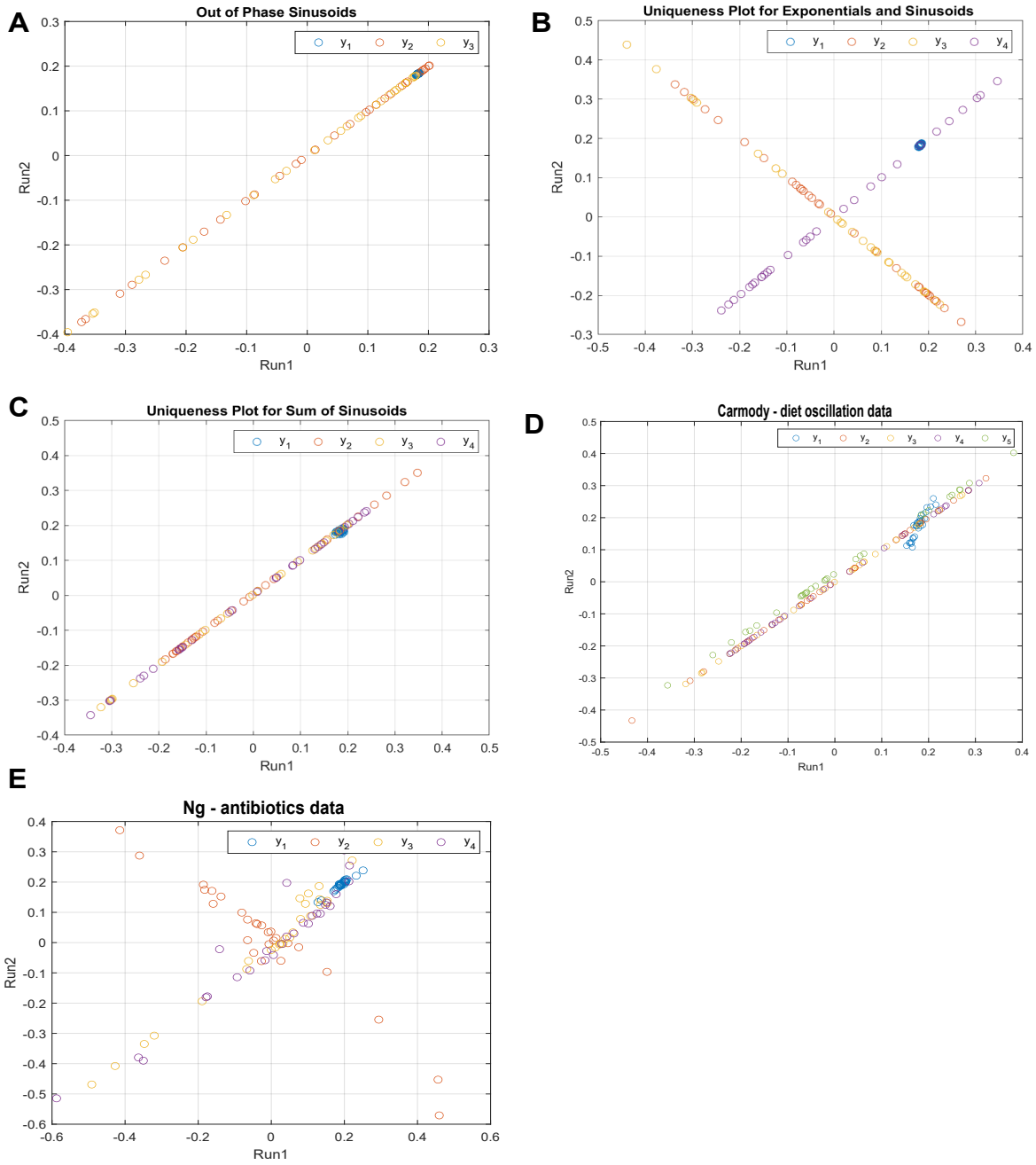
684

685

**Supplementary Figure 4.** Correlation between mean abundance of OTUs and their weight in the first loading $\Phi_1$ for the *in silico* data sets (A, B, and C) and the two experimental data sets (D and E) considered in this study.
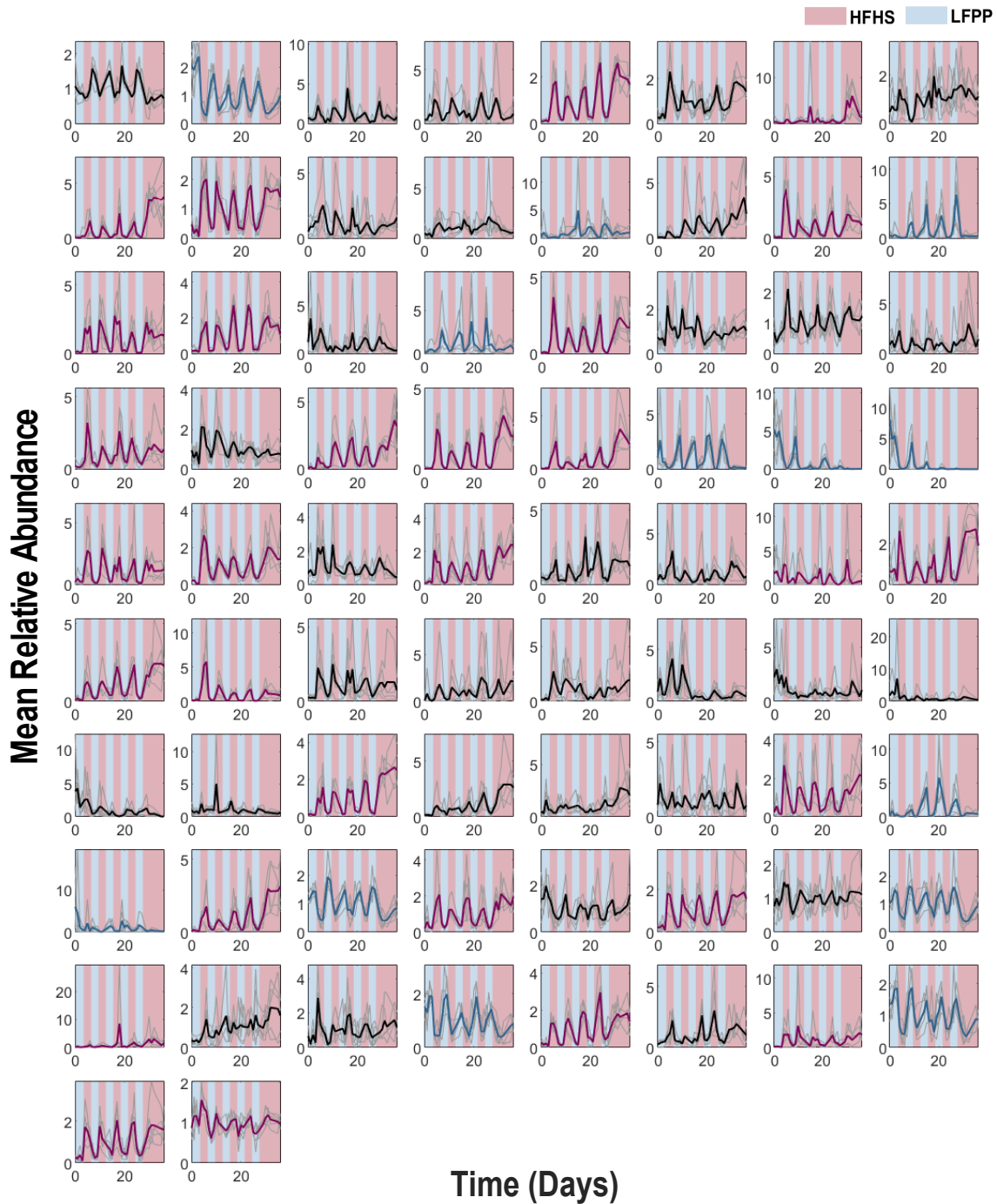
686
687
688
689
690
691
692
693
694
695

696
697 **Supplementary Figure 5.** Correlation plot showing that the ecological normal modes (ECNs) inferred using
698 EMBED are unique (up to a sign). The x- and the y-axis represent the ECNs inferred in two independent
699 runs.
700
701
702
703
704

705



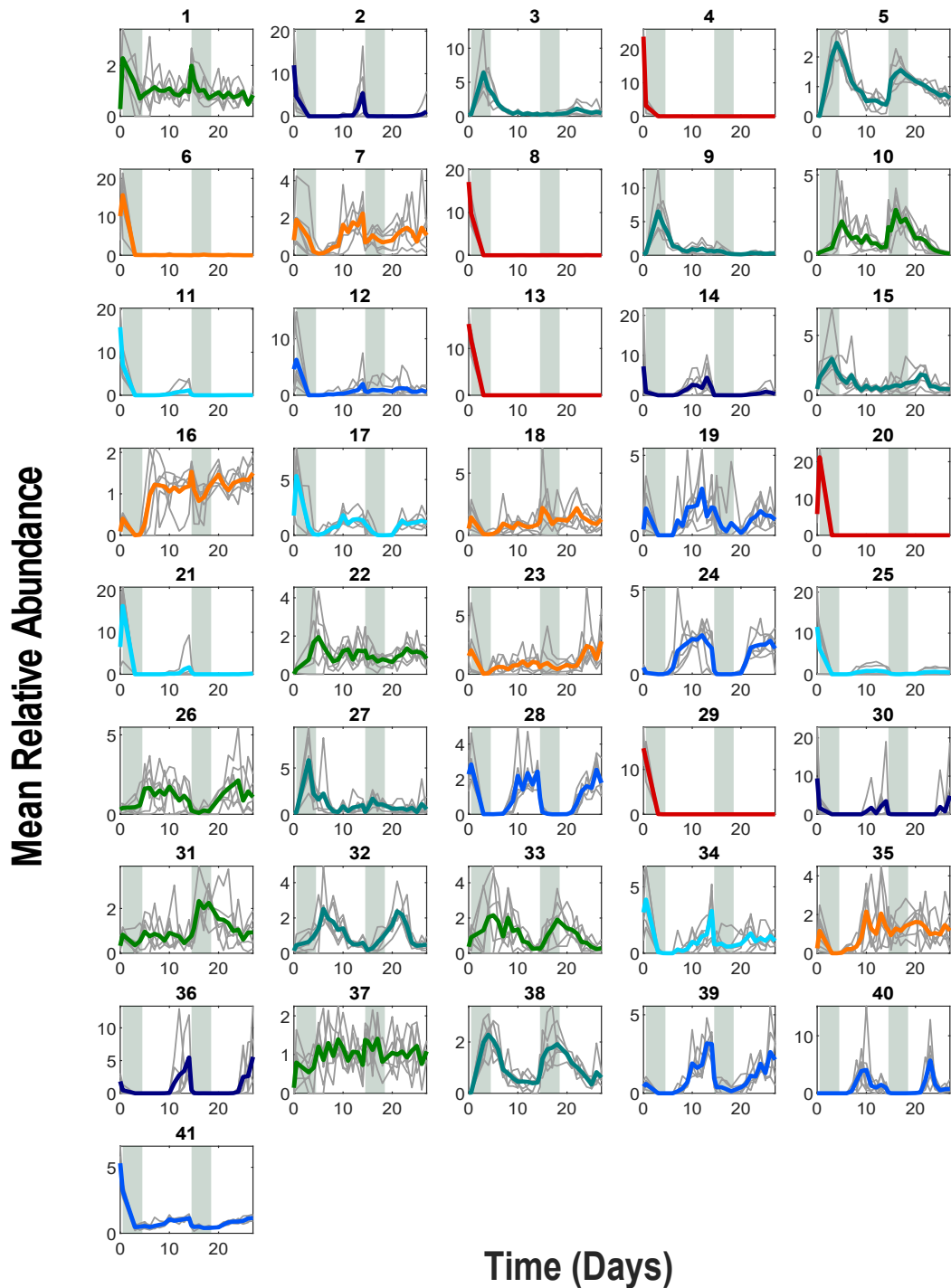**Supplementary Figure 6.** Abundance time series of individual OTUs in the diet oscillation study. The gray lines represent abundances in individual subjects. The dark lines represent averages over subjects. The colors represent the cluster identities in main text Figure 2. HFHS: High Fat High Sugar diet and LFPP: Low Fat Plat Polysaccharide diet.

**Mean Relative Abundance**

**Time (Days)**

713
714 **Supplementary Figure 7.** Abundance time series of individual OTUs in the antibiotics-treatment study.
715 The gray lines represent abundances in individual subjects. The dark lines represent averages over
716 subjects. The colors represent the cluster identities in main text Figure 3. The gray bars represent the
717 duration of time when the antibiotic was administered.
718
719