# Robust and Scalable Models of Microbiome Dynamics

Travis E. Gibson and Georg K. Gerber

Massachusetts Host Microbiome Center, Brigham and Women's Hospital, Harvard Medical School

tgibson@mit.edu
@travisegibson

## Introduction

Microbes are everywhere including in and on our bodies, and have been shown to play key roles in a variety of prevalent human diseases. Consequently, there has been intense interest in the design of bacteriotherapies or "bugs as drugs", which are communities of bacteria administered to patients for specific therapeutic applications. Central to the design of such cocktails is the knowledge (or inference) of a causal microbial interaction network and prediction of the population dynamics of the organisms. In this work we present a Bayesian nonparametric model and associated efficient inference algorithm that addresses the key conceptual and practical challenges of learning microbial dynamics from time series microbe abundance data.

## Interpretability and the Microbiome

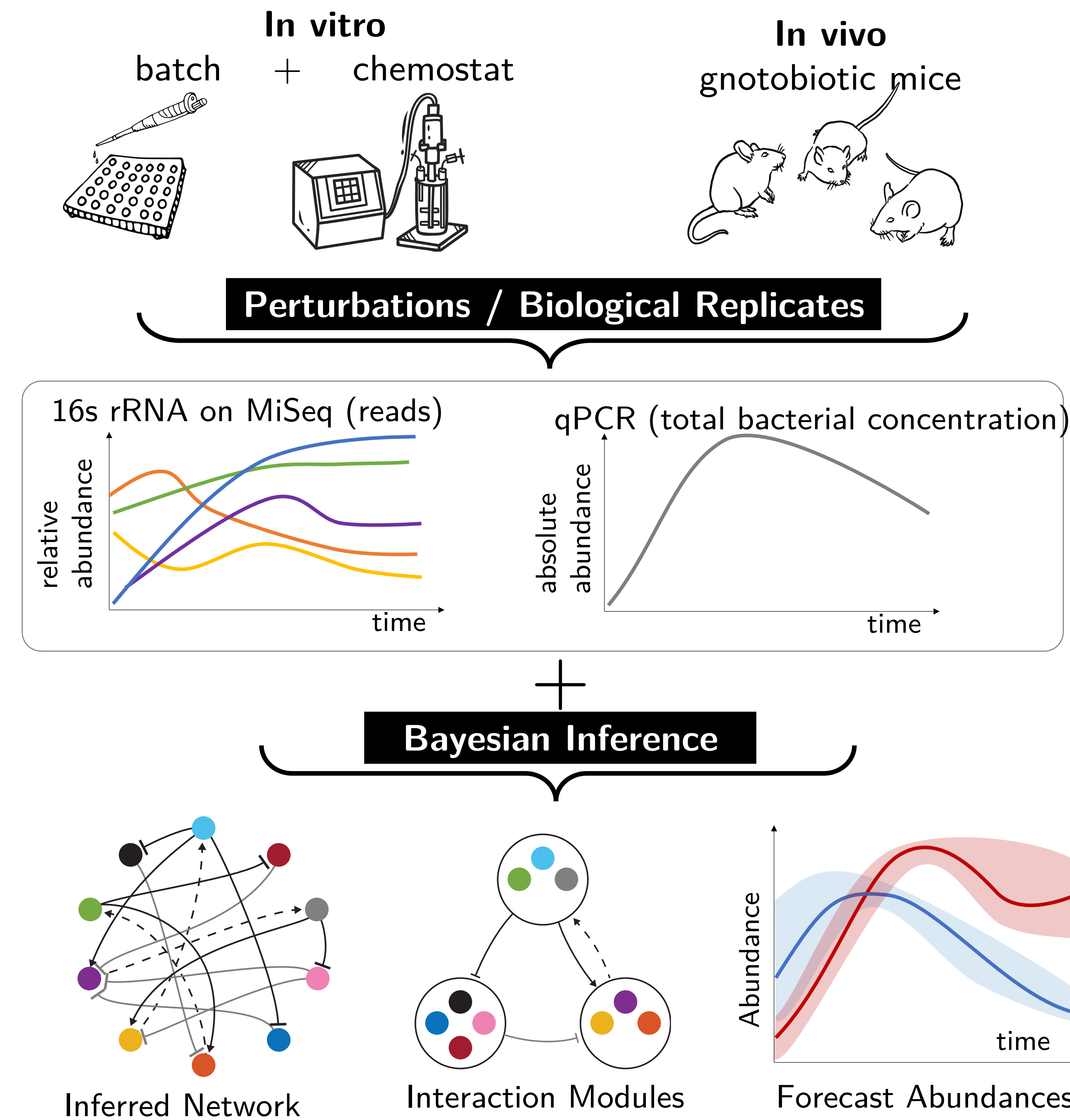Challenges associated with inference in the microbiome:

- High-dimensional (300+ strains of bacteria in the gut, **potentially 100,000 microbe-microbe interactions**).
- Temporally sparse and non-uniformly sampled data.
- High measurement noise
- Nonlinear and physically non-negative dynamics.

**Many potential interactions: how do we decide what interactions are "real" and how do we simplify the interaction landscape so as to be interpretable?**

**Contribution:**
1. Introduction of a temporally varying **auxiliary variable** technique to enable efficient inference by relaxing hard non-negativity constraint
2. Clustering of microbes into redundant **interaction modules** and **structural learning** of a compact interaction network among modules.

## Experimental Design for Inference

**In vitro**
batch + chemostat

**In vivo**
gnotobiotic mice

**Perturbations / Biological Replicates**

16s rRNA on MiSeq (reads)
qPCR (total bacterial concentration)

**Bayesian Inference**

Inferred Network   Interaction Modules   Forecast Abundances

## Dynamical Model

**Lotka-Volterra Dynamics**    Abundance of Microbe $i$ at time $t$: $\mathbf{x}_{t,i}$

$$\frac{d\mathbf{x}_{t,i}}{dt} = \boldsymbol{\alpha}_i \mathbf{x}_{t,i} + \boldsymbol{\beta}_{ii}\mathbf{x}_{t,i}^2 + \sum_{j\neq i}\boldsymbol{\beta}_{ij}\mathbf{x}_{t,i}\mathbf{x}_{t,j} + \frac{d\mathbf{w}_{t,i}}{dt}$$

growth parameter · interaction coefficient · self limiting · disturbance term

**Discrete approximation to the Lotka-Volterra dynamics**

$$\mathbf{x}_{k+1,i} = \mathbf{x}_{k,i} + \left(\boldsymbol{\alpha}_i\mathbf{x}_{k,i} + \boldsymbol{\beta}_{ii}\mathbf{x}_{k,i}^2 + \sum_{j\neq i}\boldsymbol{\beta}_{ij}\mathbf{x}_{k,i}\mathbf{x}_{k,j}\right)\Delta_k + (\mathbf{w}_{k+1,i} - \mathbf{w}_{k,i})\sqrt{\Delta_k}$$

$k$: discrete time index · discrete time step size

- Cluster the microbes into interaction modules (**Dirichlet Process** prior)
- No interactions **within** cluster, only **between** clusters.
- **Edge Selection**: add indicator variables for cluster-cluster interactions (Bayesian variable selection, structure learning for graphical models)
- Insert auxiliary variable **q** between the measurements **y**, **Q** and the state **x**

**Comment:** not allowing within cluster interactions dramatically reduces the number of inferred interactions and is consistent with the biologically important scenario of redundant functionality among sets of microbes

## Additional Model Components

**Dirichlet Process (clustering)**

cluster assignment
$$\boldsymbol{\pi}_{\mathbf{c}} \mid \alpha \sim \text{Stick}(\alpha)$$
$$\mathbf{c}_i \mid \boldsymbol{\pi}_{\mathbf{c}} \sim \text{Multinomial}(\boldsymbol{\pi}_{\mathbf{c}})$$
$$\mathbf{b}_{\mathbf{c}_i,\mathbf{c}_j} \mid \sigma_{\mathbf{b}} \sim \text{Normal}(0,\sigma_{\mathbf{b}}^2)$$
cluster interactions coefficient

**Edge Selection**
$$\mathbf{z}_{\mathbf{c}_i,\mathbf{c}_j} \mid \boldsymbol{\pi}_{\mathbf{z}} \sim \text{Bernouli}(\boldsymbol{\pi}_{\mathbf{z}})$$
indicator variable for cluster interaction

**Self Interaction**
$$\mathbf{a}_{i,1}, \mathbf{a}_{i,2} \mid \sigma_{\mathbf{a}} \sim \text{Normal}(0,\sigma_{\mathbf{a}}^2)$$

**Dynamics**
$$\mathbf{x}_{k+1,i} \mid \mathbf{x}_k, \mathbf{a}_i, \mathbf{b}, \mathbf{c}, \mathbf{z}, \sigma_{\mathbf{w}} \sim$$
$$\text{Normal}\left(\mathbf{x}_{k,i} + \mathbf{x}_{k,i}\left(\mathbf{a}_{i,1} + \mathbf{a}_{i,2}\mathbf{x}_{k,i} + \sum_{\mathbf{c}_j\neq\mathbf{c}_i}\mathbf{b}_{\mathbf{c}_i,\mathbf{c}_j}\mathbf{z}_{\mathbf{c}_i,\mathbf{c}_j}\mathbf{x}_{k,j}\right), \Delta_k\sigma_{\mathbf{w}}^2\right)$$
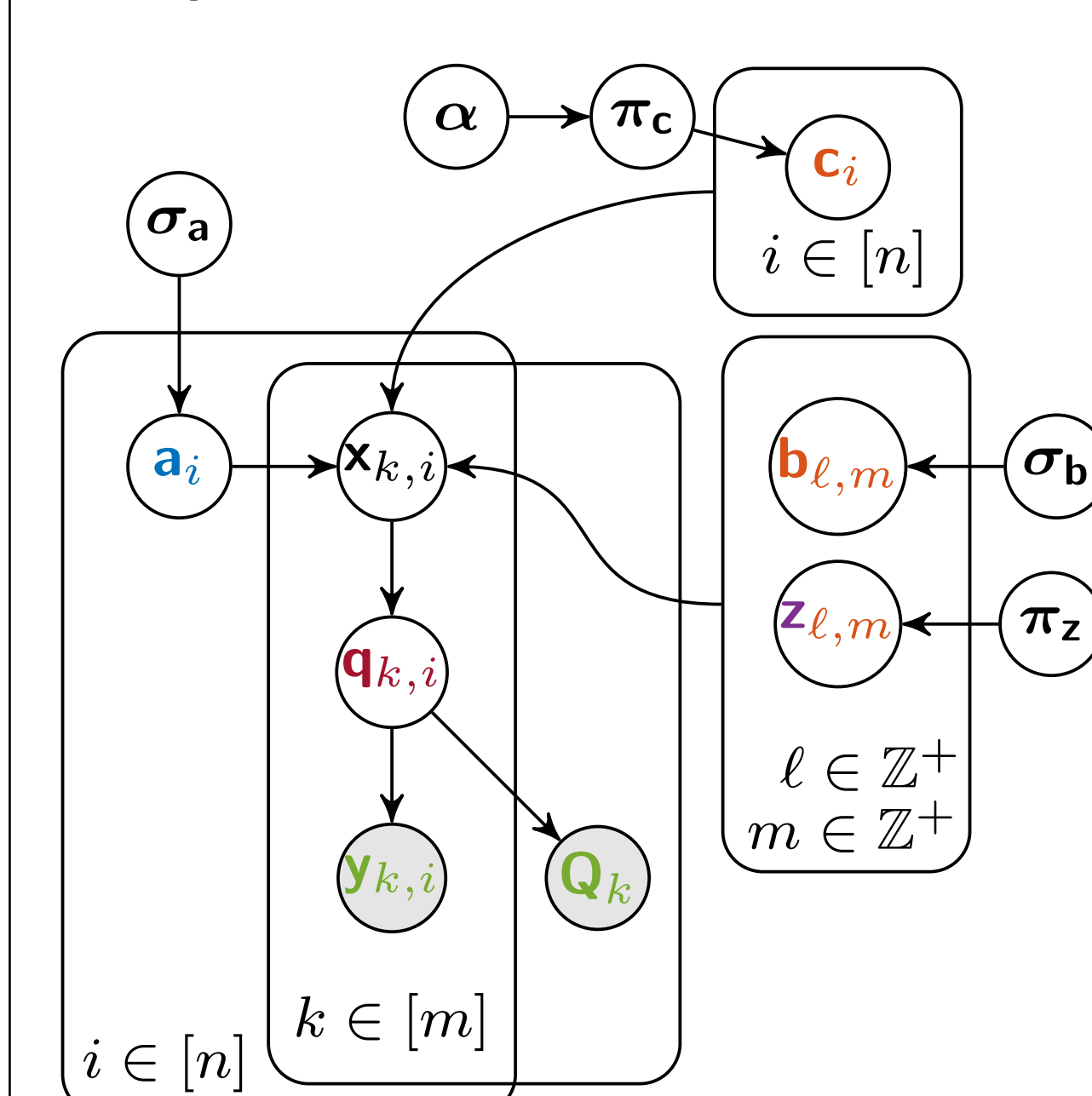
**Measurement Model**
auxiliary variable $\quad \mathbf{q}_{k,i} \mid \mathbf{x}_{k,i} \sim \text{Normal}(\mathbf{x}_{k,i}, \sigma_{\mathbf{q}}^2) \qquad \mathbf{q}_{k,i} \sim \text{Uniform}[0,\mathbf{L}]$

reads $\quad \mathbf{y}_{k,i} \mid \mathbf{q}_{k,i} \sim \text{NegBin}(\phi(\mathbf{q}_k), \epsilon(\mathbf{q}_k))$

qPCR $\quad \mathbf{Q}_k \mid \mathbf{q}_{k,i} \sim \text{Normal}\left(\sum_i \mathbf{q}_{k,i}, \sigma_{\mathbf{Q}_k}^2\right)$

Introduction of **q** allows for
- Efficient Gibbs/collapsed Gibbs sampling
- Posterior distributions for coefficients **a**, **b** are Gaussian, (direct sampling from posterior)
- Can Marginalize out in closed form the interaction coefficients **b** for both module learning and structure learning

Without **q**, dynamics would have been modeled with truncated distribution, resulting in the posteriors of **a**, **b** being truncated as well and not allowing for marginalization elsewhere in the model.
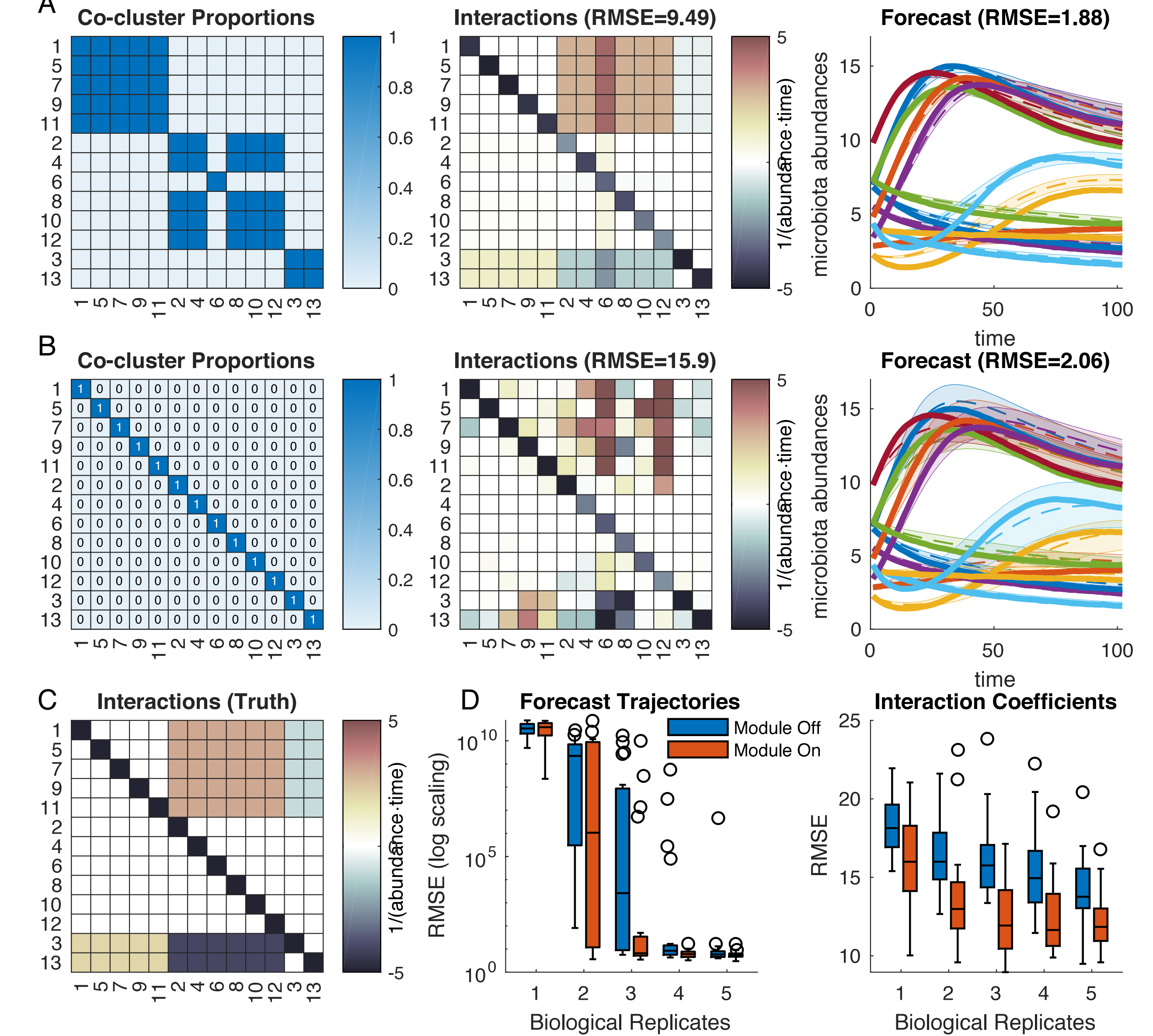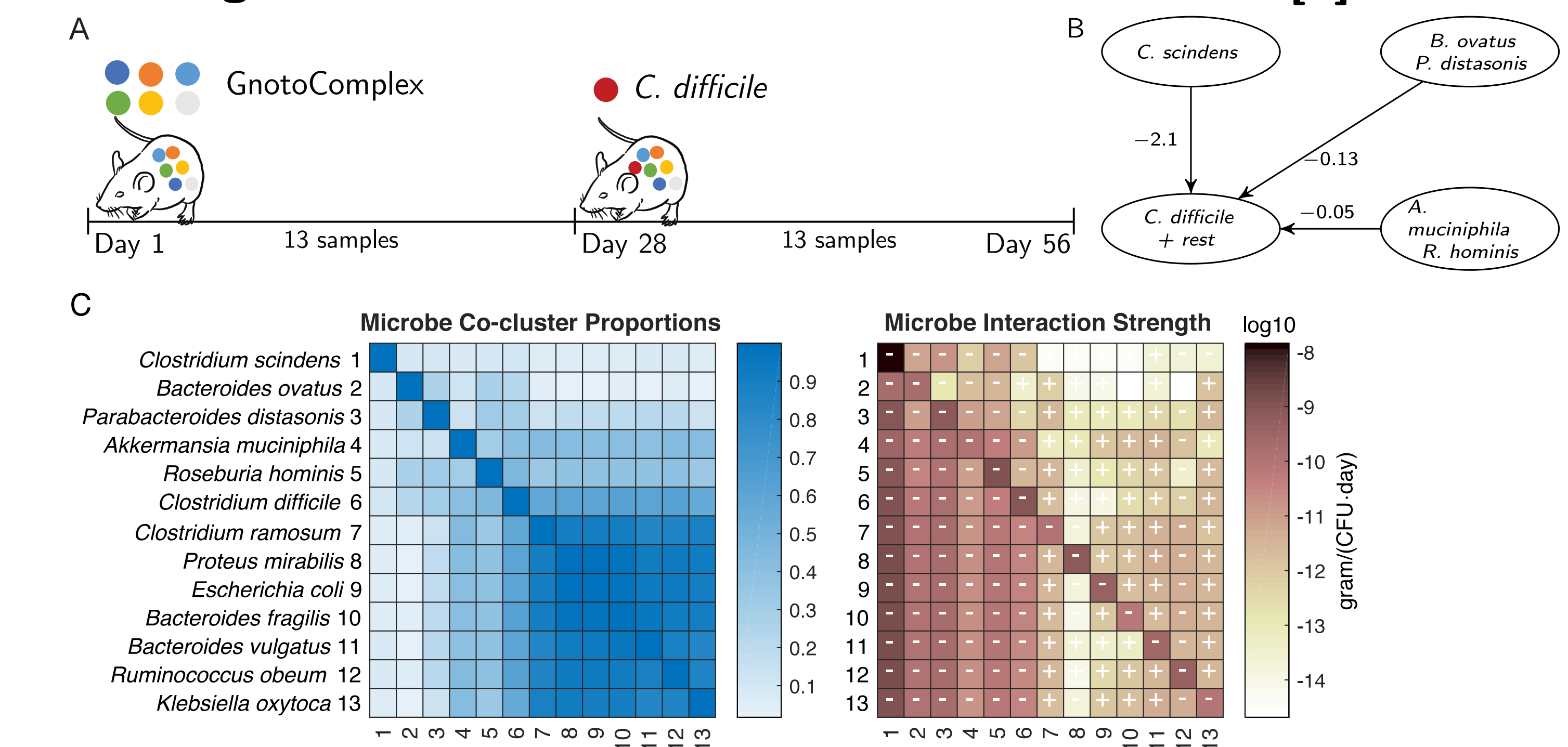
**Graphical Model**

[1] Bucci, Vanni, et al. "MDSINE: Microbial Dynamical Systems INference Engine for microbiome time-series analyses." *Genome Biology* (2016)

## Inference on Synthetic and In Vivo Time Series

**Synthetic data illustrating the use of module learning**



A Co-cluster Proportions | Interactions (RMSE=9.49) | Forecast (RMSE=1.88)
B Co-cluster Proportions | Interactions (RMSE=15.9) | Forecast (RMSE=2.06)
C Interactions (Truth) | D Forecast Trajectories | Interaction Coefficients

**Learning interaction modules from *in vivo* time series [1]**



A GnotoComplex — C. difficile — Day 1 · 13 samples · Day 28 · 13 samples · Day 56

C Microbe Co-cluster Proportions | Microbe Interaction Strength

1 Clostridium scindens
2 Bacteroides ovatus
3 Parabacteroides distasonis
4 Akkermansia muciniphila
5 Roseburia hominis
6 Clostridium difficile
7 Clostridium ramosum
8 Proteus mirabilis
9 Escherichia coli
10 Bacteroides fragilis
11 Bacteroides vulgatus
12 Ruminococcus obeum
13 Klebsiella oxytoca

## Ongoing Projects

**Microbial Interaction Engineering**

Marika Ziesack
Silver Lab, Harvard

E. coli · B. fragilis · B. theta · S. typhimurium

- Engineered to overproduce 1 amino acid
- Engineered to need 3 amino acids

**Bacteriophage Therapy**

Bryan Hsu
Silver Lab, Harvard

○ uninfected bacteria
◑ infected bacteria
bacteriophages

- Bacteriophages are "bacteria viruses"
- Phages as control knobs to modulate microbiome

**Complex Microbiota**

Massachusetts Host Microbiome Center

transplant

- Colonize mice with human fecal sample containing 300+ bacteria