

Robust and Scalable Models of Microbiome Dynamics

Travis E. Gibson¹

tgibson@mit.edu

Georg K. Gerber^{1,2}

ggerber@bwh.harvard.edu

¹ Massachusetts Host Microbiome Center
Brigham and Women's Hospital & Harvard Medical School

²Harvard-MIT Health Sciences and Technology

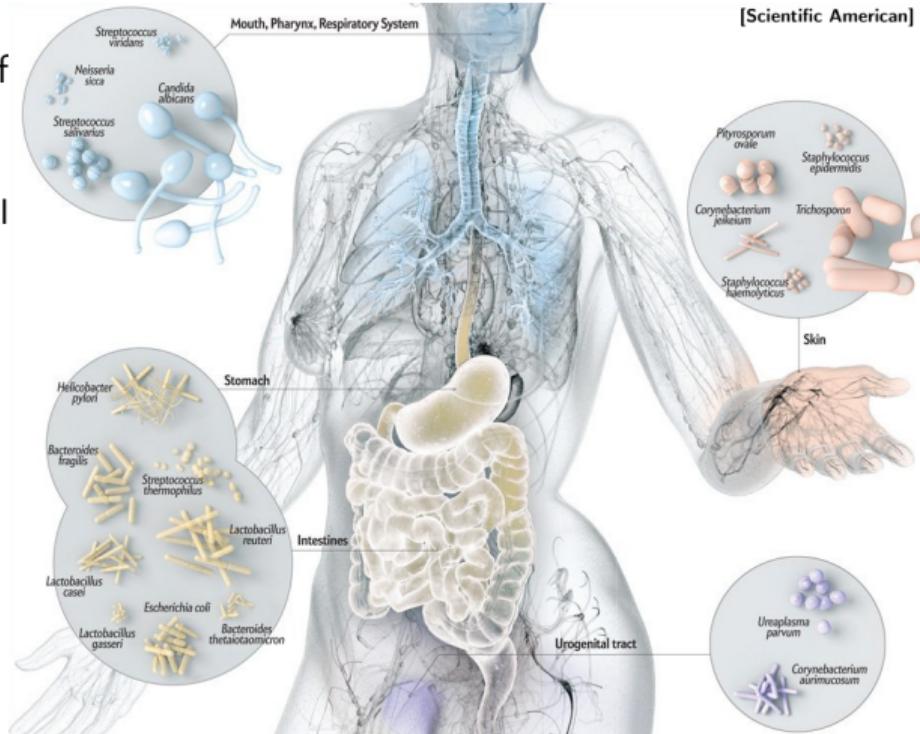
35th International Conference on Machine Learning
July 12th, 2018

Outline

- ① What is the microbiome
- ② From sequences to reads to microbial interactions
- ③ Bayesian nonparametric model for microbe dynamics

The Microbiome

- ① The **microbiome** is the aggregate of microorganisms that resides on or within any of a number of human tissues and biofluids:
 - skin, mammary glands, placenta, seminal fluid, uterus, ovarian follicles, lung, saliva, oral mucosa, conjunctiva, biliary and **gastrointestinal tracts**) [wikipedia]
- ② 10^{14} Microbes in/on your body [Sender et al. *PLoS Biology* 2016]
- ③ 3.3 million genes compared to 23,000 human genes [Qin et al. *Nature* 2010]
- ④ Play a role in a variety of human diseases:
 - infections, arthritis, food allergy, cancer, inflammatory bowel disease, neurological diseases, and obesity/diabetes



Bacteriotherapy

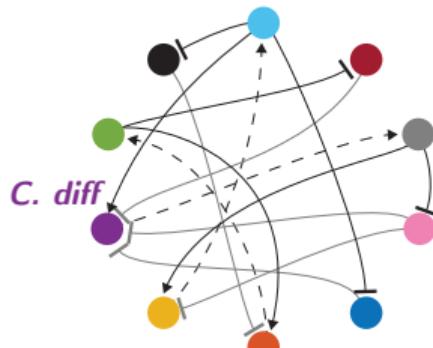
Bacteriotherapy: communities of bacteria administered to patients for specific therapeutic applications

- “bugs-as-drugs”

Clostridium difficile infection

- Causes serious diarrhea (14K deaths/yr)
- Antibiotics disrupt helpful bacteria in gut
- Increasingly difficult to treat with conventional therapies (more antibiotics): 20-30% recurrence rate

Pharmacology meets Ecology



microbial interaction network

positive microbe A produces a small molecule (metabolite) that microbe B needs

negative two microbes competing for the same niche

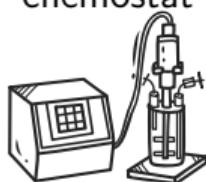
what if there were 300 bugs in the network?

Workflow in our lab

batch
experiments



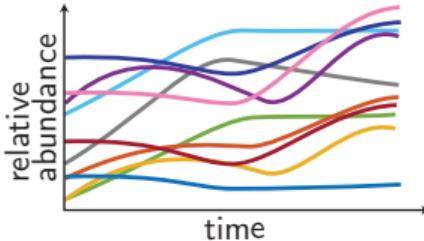
chemostat



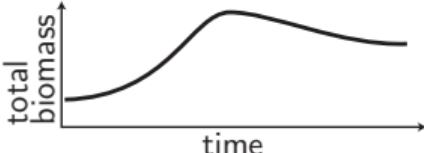
animal
experiments



- 16S rRNA on MiSeq (reads) for relative abundances of species

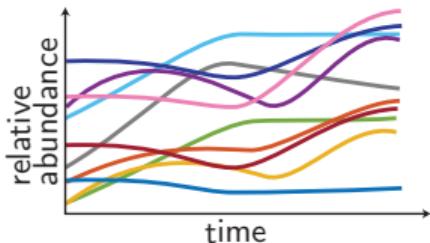


- 16S rRNA qPCR (universal primers) for bacterial biomass



Sequencing to obtain microbe relative abundances

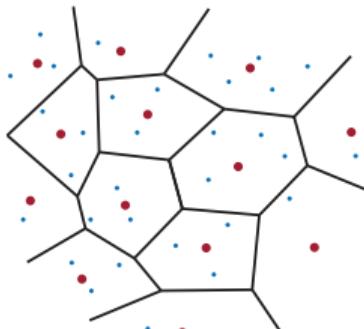
- 16S rRNA on MiSeq (reads) for relative abundances of species



1. Fastq file

```
@sequence-id  
TCGCACCTCAACGCCCTGCATATGACAAGACAGAATC .....nucleobase sequence  
+  
<> ;##=><9=AAAAAAA9#:<#<;<<<?????#= .....quality scores
```

2. Sequences are clustered



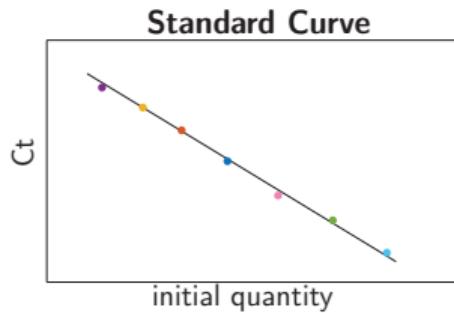
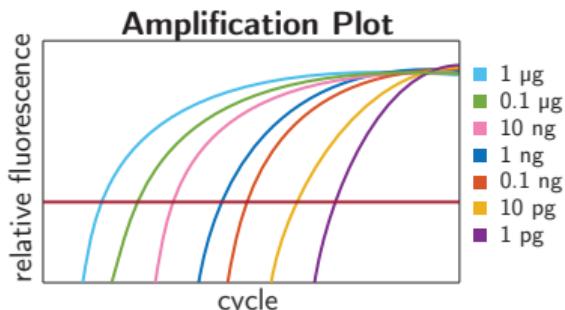
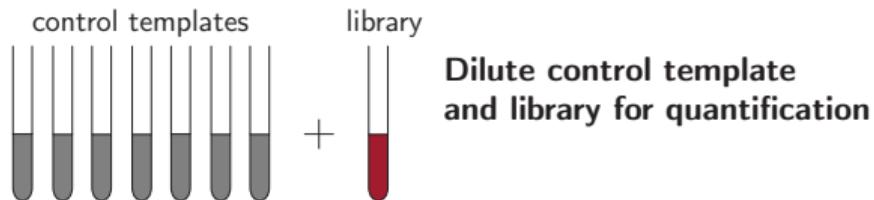
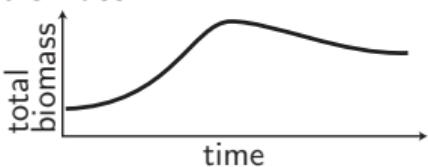
3. Read Table

	sample1	sample2	...
bacteria1	11	1004	
bacteria2	0	7	
bacteria3	301	275	
⋮	⋮	⋮	⋮

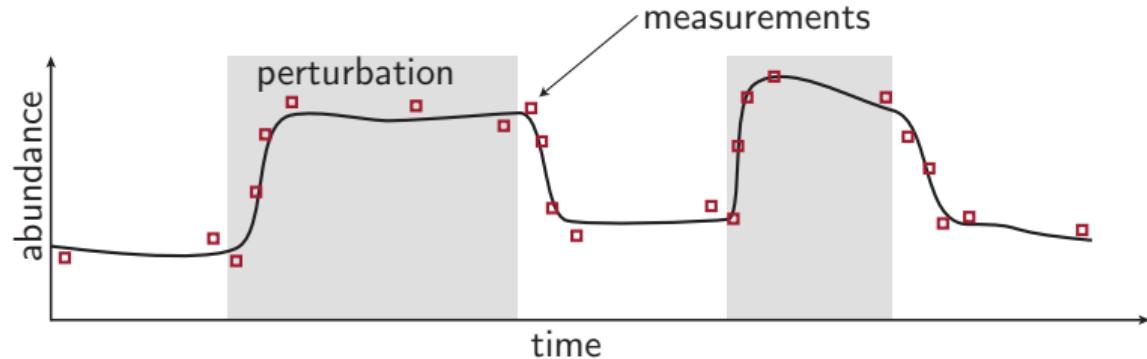
Reads ~ Negative Binomial

Quantitative PCR for total bacterial biomass

- 16S rRNA qPCR (universal primers) for bacterial biomass



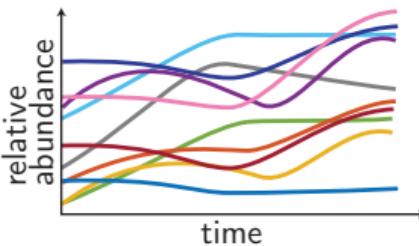
Irregular and sparse measurements



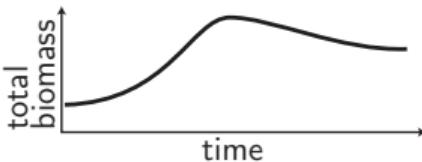
- **measurements - irregular,
sparse & noisy**

Learning microbial interaction networks

- 16S rRNA on MiSeq (reads) for relative abundances of species

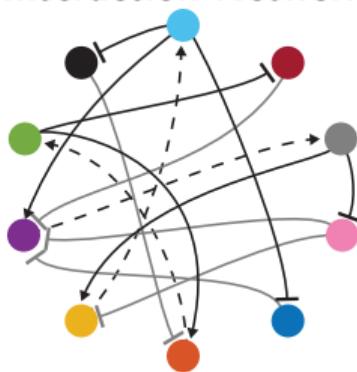


- 16S rRNA qPCR (universal primers) for bacterial biomass



- measurements - irregular, sparse & noisy

Interaction Network



Abundance of microbe i at time t : $\mathbf{x}_{t,i}$

$$\frac{d\mathbf{x}_{t,i}}{dt} = \alpha_i \mathbf{x}_{t,i} + \beta_{ii} \mathbf{x}_{t,i}^2 + \sum_{j \neq i} \beta_{ij} \mathbf{x}_{t,i} \mathbf{x}_{t,j}$$

growth, self limiting, interaction

Previous literature specific to the microbiome

- Log transform dynamics → Linear Regression + L2 [Stein et al. *PLoS Comput Biology* 2013]
- Sparse linear regression with bootstrap aggregation [Fisher et al. *PLoS One* 2014]
- Bayesian model with deterministic dynamics (independent filtering) [Bucci et al. *Genome Biology* 2016]
- Extended Kalman Filter [Alshawaqfeh et al. *BMC Genomics* 2017]

Goal with our model and short literature review

Three main contributions in our model

① Clustering (interaction modules)

- Dirichlet Process (DP)

[Rasmussen, *Advances in Information Processing Systems 2000*]

[Neal, *Journal of computational and graphical statistics, 2000*]

② Edge selection (structure learning, variable selection)

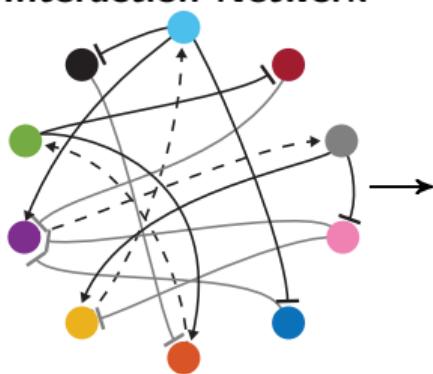
- Bayesian Networks

[George and McCulloch, *Journal of the ASA, 1993*]

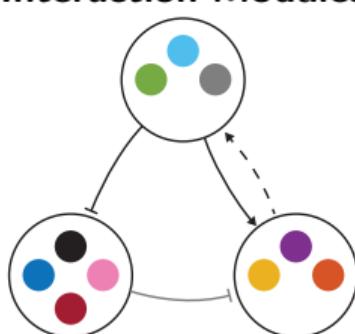
[Heckerman, *A Tutorial on Learning with Bayesian Networks, 2008*]

③ Introduction of an auxiliary variable between the measurement model and latent state

Interaction Network



Interaction Modules



- 300 species
- 90,000 interactions
- Redundant gene function

Back to the basic time series model

- Abundance of microbe i at time t : $\mathbf{x}_{t,i}$

$$\frac{d\mathbf{x}_{t,i}}{dt} = \alpha_i \mathbf{x}_{t,i} + \beta_{ii} \mathbf{x}_{t,i}^2 + \sum_{j \neq i} \beta_{ij} \mathbf{x}_{t,i} \mathbf{x}_{t,j} + \frac{dw_{t,i}}{dt}$$

growth, self limiting, interaction, stochastic disturbance

- Convert to discrete time

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

discrete time step size

Next we discuss the three main ingredients to our model

- ① Clustering (interaction modules)
- ② Edge selection (structure learning, variable selection)
- ③ Introduction of an auxiliary variable between the measurement model

Complete Model

Dirichlet Process

$$\pi_c \mid \alpha \sim \text{Stick}(\alpha)$$

$$c_i \mid \pi_c \sim \text{Multinomial}(\pi_c)$$

$$b_{c_i, c_j} \mid \sigma_b \sim \text{Normal}(0, \sigma_b^2)$$

Dynamics

$$x_{k+1,i} \mid x_k, a_i, b, c, z, \sigma_w \sim$$

$$\text{Normal}\left(x_{k,i} + x_{k,i} \left(a_{i,1} + a_{i,2} x_{k,i} + \sum_{c_j \neq c_i} b_{c_i, c_j} z_{c_i, c_j} x_{k,j} \right), \Delta_k \sigma_w^2 \right)$$

Constraint and Measurement Model

aux $q_{k,i} \mid x_{k,i} \sim \text{Normal}(x_{k,i}, \sigma_q^2)$

reads $y_{k,i} \mid q_{k,i} \sim \text{NegBin}(\phi(q_k), \epsilon(q_k))$

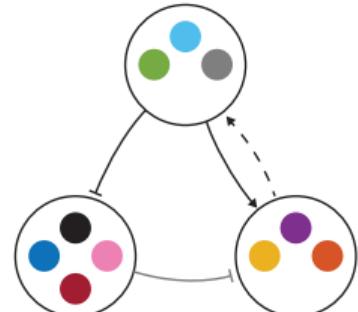
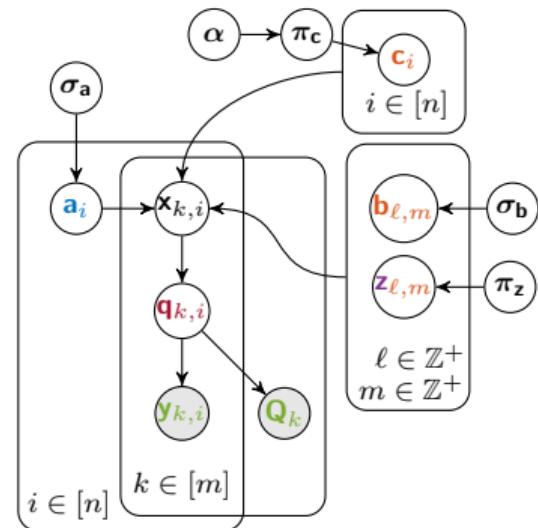
qPCR $Q_k \mid q_{k,i} \sim \text{Normal} \left(\sum_i q_{k,i}, \sigma_{Q_k}^2 \right)$

Edge Selection (Structure)

$$z_{c_i, c_j} \mid \pi_z \sim \text{Bernoulli}(\pi_z)$$

Self Interactions

$$a_{i,1}, a_{i,2} \mid \sigma_a \sim \text{Normal}(0, \sigma_a^2)$$



Simplified model unraveled in time - auxiliary variable

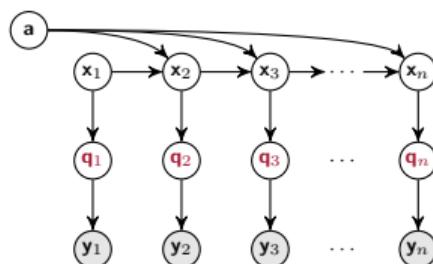
$$x_{t+1,i} | x_t, a \sim \text{Normal}(a_i^T f(x_t), \sigma_{x_i}^2)$$

$$q_{k,i} | x_{k,i} \sim \text{Normal}(x_{k,i}, \sigma_q^2)$$

$$q_{k,i} \sim \text{Uniform}[0, L]$$

$$y_{k,i} | \sigma_y, q_{k,i} \sim \text{Normal}_{\geq 0}(q_{k,i}, \sigma_y^2)$$

$$a_i \sim \text{Normal}(0, \sigma_{a_i}^2)$$



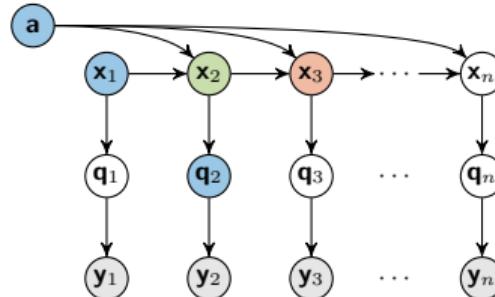
Prior on q is positive, relaxing the distribution on the dynamics for x

Parameter inference Gibbs step: $a^{(g+1)} \sim p_{a|x}(\cdot | x^{(g)})$

- Direct sampling from the posterior possible (Bayesian Regression!)

Sampling for other variables

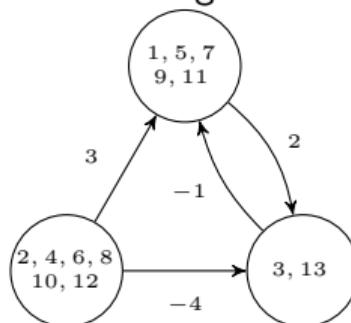
- Collapsed Gibbs sampling for Dirichlet Process and Edge Selection (integrate out a)
- Filtering is still challenging but easy to design proposals for (MH)



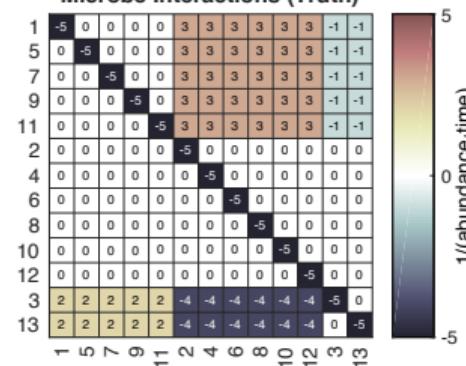
Synthetic experiment

- Comparing inference with and without clustering enabled

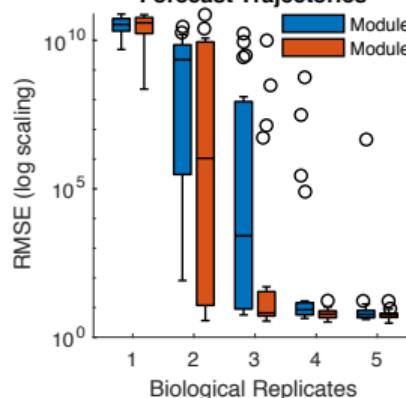
- Ground truth interaction network



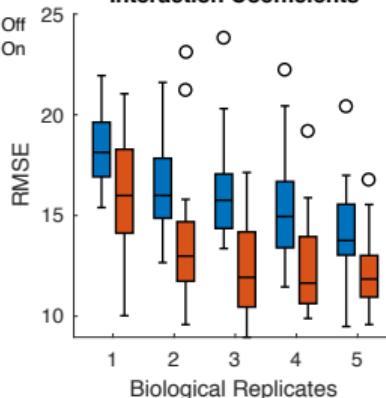
C Microbe Interactions (Truth)



D Forecast Trajectories

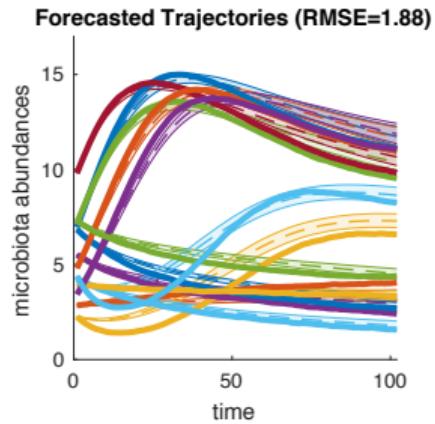
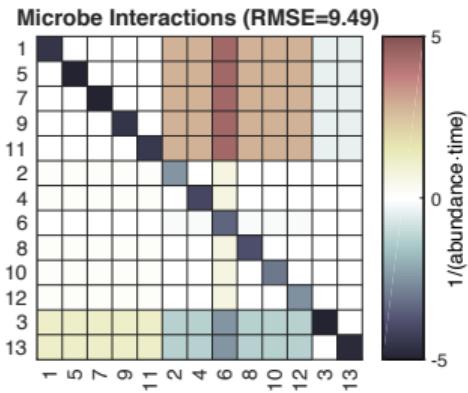
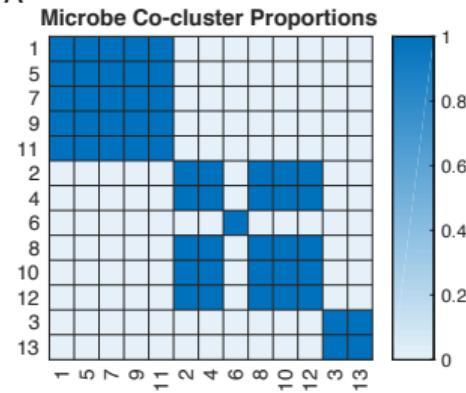


Interaction Coefficients

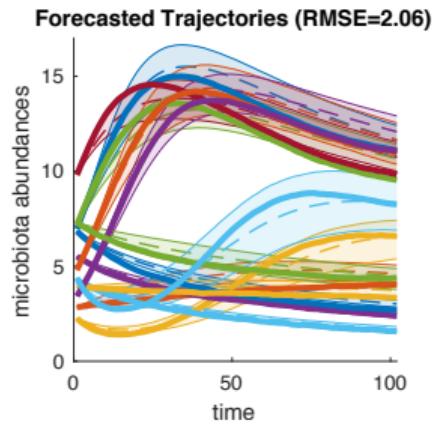
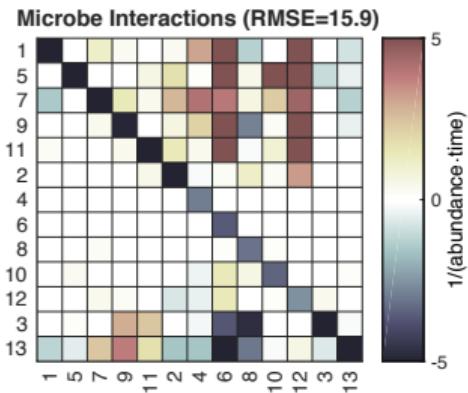
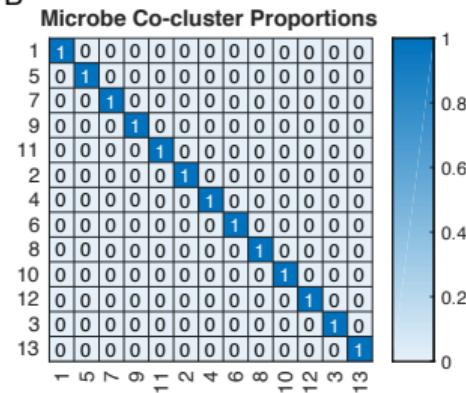


Synthetic experiment continued

A

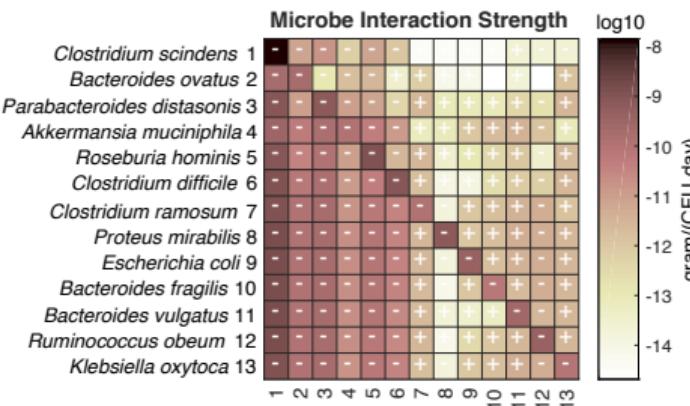
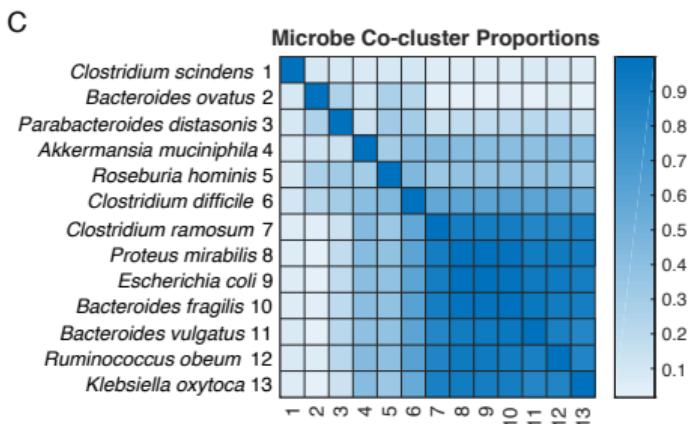
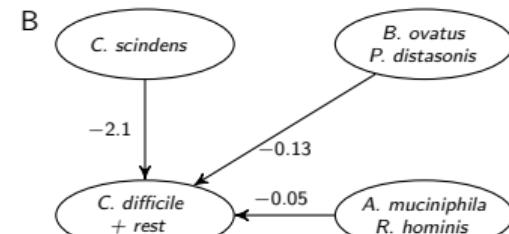
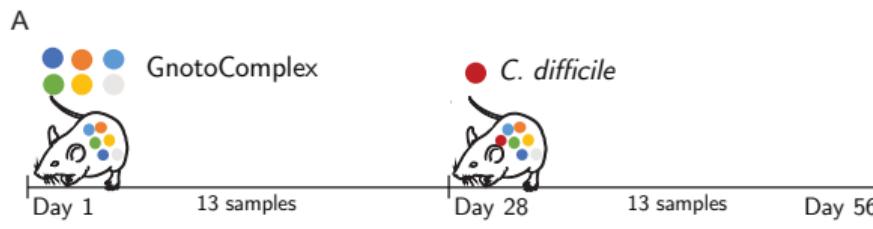


B



In vivo infection experiment

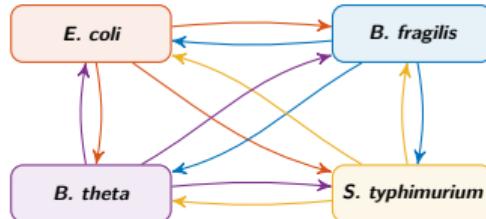
- Data from *C. difficile* challenge experiment performed with 5 germ free mice [Bucci et al. *Genome Biology* 2016]



Three ongoing projects using the model



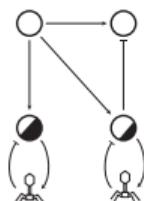
Marika Ziesack
Silver Lab, Harvard



- Microbes engineered to overproduces one amino acid
- Microbes engineered to need three amino acids
- Compare inference on WT and engineered strains to prove that engineering was performed.



Bryan Hsu
Silver Lab, Harvard



○ uninfected bacteria
● infected bacteria
✖ bacteriophages

- Bacteriophages are “bacteria viruses”
- Using phages as control knobs to modulate microbiome
- Unlike antibiotics, phages are host specific and can be present in equilibrium with host

Massachusetts
Host Microbiome
Center



- Transplant complex bacteria into germ free mice (300+ bacteria)

Conclusions

- ① Dynamical systems model for microbial dynamics based on what we term interaction modules (probabilistic clusters of latent variables with redundant interaction structure)
- ② Bayesian formulation of the stochastic dynamical systems model that propagates measurement and latent state uncertainty throughout the model
- ③ Introduction of a temporally varying auxiliary variable technique to enable efficient inference by relaxing the hard non-negativity constraint on states

Time series metagenomics is a growing area for ML applications

- Metagenomics is a “Big Data” problem, but number of time series samples is still sparse in comparison to the number of latent dynamical coefficients we need to learn.

Poster #2

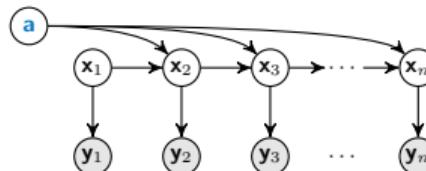
Backup Slides

Simple example without the intermediate auxiliary variable

$$\mathbf{x}_{t+1,i} \mid \mathbf{x}_t, \mathbf{a} \sim \text{Normal}_{\geq 0}(\mathbf{a}_i^\top f(\mathbf{x}_t), \sigma_{x_i}^2)$$

$$\mathbf{y}_{t,i} \mid \mathbf{x}_{t,i} \sim \text{Normal}_{\geq 0}(\mathbf{x}_{t,i}, \sigma_{y_i}^2)$$

$$\mathbf{a}_i \sim \text{Normal}(0, \sigma_{\mathbf{a}_i}^2)$$



Note the truncated distributions for \mathbf{x} and \mathbf{y}

Parameter inference Gibbs step: $\mathbf{a}^{(g+1)} \sim p_{\mathbf{a}|\mathbf{x}}(\cdot \mid \mathbf{x}^{(g)})$

$$p_{\mathbf{a}|\mathbf{x}} \propto p_{\mathbf{x}|\mathbf{a}} p_{\mathbf{a}}$$
$$\text{Normal}_{\geq 0}(\mathbf{x}; \mu(\mathbf{a}, \mathbf{x}), \sigma^2)$$
$$\text{Normal}(\mathbf{a}; 0, \sigma^2)$$
$$= \frac{e^{-\frac{1}{2\sigma^2}(\mathbf{x}-\mu(\mathbf{a}, \mathbf{x}))^2}}{\sigma\sqrt{2\pi}\left(\Phi(\infty) - \Phi\left(-\frac{\mu(\mathbf{a}, \mathbf{x})}{\sigma}\right)\right)} \frac{e^{-\frac{1}{2\sigma^2}\mathbf{a}^2}}{\sigma\sqrt{2\pi}}$$

Sampling for other variables

- Filtering (sampling from posterior of \mathbf{x}) is challenging
- Can not use collapsed Gibbs sampling for Dirichlet Process or Edge Selection

Negative Binomial

$$\mathbf{y}_{k,i} \mid \mathbf{q}_k \sim \text{NegBin}(\phi(\mathbf{q}_k, r_k), \epsilon(\mathbf{q}_k, a_0, a_1))$$

$$\phi(q_k, r_k) = r_k \frac{q_{k,i}}{\sum_i q_{k,i}} \tag{1}$$

$$\epsilon(q_k, a_0, a_1) = \frac{a_0}{q_{k,i}/\sum_i q_{k,i}} + a_1 \tag{2}$$

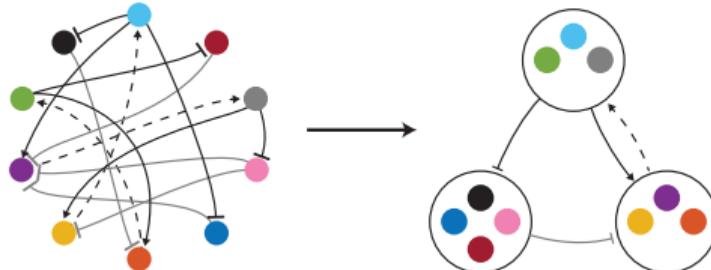
$$\begin{aligned} \text{NegBin}(y; \phi, \epsilon) &= \frac{\Gamma(r+y)}{y! \Gamma(r)} \left(\frac{\phi}{r+\phi} \right)^y \left(\frac{r}{r+\phi} \right)^r \\ r &= \frac{1}{\epsilon} \end{aligned}$$

Priors

- We use conjugate priors on many variables (e.g., the variance terms ($\sigma_a^2, \sigma_b^2, \sigma_w^2$) have Inverse-Chi-squared priors).
- The module assignments, \mathbf{c} , are also updated by a standard Gibbs sampling approach for Dirichlet Processes.
- For the concentration parameter α , which has a Gamma prior on α , we use the sampling method described by Escobar and West 1995.

Backup: Alternative View of Bacteriotherapy

Interpretable inference for bacteriotherapy design



Predict microbial abundances (phenotype) in the presence of the bacteriotherapy

