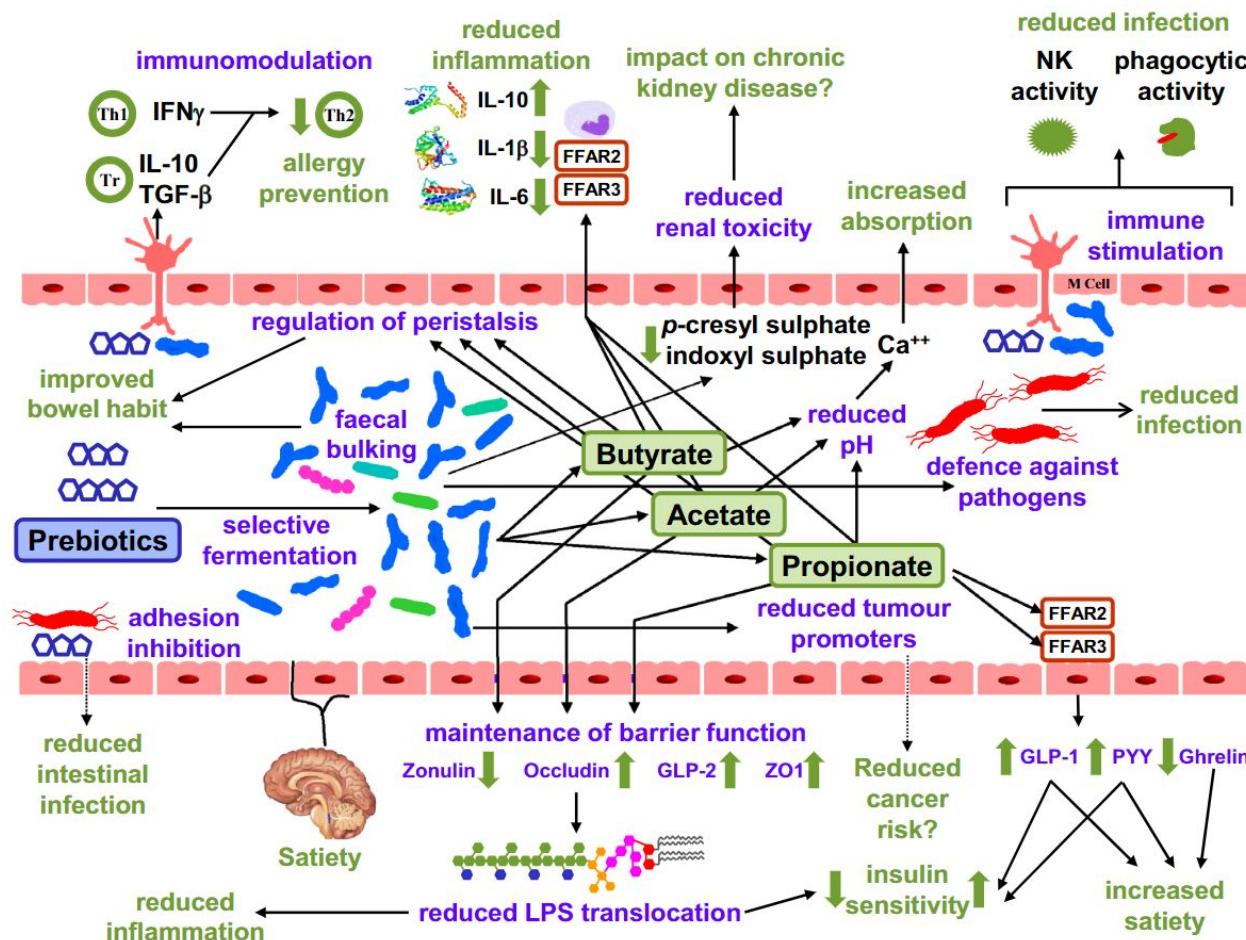


# Characterizing sources of variation in bacterial sequence count data

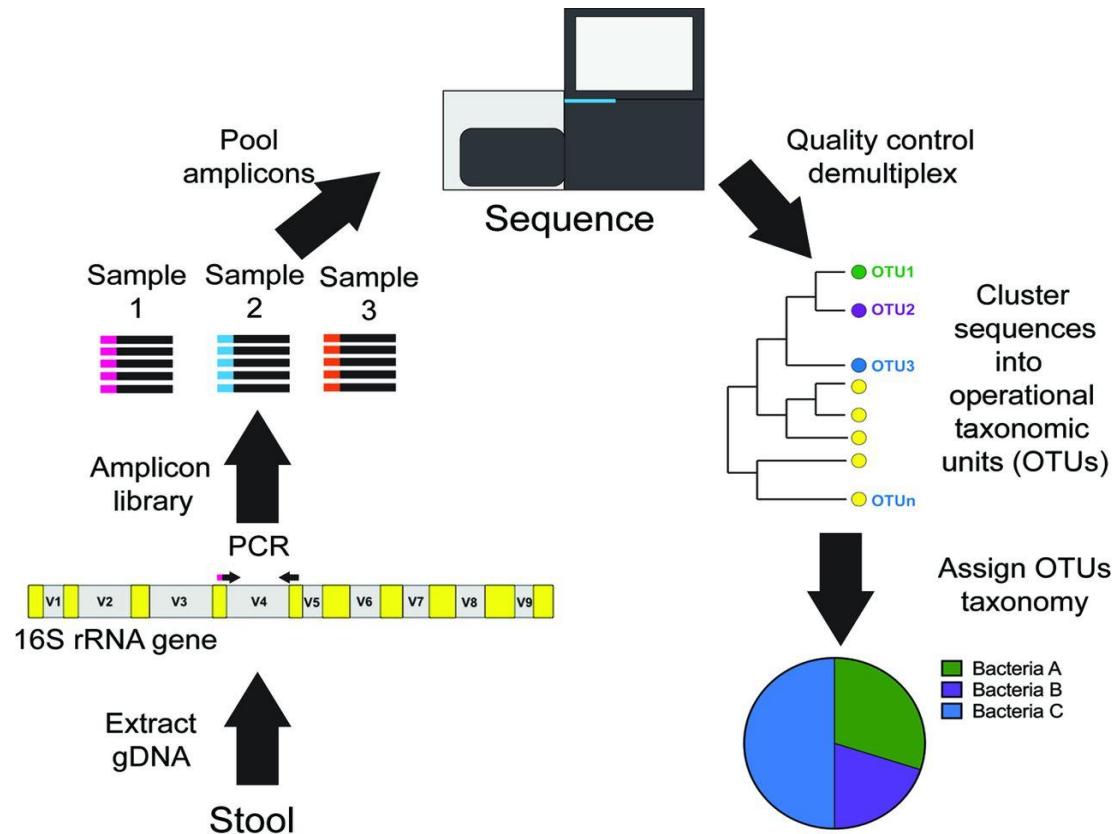
Nov. 13, 2019

The human gut microbiome plays roles in metabolism and host immune response.



Mohajeri et al. "The role of the microbiome for human health: from basic science to clinical applications." *European Journal of Nutrition* (2018)

We can profile the composition of the gut microbiome by 16S rRNA sequencing.

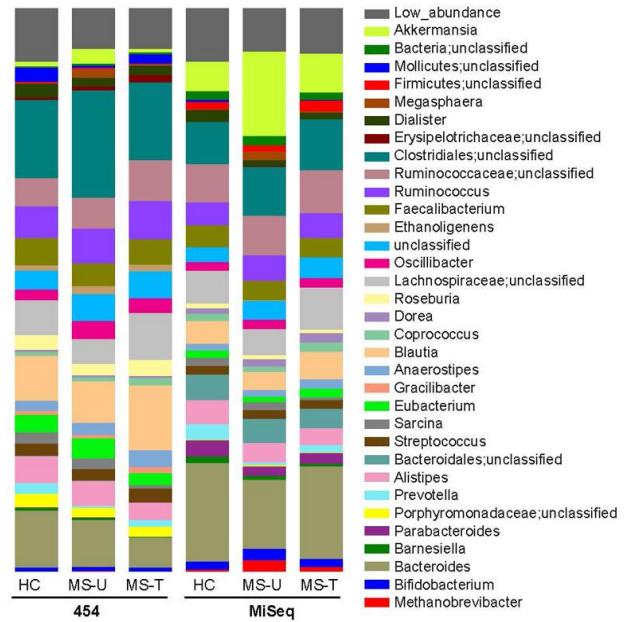


Koh, A. "Potential for Monitoring Gut Microbiota for Diagnosing Infections and Graft-versus-Host Disease in Cancer and Stem Cell Transplant Patients." *Clinical Chemistry* (2017).

The data are counts...

	Phylum 1	Phylum 2	Phylum 3	...	Phylum 30	Phylum 31	Phylum 32	...	Phylum 48	Phylum 49	Phylum 50
<b>Sample 1</b>	9351	18740	2548		14	47	0		29	15	0
<b>Sample 2</b>	28210	5029	449		0	9	106		70	0	0
<b>Sample 3</b>	31591	10762	559		40	191	13		0	0	0
<b>Sample 4</b>	6542	15030	4259		4	20	0		17	0	0
<b>Sample 5</b>	8333	4367	1157		7	435	208		9	0	0
<b>Sample 6</b>	12278	1426	379		0	16	108		77	0	0
<b>Sample 7</b>	7759	311	144		147	245	55		0	0	0
<b>Sample 8</b>	20609	56	4		101	89	13		0	0	0
<b>Sample 9</b>	28193	19442	321		0	5	50		0	0	0
<b>Sample 10</b>	16120	6318	257		59	28	90		5	0	0
<b>Sample 11</b>	552	98	951		147	0	9		52	0	55
<b>Sample 12</b>	1020	13553	1905		23	239	0		6	63	0
<b>Sample 13</b>	3163	610	1646		631	31	8		0	0	0
<b>Sample 14</b>	30054	909	55		39	1049	782		0	0	0
<b>Sample 15</b>	4122	94	263		38	0	102		14	0	0
<b>Sample 16</b>	12655	9327	385		3	105	0		0	0	0
<b>Sample 17</b>	26461	152	4		0	0	0		0	0	0
<b>Sample 18</b>	7707	7102	332		3	37	7		10	39	0
<b>Sample 19</b>	7181	7701	183		634	14	91		156	0	0
<b>Sample 20</b>	17341	3686	115		14	108	23		0	0	0

...that transmit proportional information.



**HC** healthy control

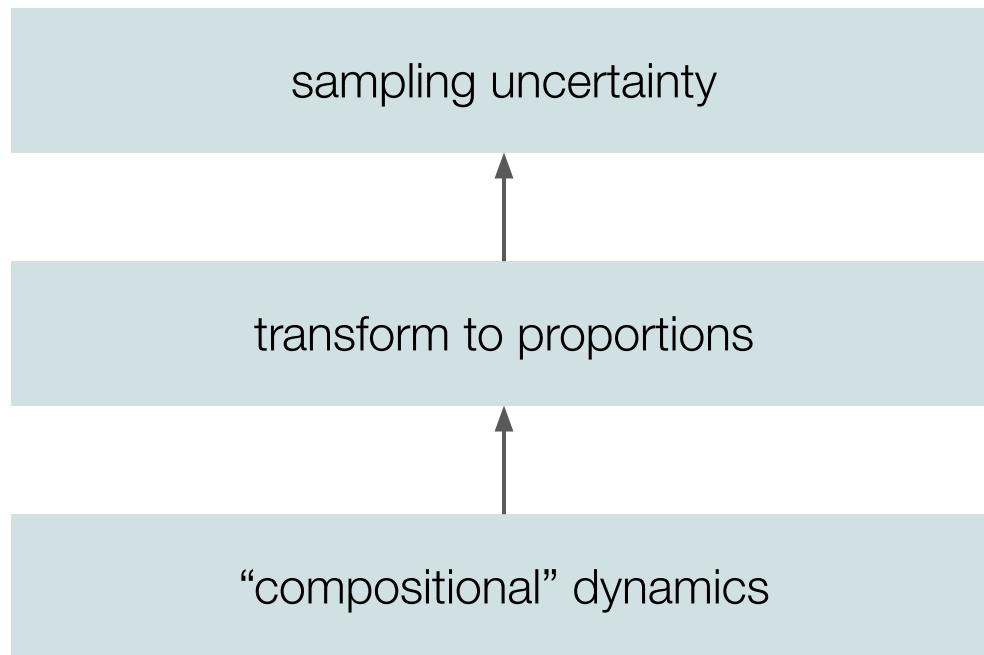
**MS-U** multiple sclerosis (untreated)

**MS-T** multiple sclerosis (treated)

Jangi et al. "Alterations of the human gut microbiome in multiple sclerosis." Nature Communications (2016)

We're interested in biological and technical variation between samples.

Count and proportionality considerations motivate a multinomial logistic normal model.



# Aims

- A1** Performing batch effect correction in microbial sequence count data
- A2** Characterizing population variation in gut microbial and functional “dynamics”
- A3** Testing association of dynamics and host fitness

# Aims

- A1** Performing batch effect correction in microbial sequence count data
- A2** Characterizing population variation in gut microbial and functional “dynamics”
- A3** Testing association of dynamics and host fitness

# 16S batch effect correction methods repurpose approaches from microarray and RNA-seq data.

	Compositional representation?	Model sampling uncertainty?	Model batch?
<b>Linear models</b>			
limma (Ritchie et al, 2015)	✓		✓
ComBat (Leek and Storey, 2007)	✓		✓
SVA (Leek et al., 2012)	✓		✓
<b>Poisson/NB GLMs</b>			
edgeR (Robinson et al., 2010)		✓	✓
DESeq2 (Love et al., 2014)		✓	✓
<b>Differential expression in the microbiome</b>			
ALDEx2 (Fernandes et al., 2014)	✓	✓	
BDMMA (Dai et al., 2018)	✓	✓	✓
McLaren et al., 2019	✓		✓
<b>Non-parametric</b>			
Percentile normalization (Gibbons et al., 2016)			✓

- 1 Enhance a classic approach by modeling sampling uncertainty & framing correction within compositional data POV

$$Y_{ijg} = \alpha_g + X\beta_g + \gamma_{ig} + \delta_{ig}\varepsilon_{ijg}$$

sample  
batch  
gene

baseline gene expression

covariate effects

batch effect

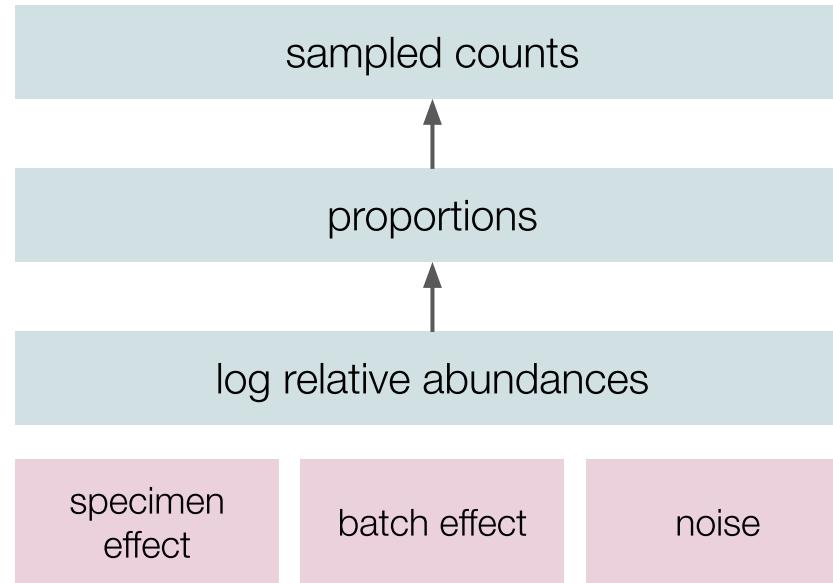
batch scatter x baseline scatter

The ComBat model

Johnson, et al. "Adjusting batch effects in microarray expression using empirical Bayes methods." Biostatistics (2007)

- 2 Introduce a measure for quantifying goodness-of-correction
- 3 Evaluate performance on simulated and real data sets

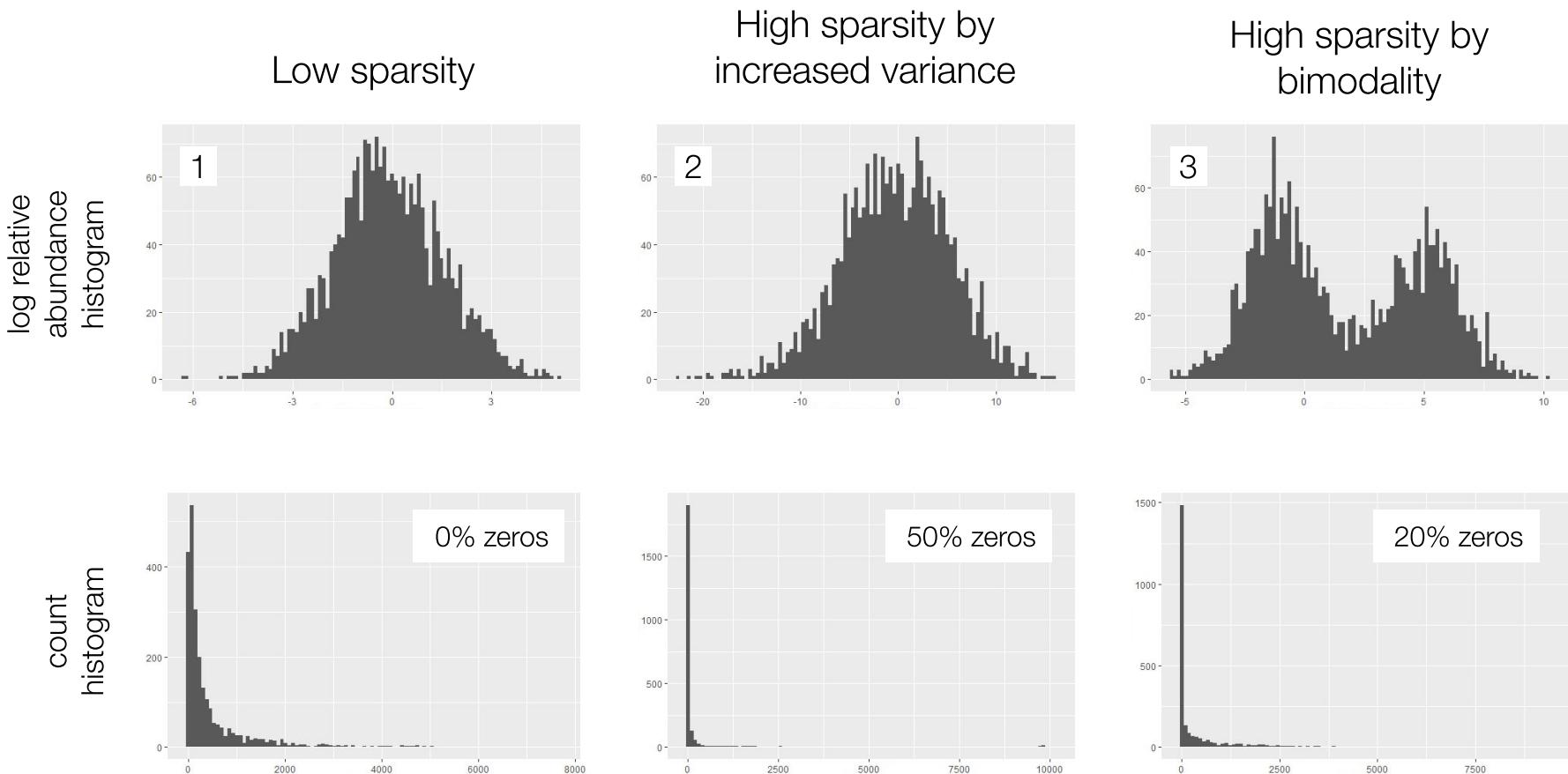
# A simulation model in outline



# A relative scoring scheme

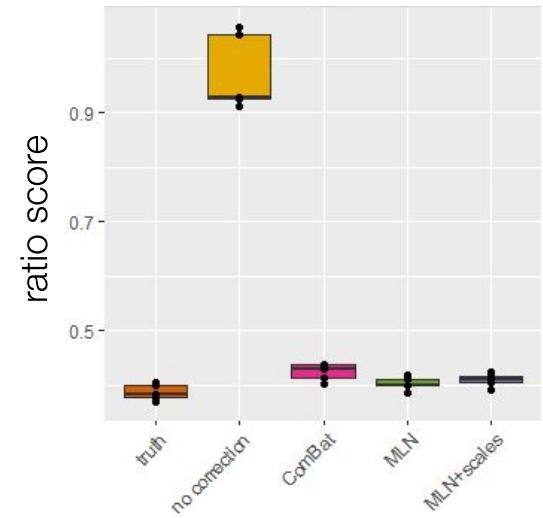
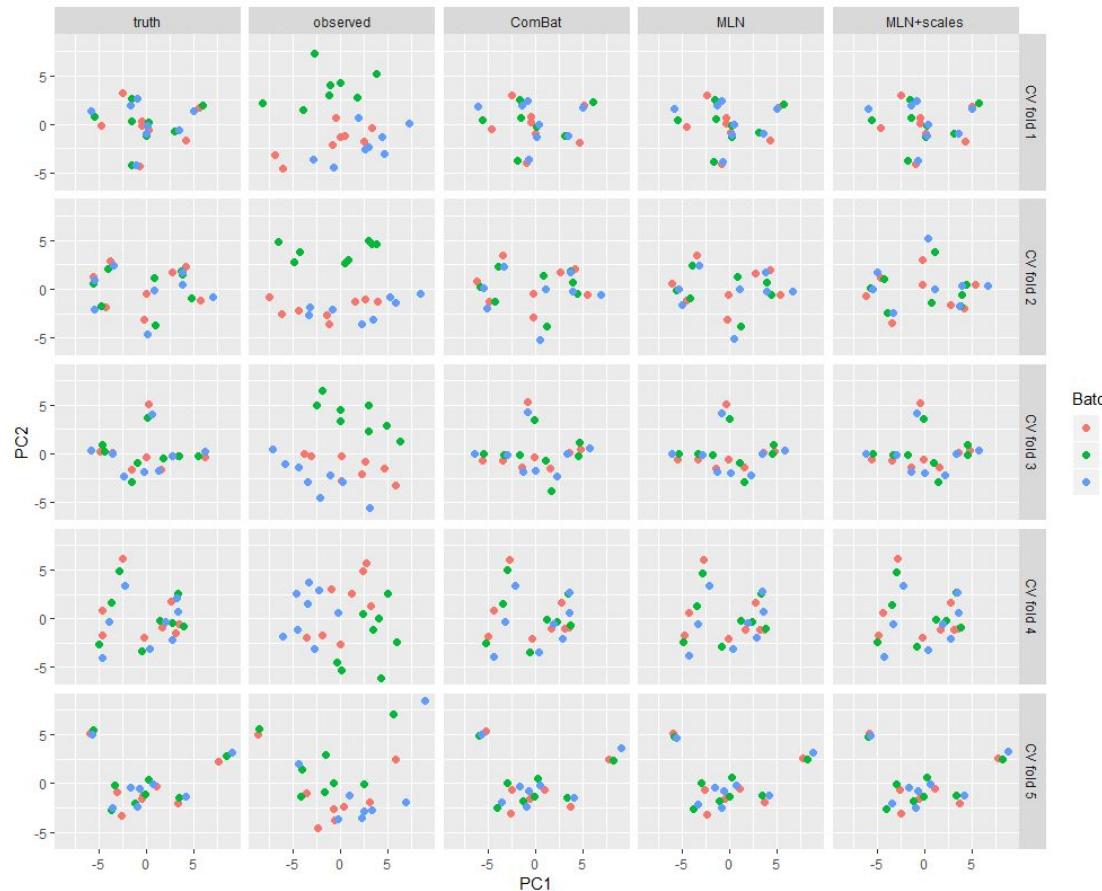
$$\text{score} = \frac{\text{average distance between batches}}{\text{average distance between specimens}}$$

# Simulating microbial sequence count data



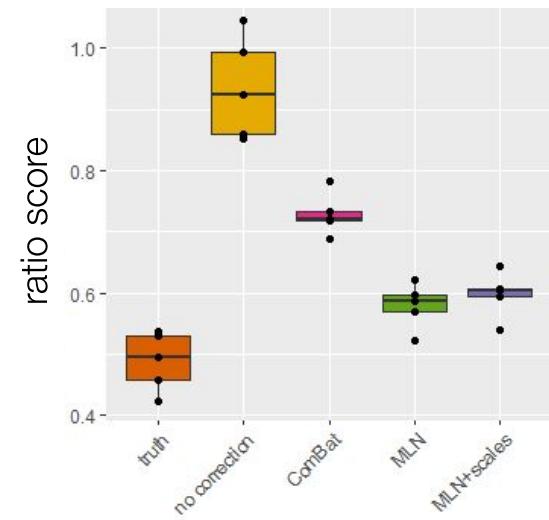
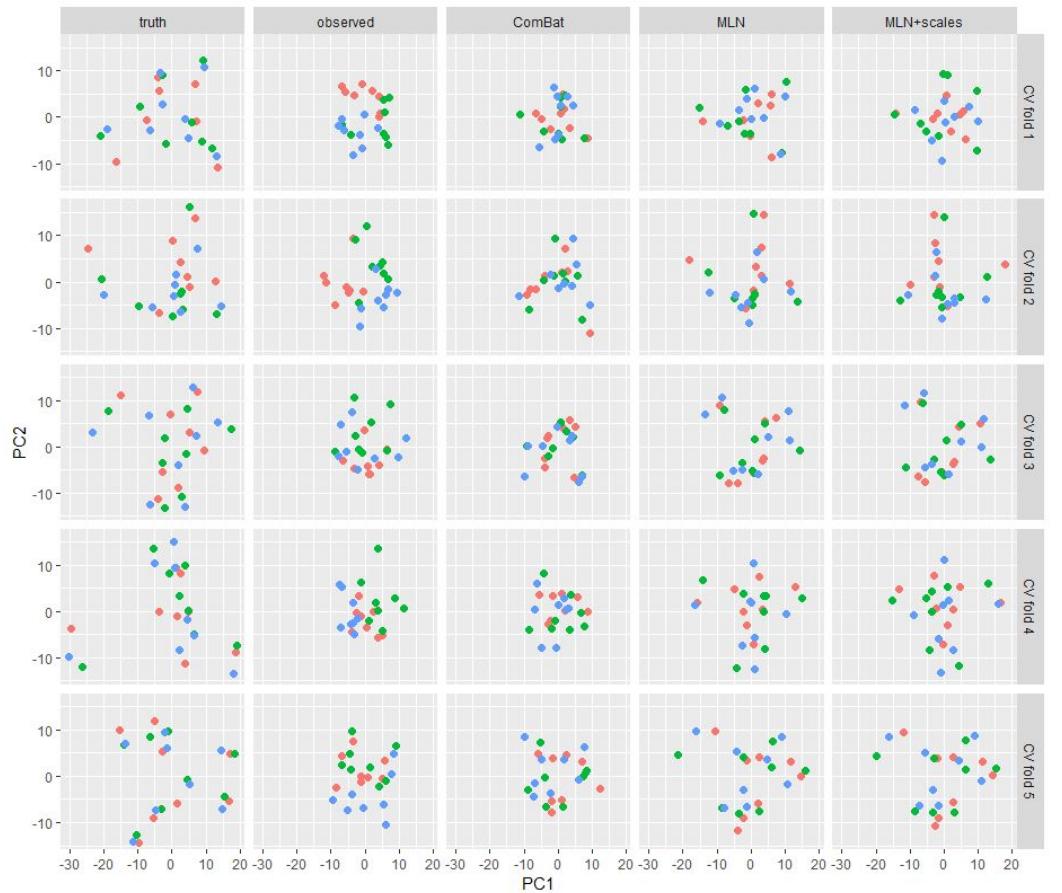
# Results on simulated data: low sparsity

Batch correction over 5 folds



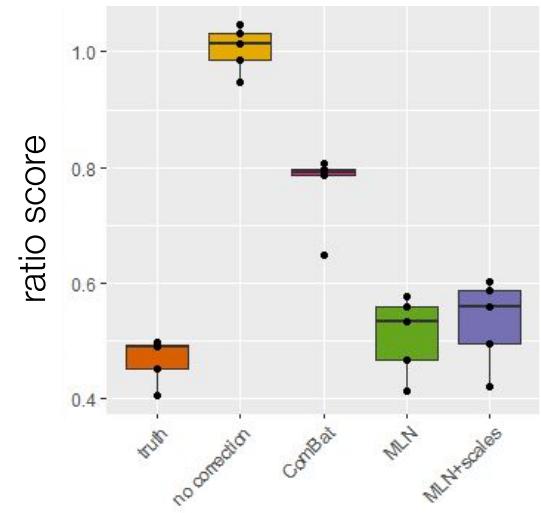
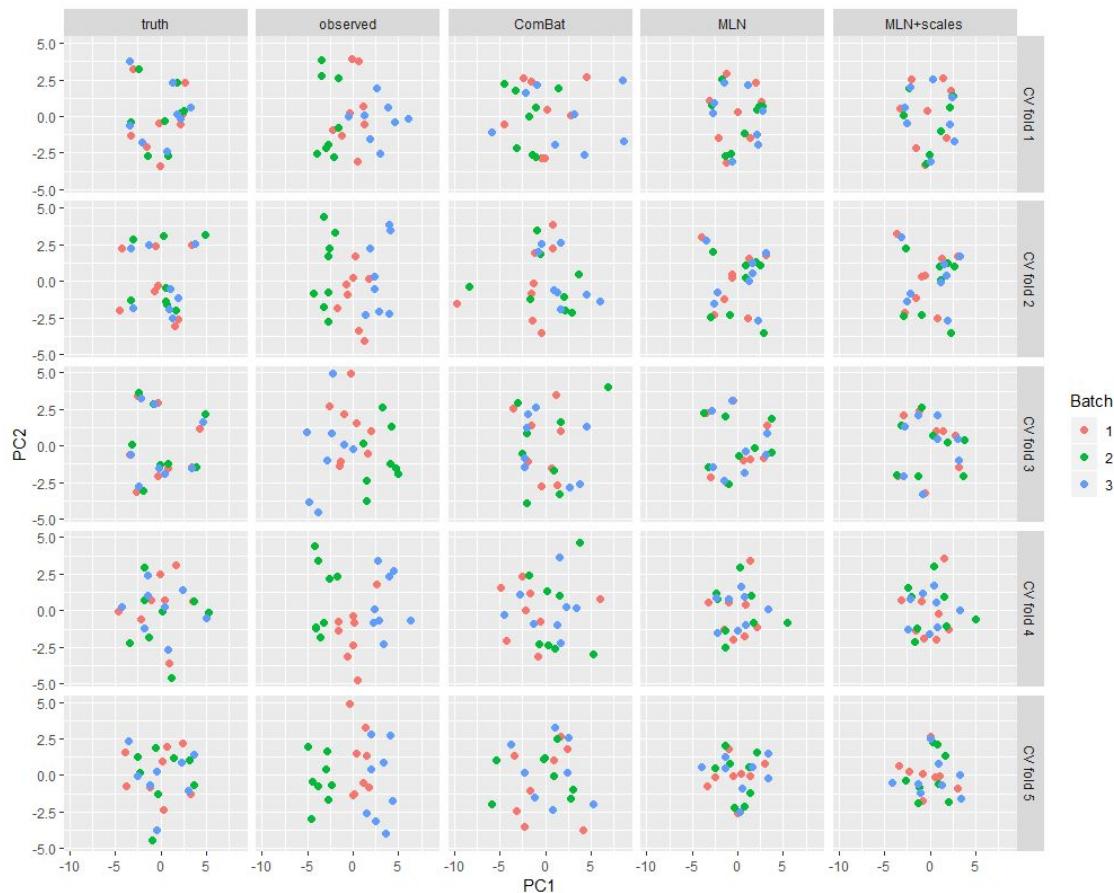
# Results on simulated data: high sparsity

Batch correction over 5 folds



# Results on simulated data: high sparsity (bimodal)

Batch correction over 5 folds



# Real data set

## Pediatric Obesity Metabolism & Microbiome Study (POMMS) baseline

20 patients

Replicates  
One  
Extraction  
Qigen DNeasy  
PowerSoil Kit  
Sequencing  
MiSeq 250b PE

20 patients

Replicates  
Two  
Extraction  
MagAttract 96-well  
plate  
Sequencing  
MiniSeq 150b PE

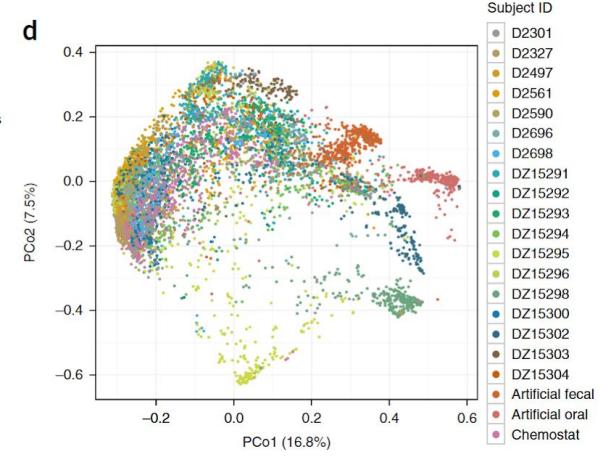
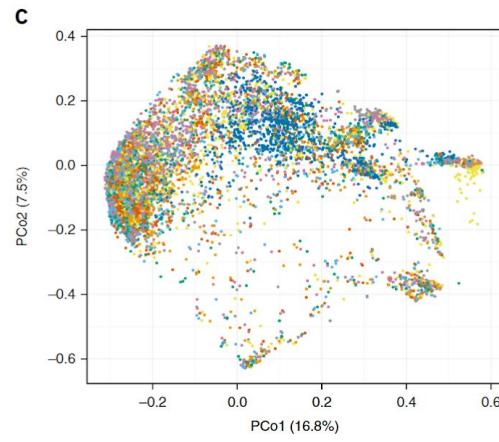
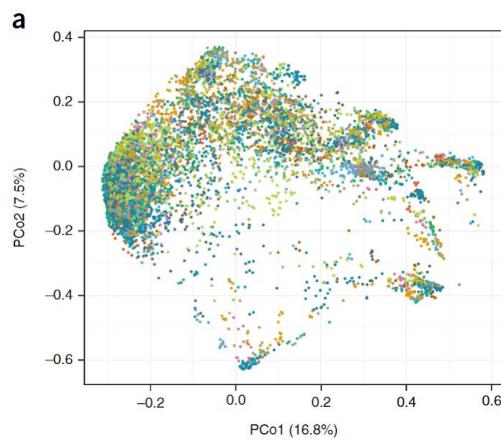
20 patients

Replicates  
Two  
Extraction  
MagAttract 96-well  
plate  
Sequencing  
MiSeq 250b PE

# Real data set

## Microbiome Quality Control Project data

Sinha et al. Nature Biotechnology (2017)

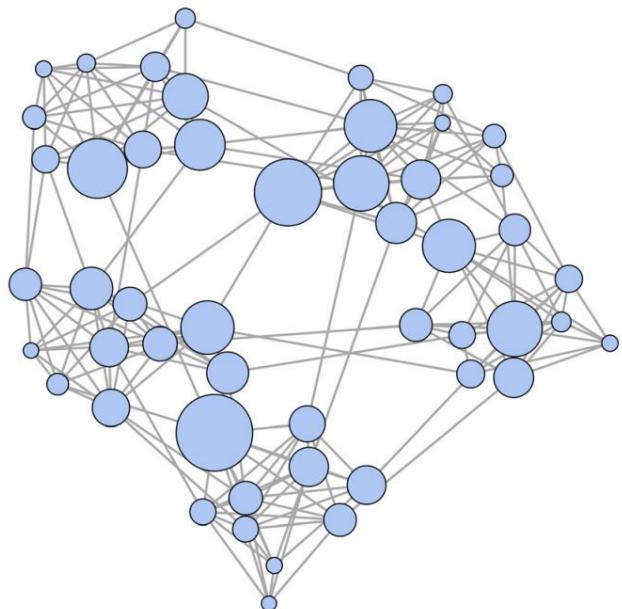


# Aims

- A1** Performing batch effect correction in microbial sequence count data
- A2** Characterizing population variation in gut microbial and functional “dynamics”
- A3** Testing association of dynamics and host fitness

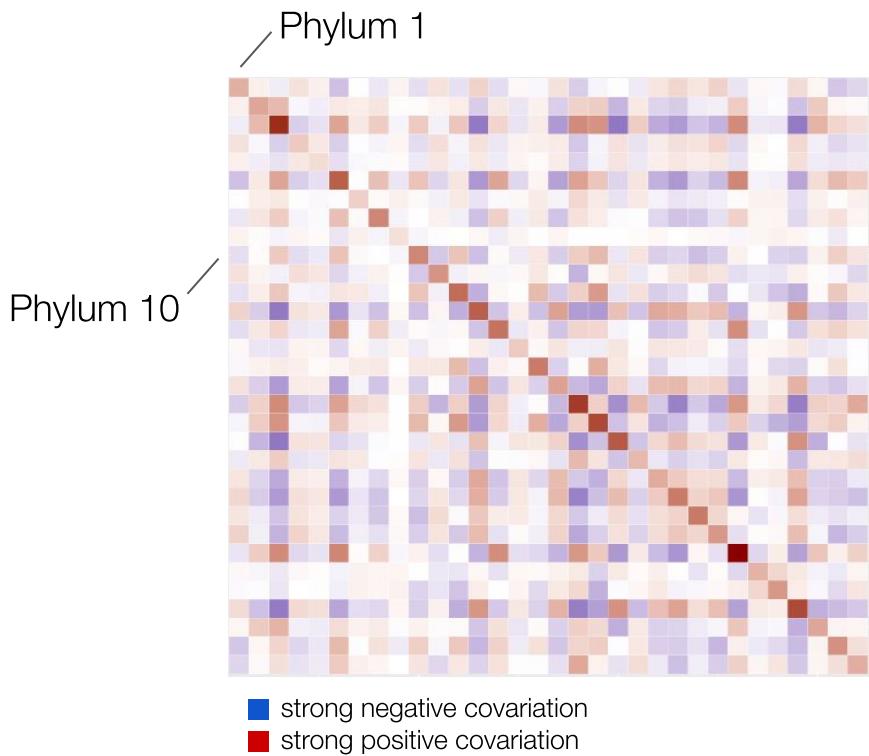
# Do individuals' gut microbiomes evolve according to similar rules?

Interactions between microbes are likely to be context-dependent.



Layeghifard et al. "Disentangling Interactions in the Microbiome: A Network Perspective." Cell Press (2017).

Associations are coarser but easier to characterize.

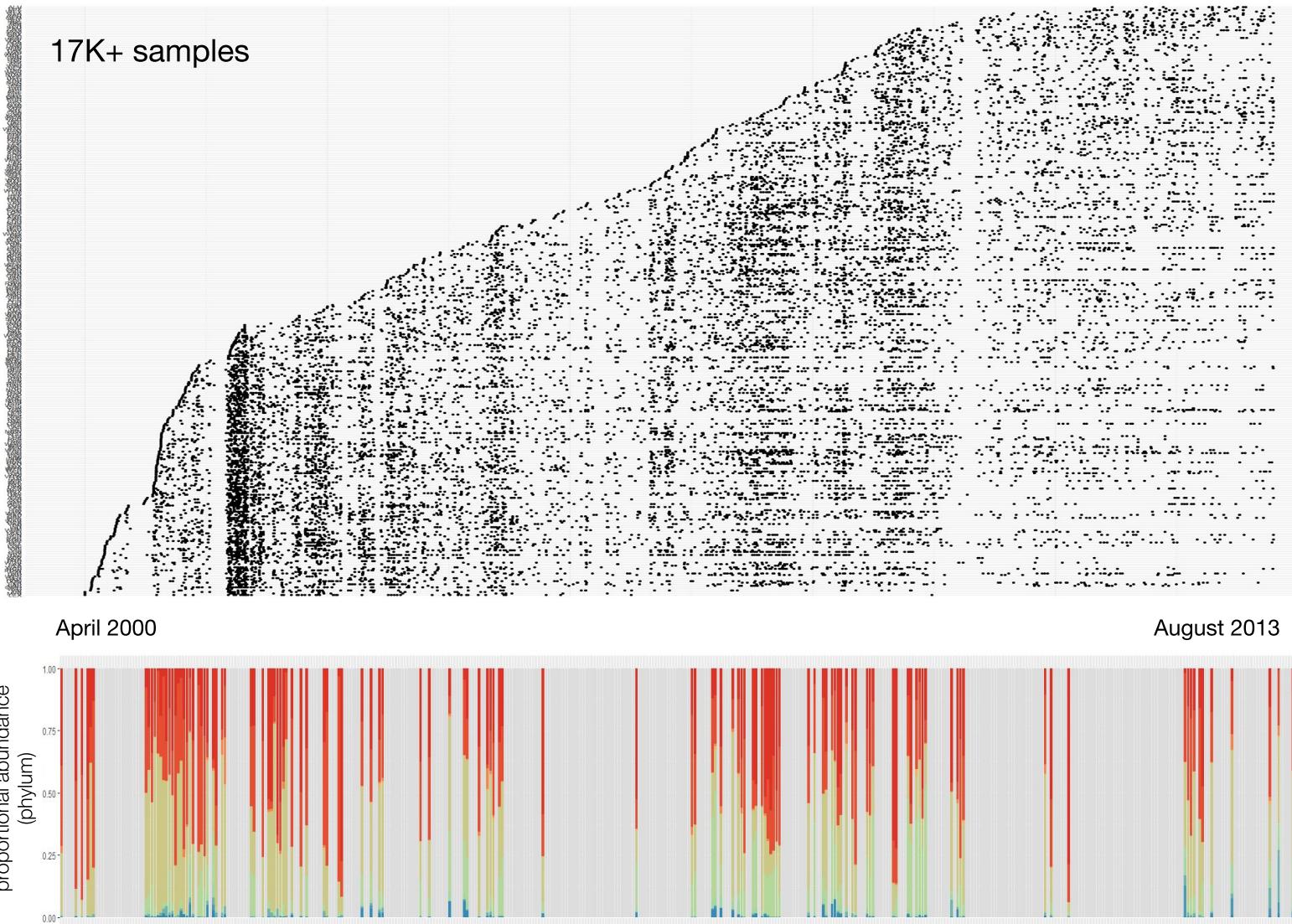


We'll call these "dynamics."

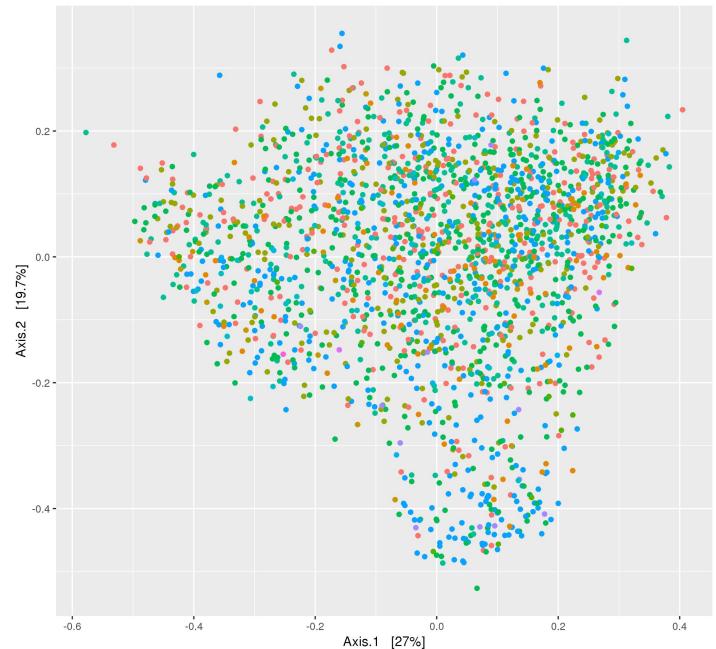
- 1** Amboseli Baboon Research Project 16S data set
- 2** Model & inference framework for estimating microbial covariance
- 3** Considerations: individual modeling, filtration, parameterization
- 4** Quantifying distances between individuals
- 5** Initial results



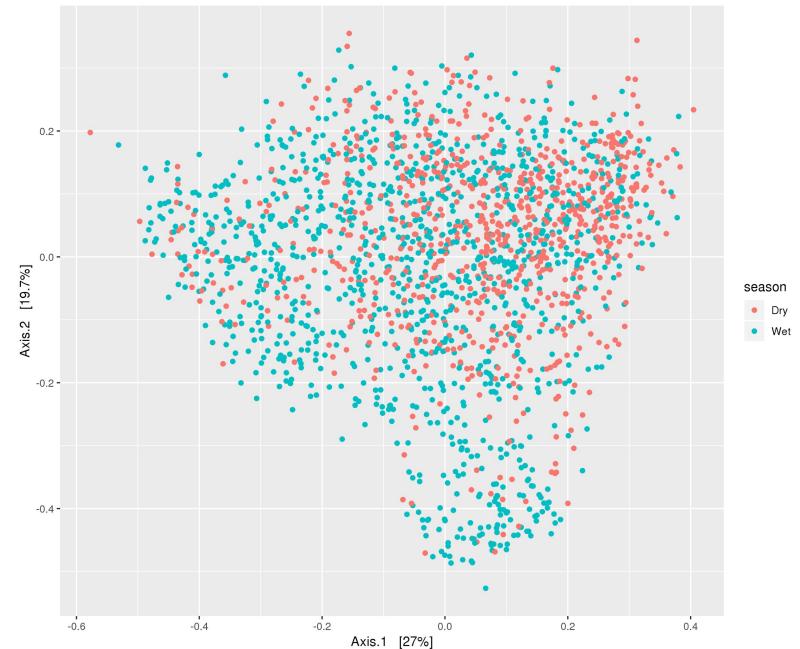
Amboseli Baboon Research Project 16S data set



# High-level associations of composition with social group and season are visible

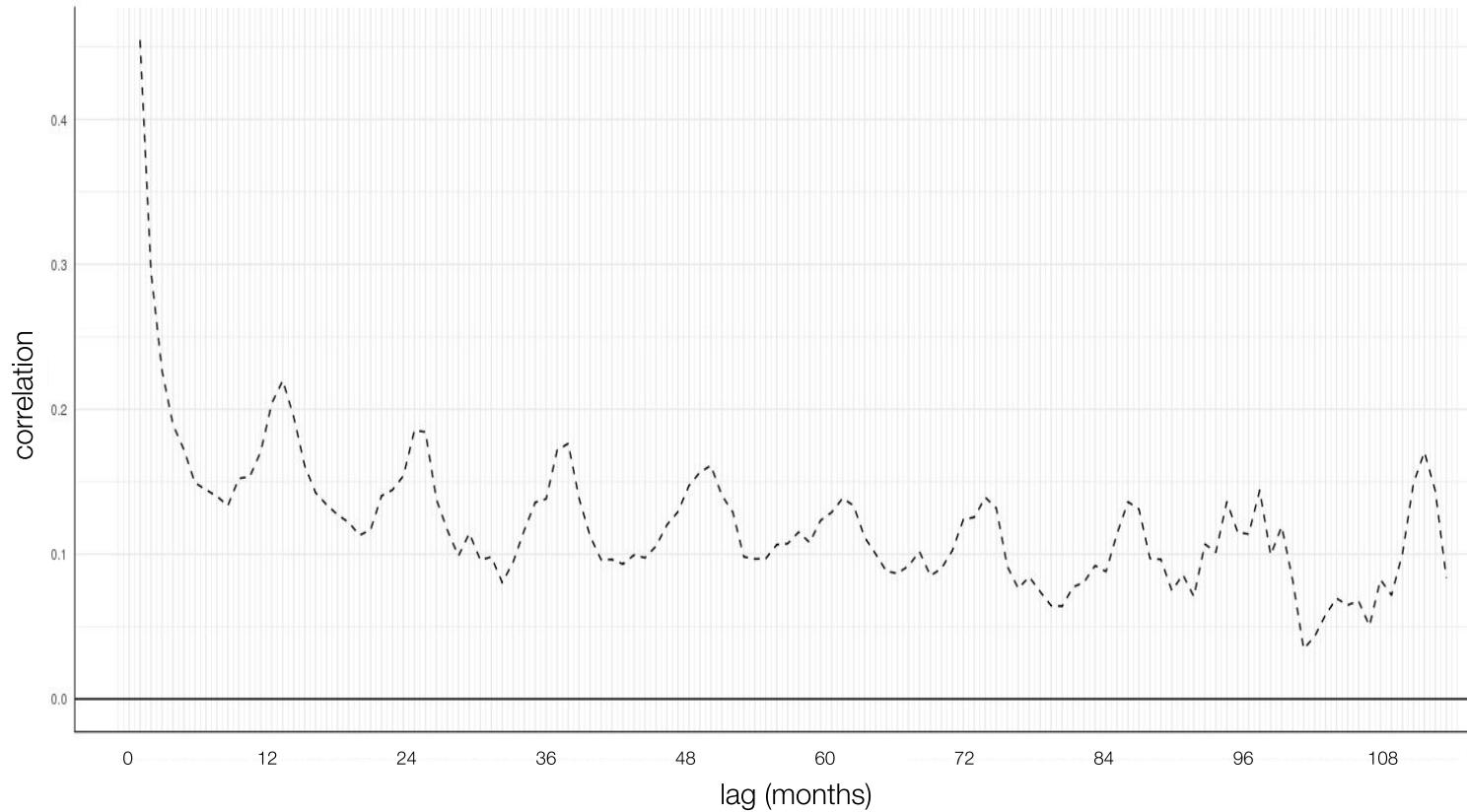


group at collection  
dissimilarity by Bray-Curtis

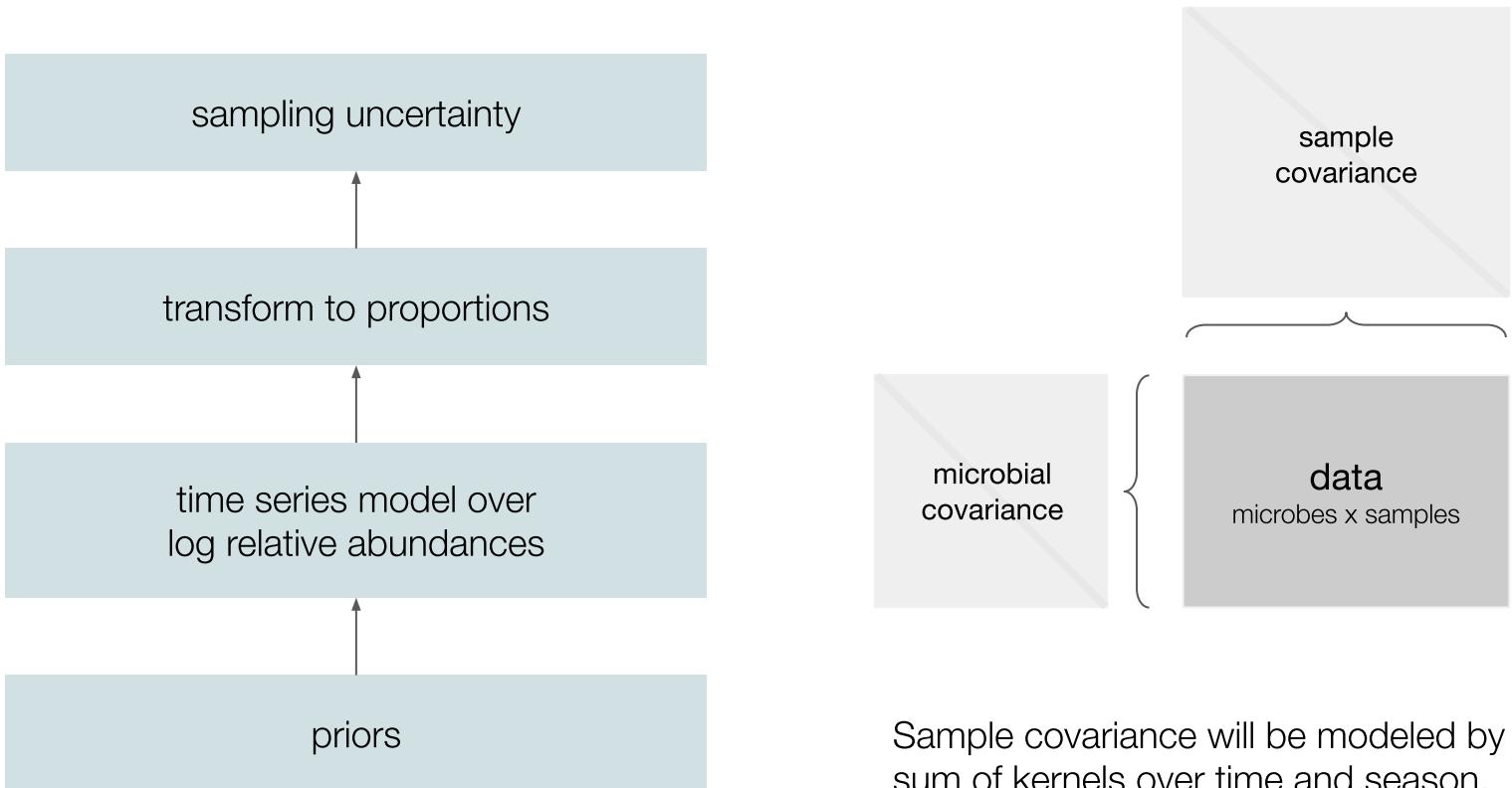


season  
dissimilarity by Bray-Curtis

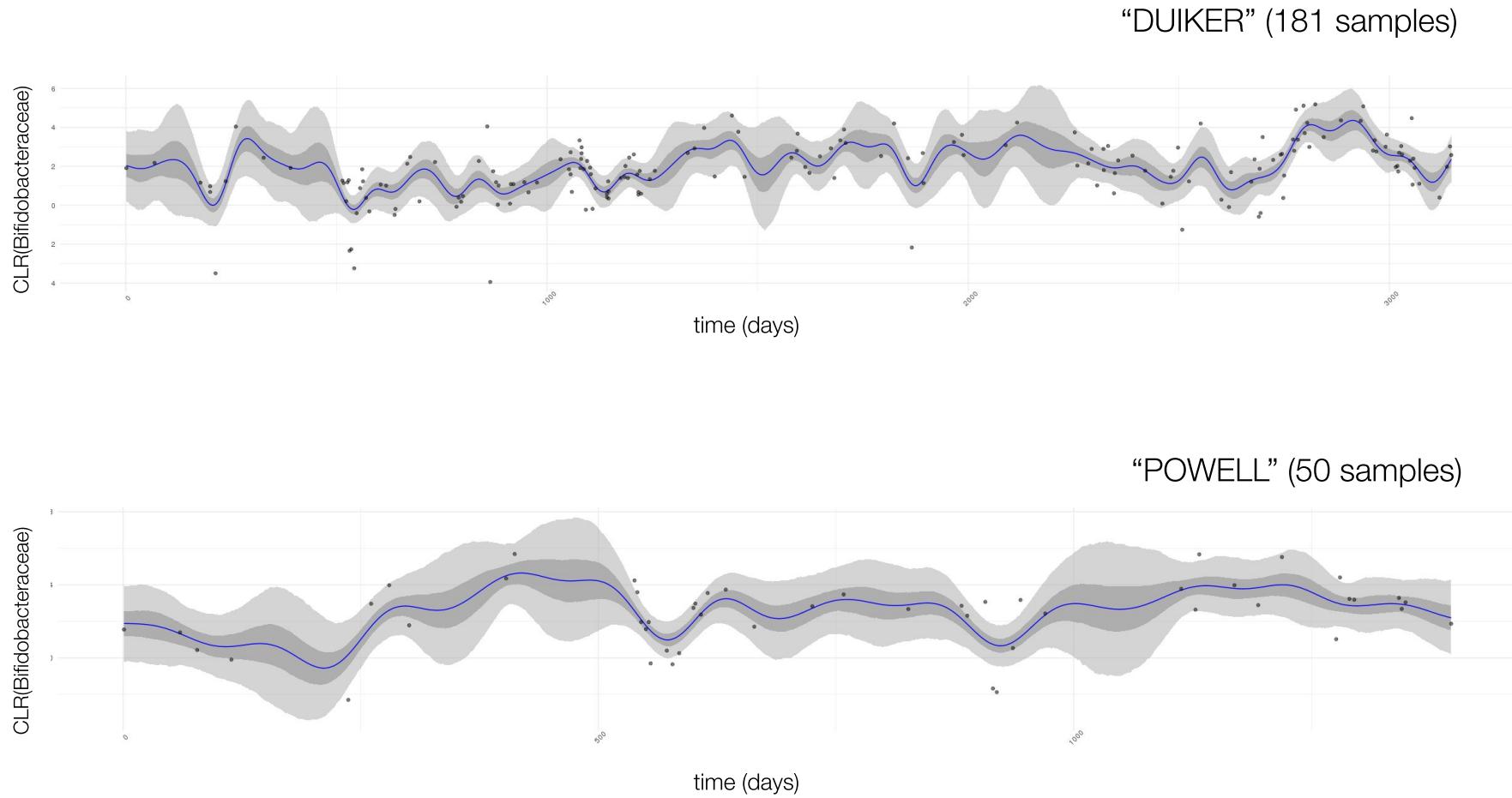
Autocorrelation plots show months-long memory and seasonal fluctuation



We will learn a covariance matrix for each individual via a MLN model housing a Gaussian process

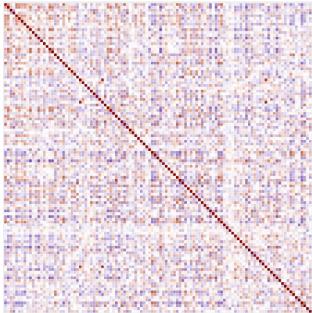


# Differences in sample number are a consideration

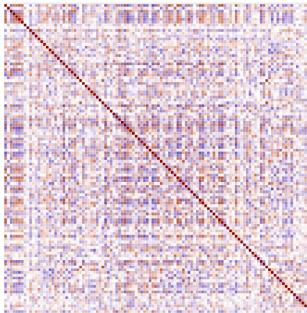


We can quantifying distance between individuals via the Riemannian metric

Individual A

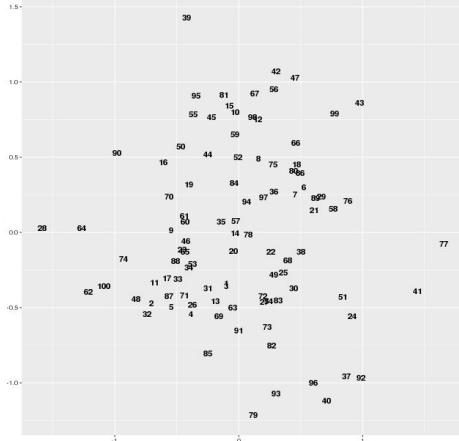
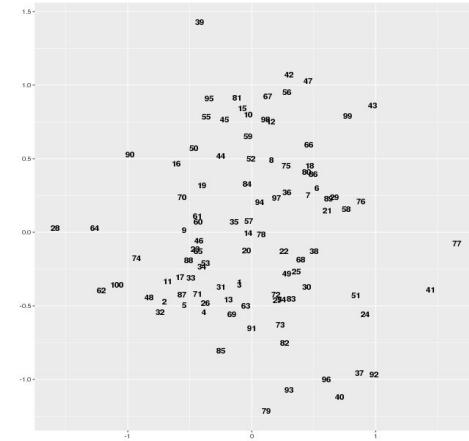
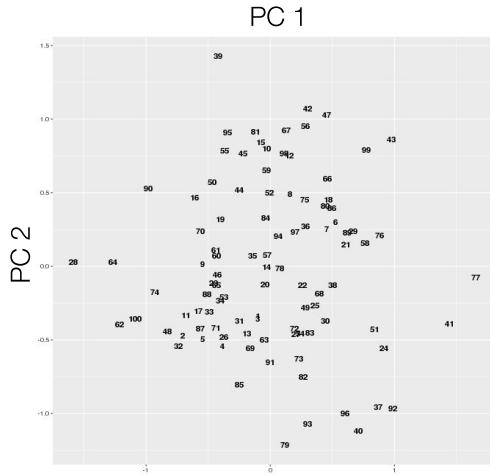


Individual B



$$d(\mathbf{A}, \mathbf{B}) = \sqrt{\text{tr} \left( \ln^2 \left( \sqrt{\mathbf{A}^{-1}} \mathbf{B} \sqrt{\mathbf{A}^{-1}} \right) \right)}$$

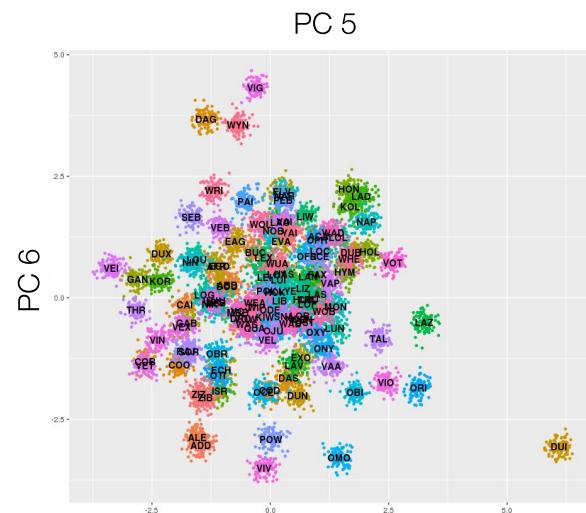
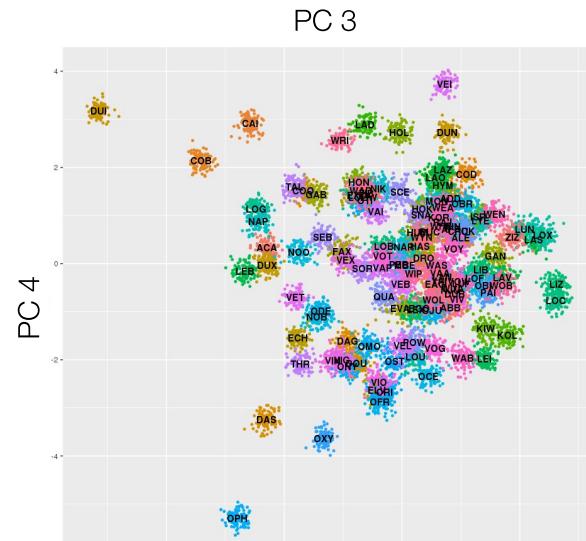
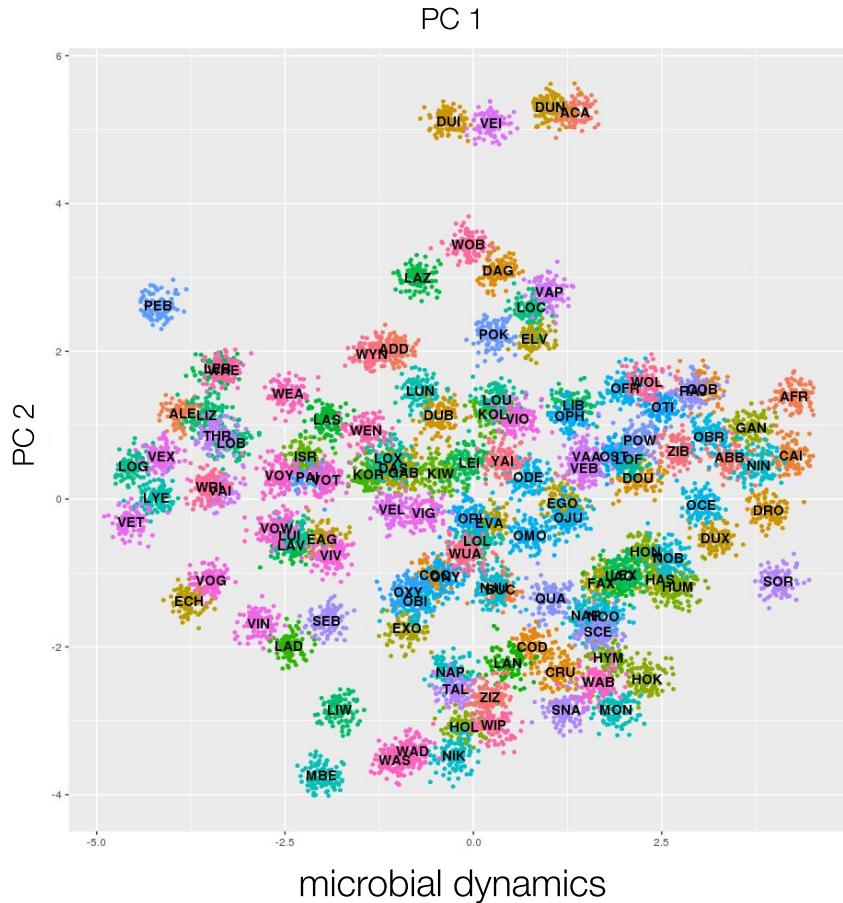
Distances are conserved across common log ratio representations



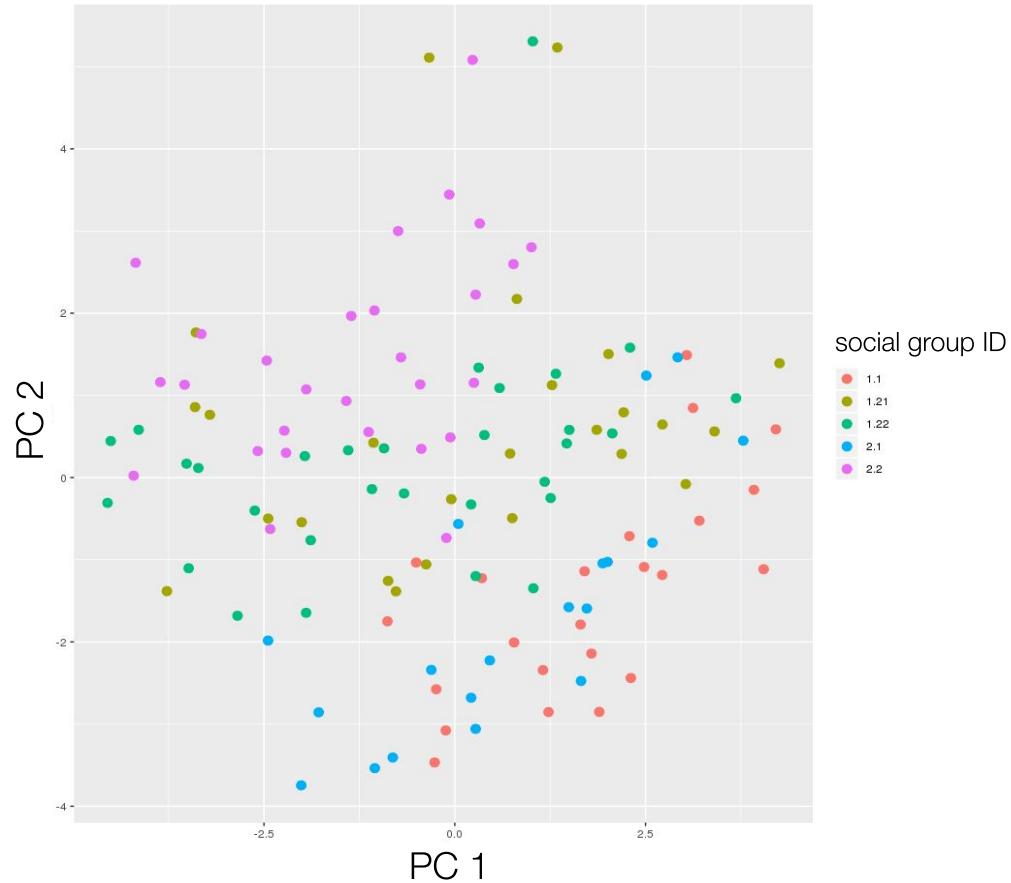
ordination in  $\text{CLR}_{D-1}$

# Visualizing distances in terms of microbial “dynamics”

## Genus level

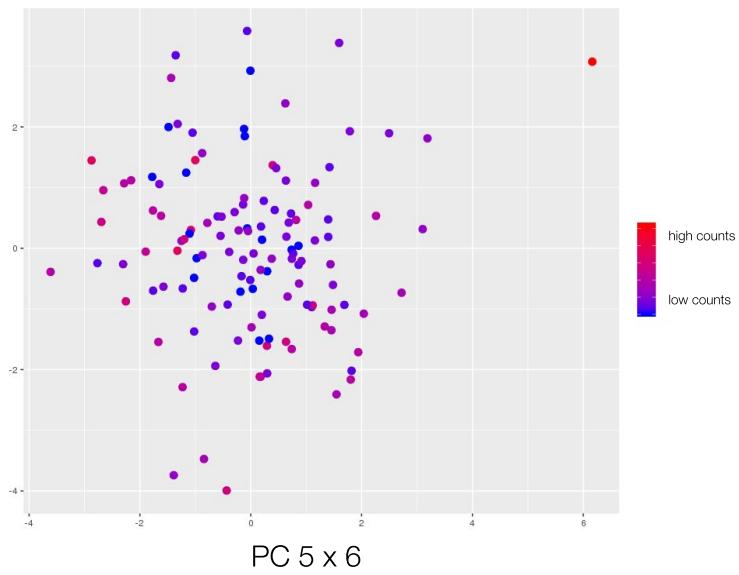
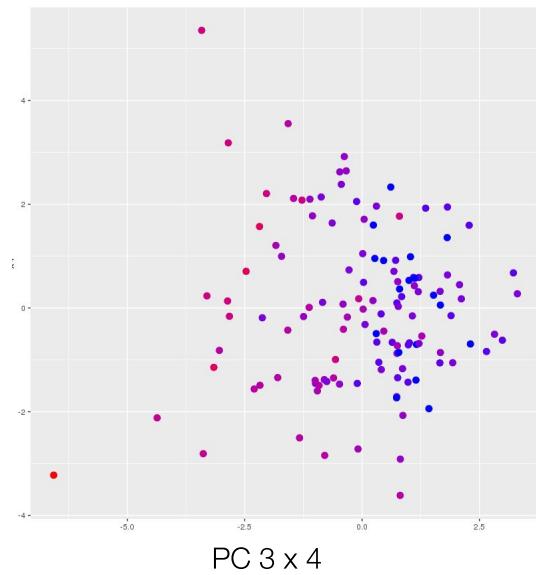
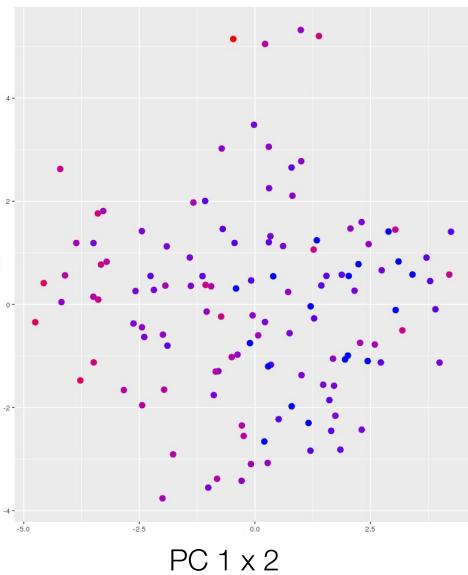


# Social group associates with distances in terms of dynamics



Note: Here individuals are represented by MAP estimate only.

# Sampling characteristics may influence estimates of dynamics



# Can we answer whether individuals have shared or unique dynamics?

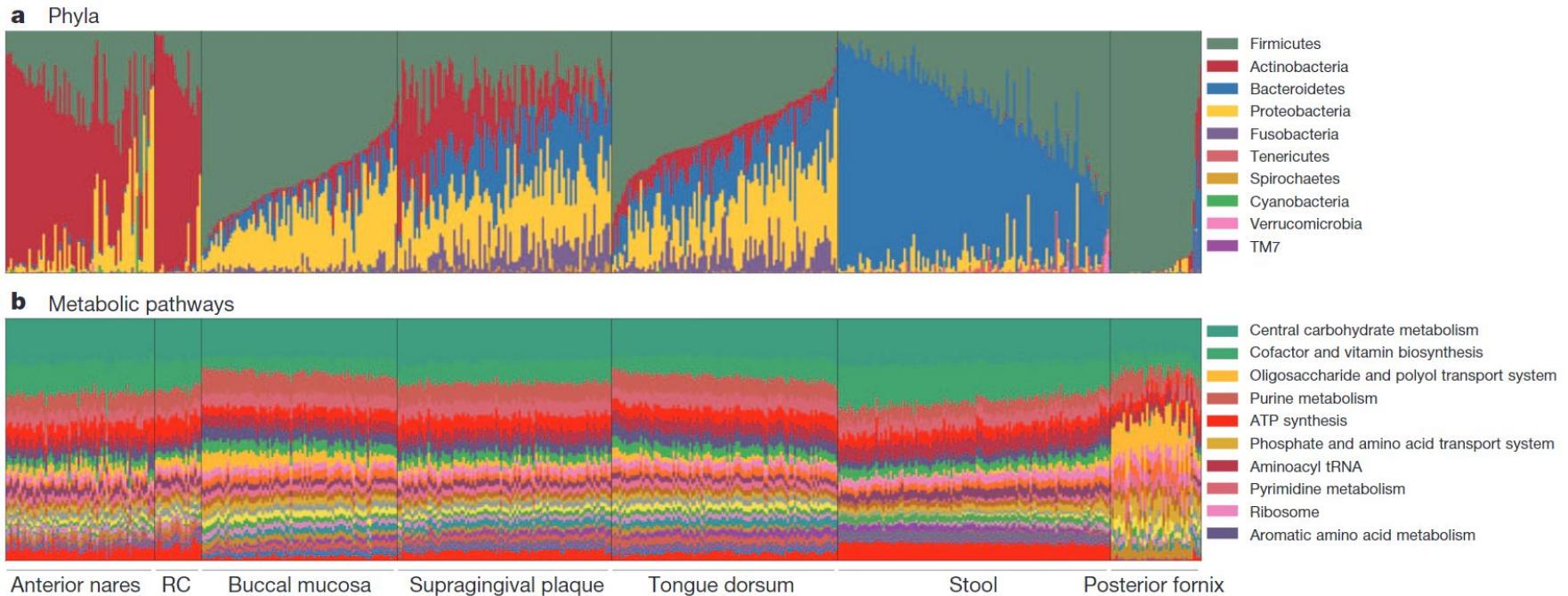
- Individuals can be resolved from each other in terms of their dynamics
- Variation appears continuous
- Need an interpretation for extreme difference
- Social group associates with covariance between microbes

These indicate that shared exposures could explain much of observed distances.

## Further questions

- Are there subsets of microbial relationships in subsets of individuals that are shared?
- Are rare taxa the primary drivers of dissimilarity between individuals?
- Do functional metagenomics data reveal greater similarity between individuals than microbial abundances?

# Comparing functional dynamics to compositional dynamics



Structure, function and diversity of the healthy human microbiome  
The Human Microbiome Project Consortium in Nature (2012)

# Aims

- A1** Performing batch effect correction in microbial sequence count data
- A2** Characterizing population variation in gut microbial and functional “dynamics”
- A3** Testing association of dynamics and host fitness

# Do individuals with atypical dynamics also have atypical fitness outcomes?

- 1 Define *fitness*: normalized number of offspring
- 2 Define a kernel ( $K$ ) between individuals based on similarity of dynamics
- 3 Model out other contributors to fitness (e.g. rank)
- 4 Interpret  $\tau / (\alpha + \tau)$  as variance explained

$$g\left(\begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}\right) \sim N(X\beta, \alpha I + \tau K)$$

$y_i$  fitness outcome of individual i

$X\beta$  baseline outcome explained by covariates

$\alpha I$  scaled noise

$\tau K$  variation due to similarity of dynamics

```
for  $\sigma$  in  $\{\sigma_1, \dots, \sigma_p\}$  do
    for each  $k$  of  $K$  folds do
        let  $Y_{\text{train}}^{(k)} \sim N(X\beta, \alpha I + \tau K)$  ;
        optimize  $(\alpha, \tau)$  ;
        let  $\hat{Y}_{\text{test}}^{(k)} = \mathbb{E}(Y_{\text{test}}^{(k)} | Y_{\text{train}}^{(k)})$  ;
        let  $\text{err}_{\text{pred}}^{(k)} = \|\hat{Y}_{\text{test}}^{(k)} - Y_{\text{train}}^{(k)}\|$  ;
    end
end
choose  $(\alpha, \tau, \sigma)$  that minimize  $\text{err}_{\text{pred}}$ 
```

# **A1** Supplement

# Common log relative abundance transformations

$$\text{alr}(x_1, \dots, x_D) = \mathbf{y} = \left( \ln \frac{x_1}{x_D}, \dots, \ln \frac{x_{D-1}}{x_D} \right)$$

$$\text{alr}^{-1}(\mathbf{y}) = \mathbf{x} = \mathcal{C}(\exp(\mathbf{y}; 0))$$

$$\text{clr}(x_1, \dots, x_D) = \mathbf{y} = \left( \ln \frac{x_1}{g(\mathbf{x})}, \dots, \ln \frac{x_D}{g(\mathbf{x})} \right) \quad g(\mathbf{x}) = \sqrt[D]{x_1 \dots x_D}$$

$$\text{clr}^{-1}(\mathbf{y}) = \mathbf{x} = \mathcal{C}(\exp(\mathbf{y}))$$

$$\text{ilr}(x_1, \dots, x_D) = \mathbf{y} = \ln(\mathbf{x}) \cdot \mathbf{V}$$

$$\text{ilr}^{-1}(\mathbf{y}) = \mathbf{x} = \mathcal{C}(\exp(\mathbf{x} \cdot \mathbf{V}^T))$$

# Multinomial logistic normal model formulations used

linear regression

$$Y_j \sim \text{Multinomial}(\pi_j)$$

$$\pi_j = \text{ALR}_d^{-1}(\eta_j)$$

$$\eta_j \sim N(\Lambda X_j, \Sigma)$$

$$\Lambda \sim N(\Theta, \Sigma, \Gamma)$$

$$\Sigma \sim W^{-1}(\Xi, v)$$

linear regression+scales

$$Y_j \sim \text{Multinomial}(\pi_j)$$

$$\pi_j = \text{ALR}_d^{-1}(\eta_j)$$

$$\eta_j \sim N(\Lambda X_j, \Sigma)$$

$$\Lambda \sim N(\Theta, \Sigma, \Gamma)$$

$$\Sigma \sim W^{-1}(\Xi, v)$$

$$\Gamma = \sigma_0 \Gamma_0 + \sum_{i=1}^q \sigma_i \Gamma_i$$

per-batch scatter

## Batch effect simulation scheme

$$Y_j \sim \text{Multinomial}(\pi_j)$$

$$\pi_j = \text{ALR}_d^{-1}(\eta_j)$$

$$\eta = MZ + BX + E \quad (d - 1 \times n)$$

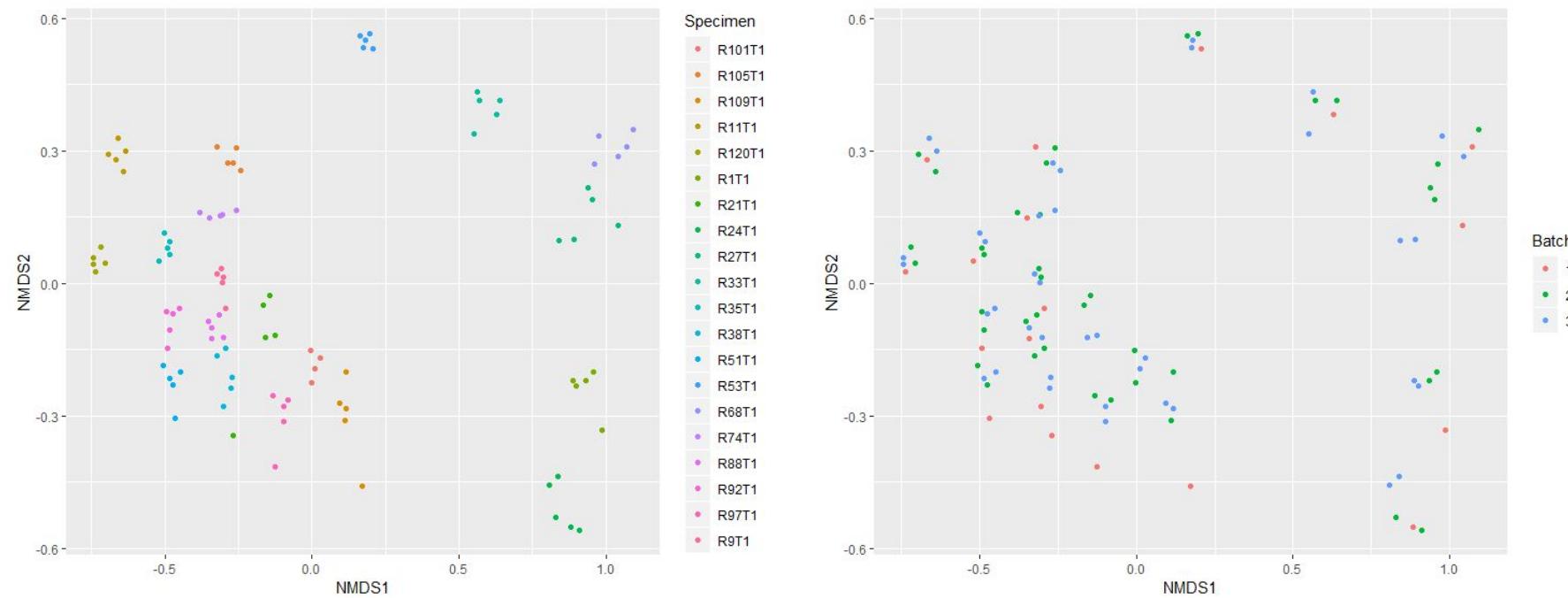
$$M \sim N(0, \Sigma_M, I_m) \quad (d - 1 \times m) \quad m = \# \text{ individuals}$$

$$B \sim N(0, \Sigma_B, \Lambda) \quad (d - 1 \times q) \quad q = \# \text{ batches}$$

$$E \sim N(0, \Sigma_E, I_n) \quad (d - 1 \times n) \quad n = \# \text{ total samples}$$

# Visualizing the POMMS baseline batch data set

Dissimilarity by Bray-Curtis



# **A2** Supplement

# MLN + dynamic linear model

$$Y_t \sim \text{Multinomial}(\pi_t)$$

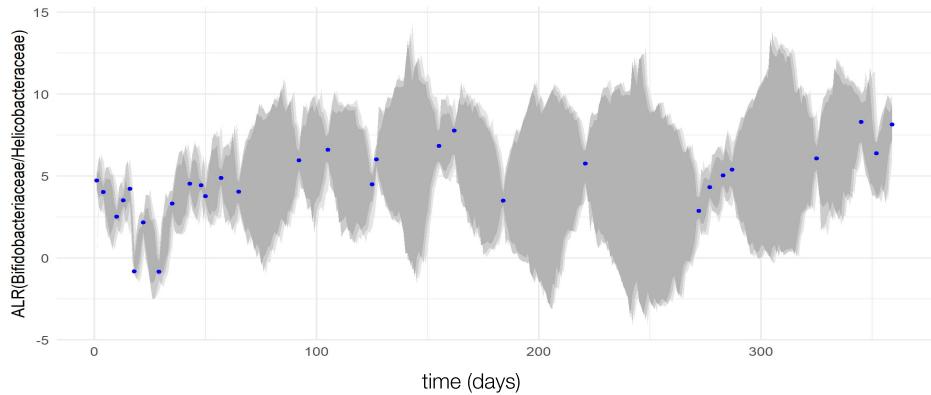
$$\pi_t = \text{ALR}_d^{-1}(\eta_t)$$

$$\eta_t^T = F_t^T \Theta_t + v_t^T \quad v_t \sim N(0, \gamma_t \Sigma)$$

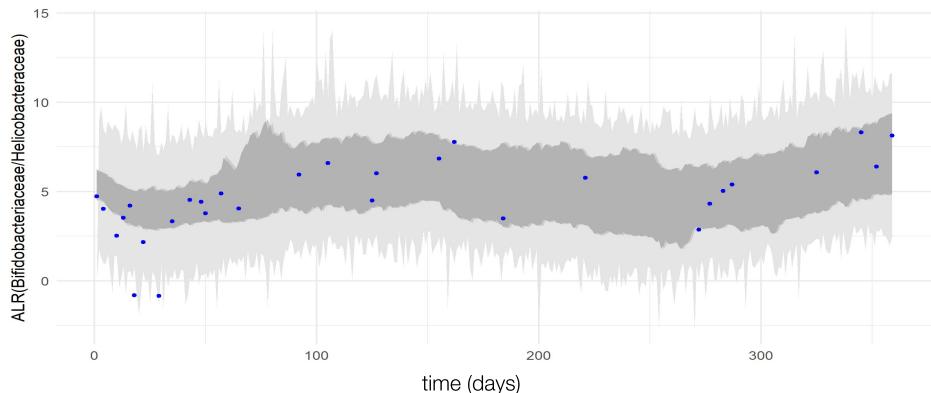
$$\Theta_t = G_t \Theta_{t-1} + w_t \quad w_t \sim N(0, W_t, \Sigma)$$

$$\Theta_0 \sim N(M_0, C_0, \Sigma)$$

$$\Sigma \sim \text{Wishart}^{-1}(\Xi, v)$$



Fixed low noise in  $v_t$



Optimized noise in  $v_t$  and  $w_t$

# MLN + Gaussian process

$$Y_t \sim \text{Multinomial}(\pi_t)$$

$$\pi_t = \text{ALR}_d^{-1}(\eta_t)$$

$$\eta \sim \mathcal{N}(\Lambda[X], \Sigma, I_N), \quad \text{i.e. } \Sigma \otimes I_N$$

$$\Lambda[X] \sim \text{GP}(\Theta[X], \Sigma, \Gamma[X]), \quad \Sigma \otimes \Gamma$$

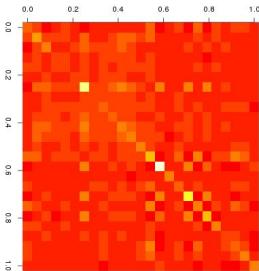
$$\Sigma \sim \text{Wishart}^{-1}(\Xi, v)$$

$\Gamma$ : squared exp. kernel (time) +  
periodic kernel (season) +  
white noise (technical noise)

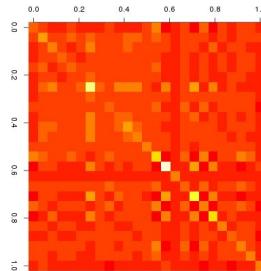
# Downsampling can inform a lower informative sample number

EX: Start with all samples for a well-sampled individual, progressively downsample, and fit the model.

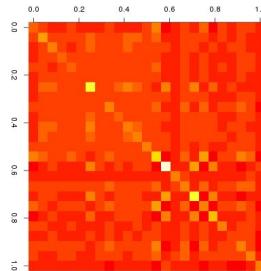
100% samples (181)



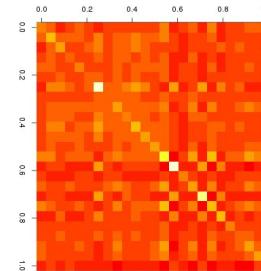
80% of samples



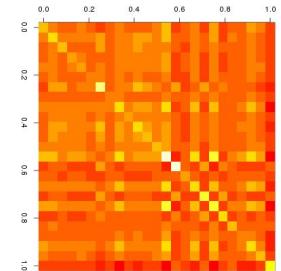
60% of samples



40% of samples



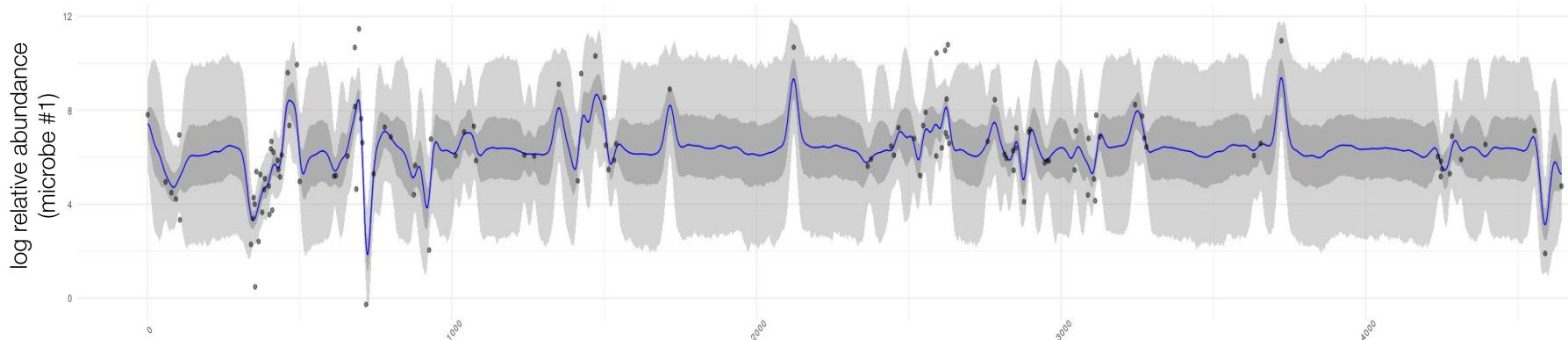
20% of samples



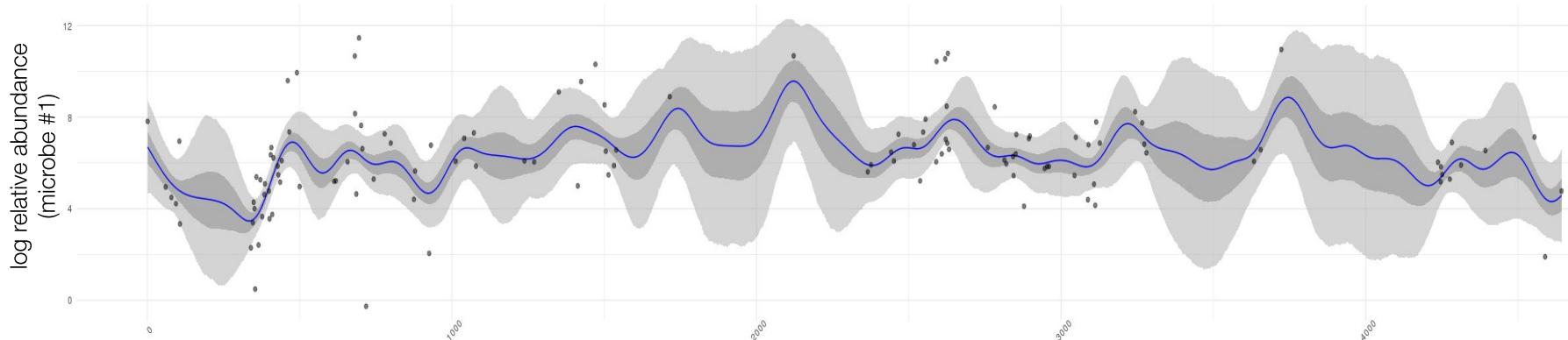
Determine minimum sample number where model fits exhibit undesirable behavior.

# Posterior predictive fits, varying choice of GP bandwidth

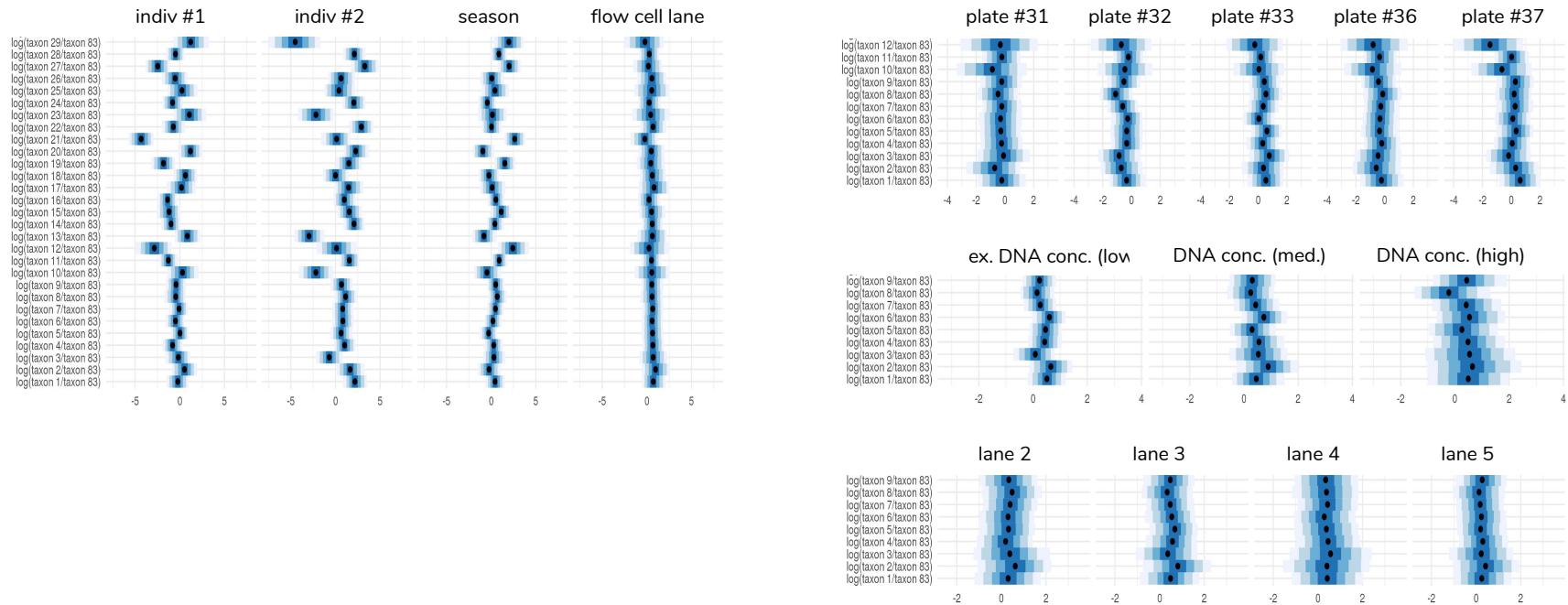
~30 day information window



~180 day information window



Regression shows reliable effects of individual ID and season on log relative abundance but not batch processing variables

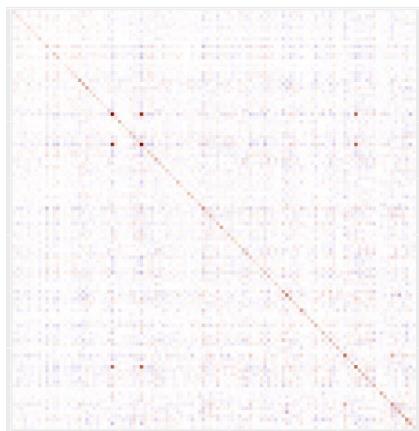


Excerpts of all variables and posterior intervals shown. Intervals that spanned zero were judged as *no effect*.

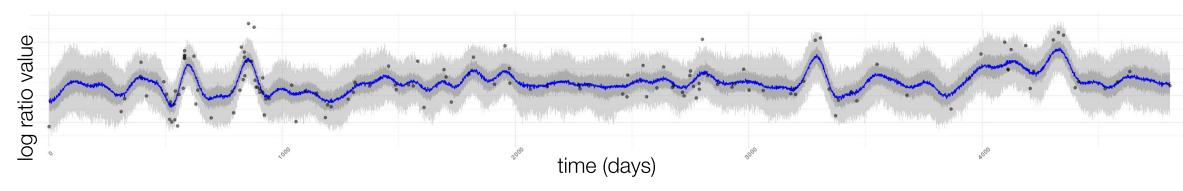
# Correlated taxa seem reasonable

Shown: 3 highly variable, strongly (+) correlated Firmicutes

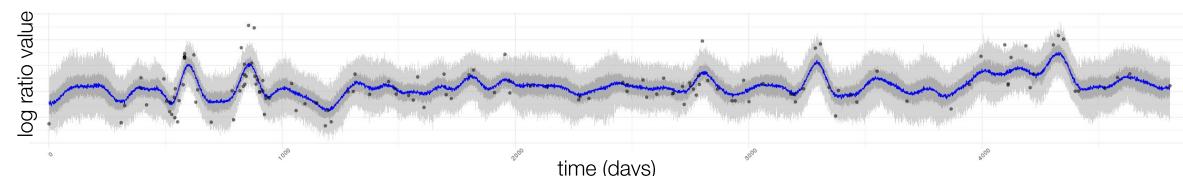
Individual "VET"



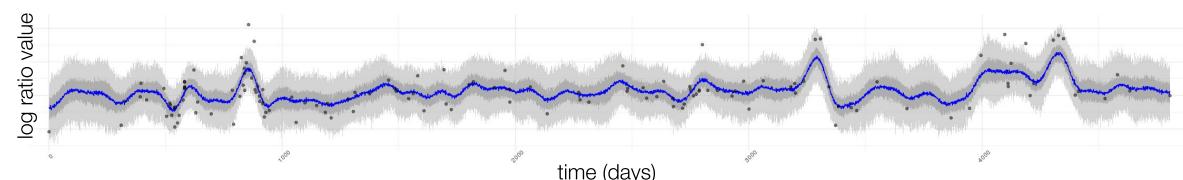
CLR(Planococcaceae/NA)



CLR(Planococcaceae/Solibacillus)

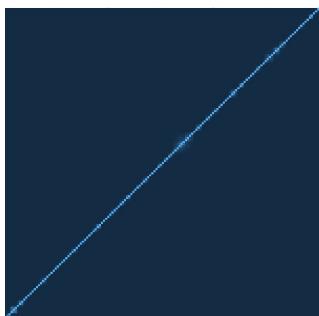


CLR(Bacillaceae/Bacillus)

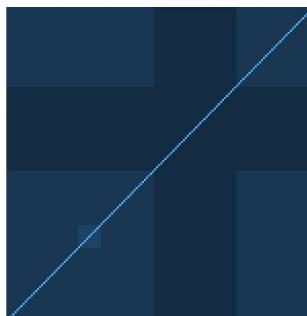


# A variance component model motivates relative contributions of GP kernels

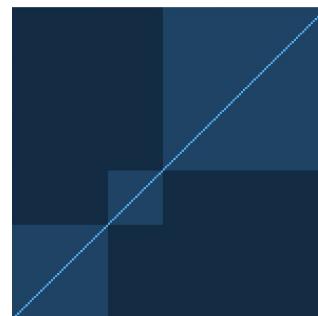
- 1 Variance components were constructed as in these toy examples



weekly sample correlation



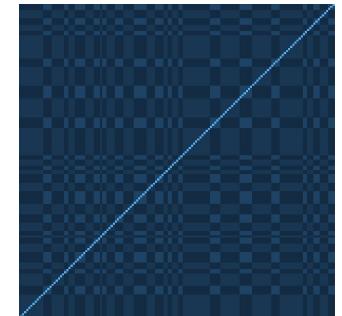
age correlation



individual correlation



group correlation



seasonal correlation

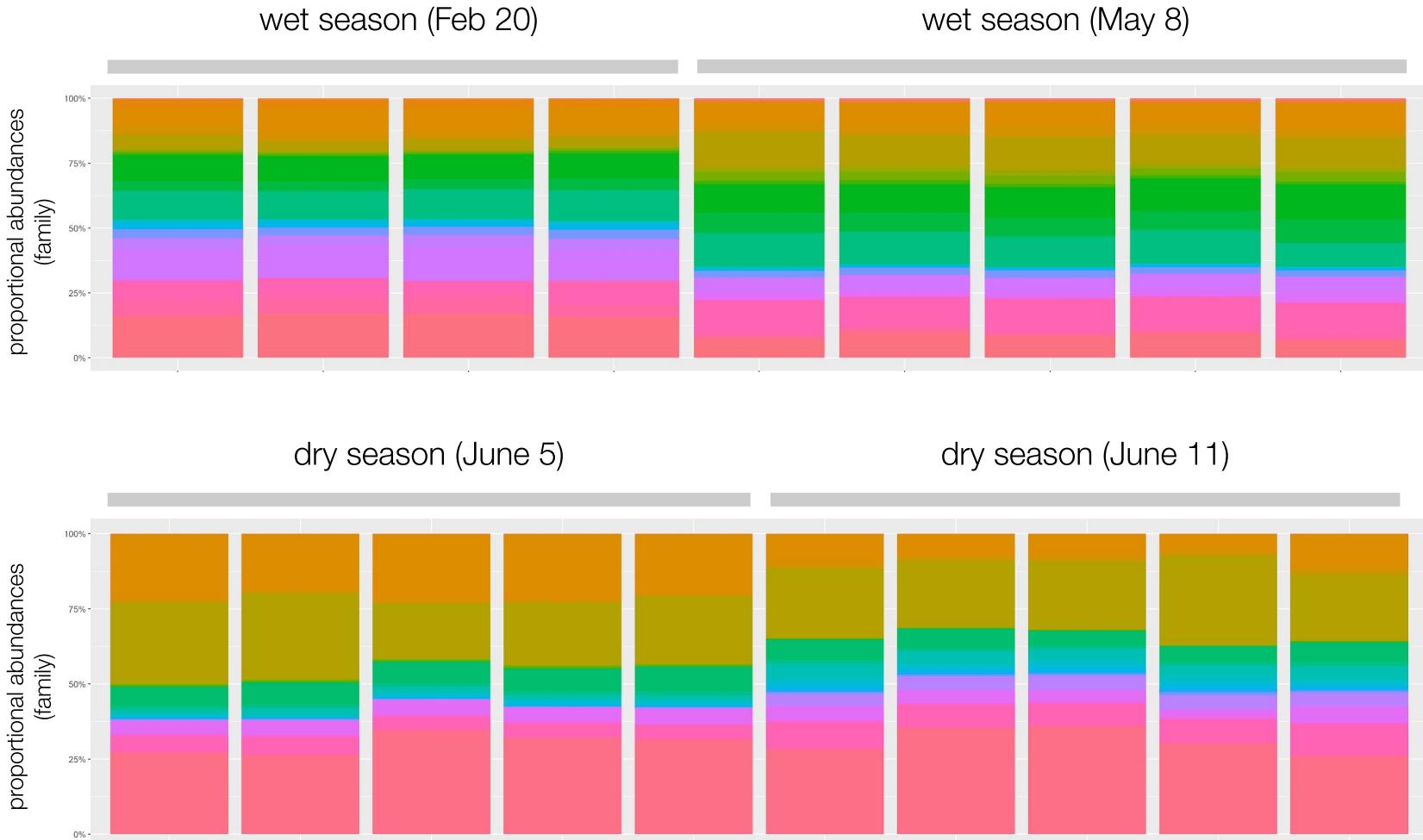
- 2 Posterior density of this simple model was optimized to give component scales  $\gamma_1 \dots \gamma_k$

$$\eta \sim \text{MN}(M, U, V)$$

$$U \sim W^{-1}(v, \Psi)$$

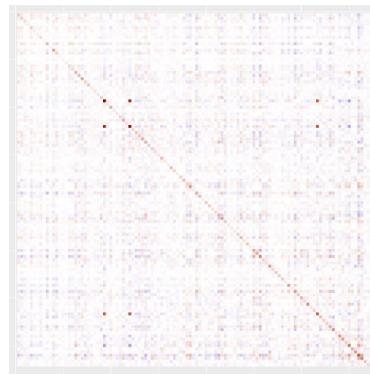
$$V = \gamma_1 V_1 + \gamma_2 V_2 + \dots + \gamma_k V_k$$

# Variation in replicates can inform noise model

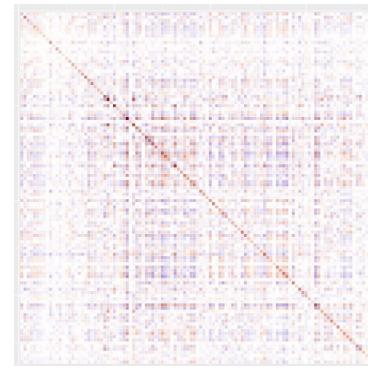


## Microbial dynamics PC 1 may capture *net covariation*

individual  
“LOGAN” from  
extreme lower  
end of PC 1

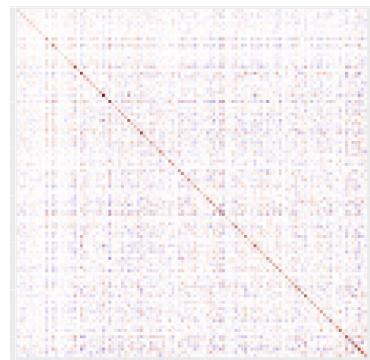


individual  
“AFRICA” from  
extreme upper  
end of PC 1

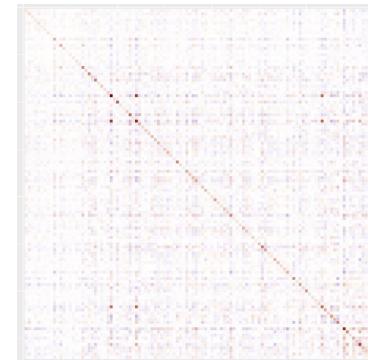


## Microbial dynamics PC 2 may capture *net variation*

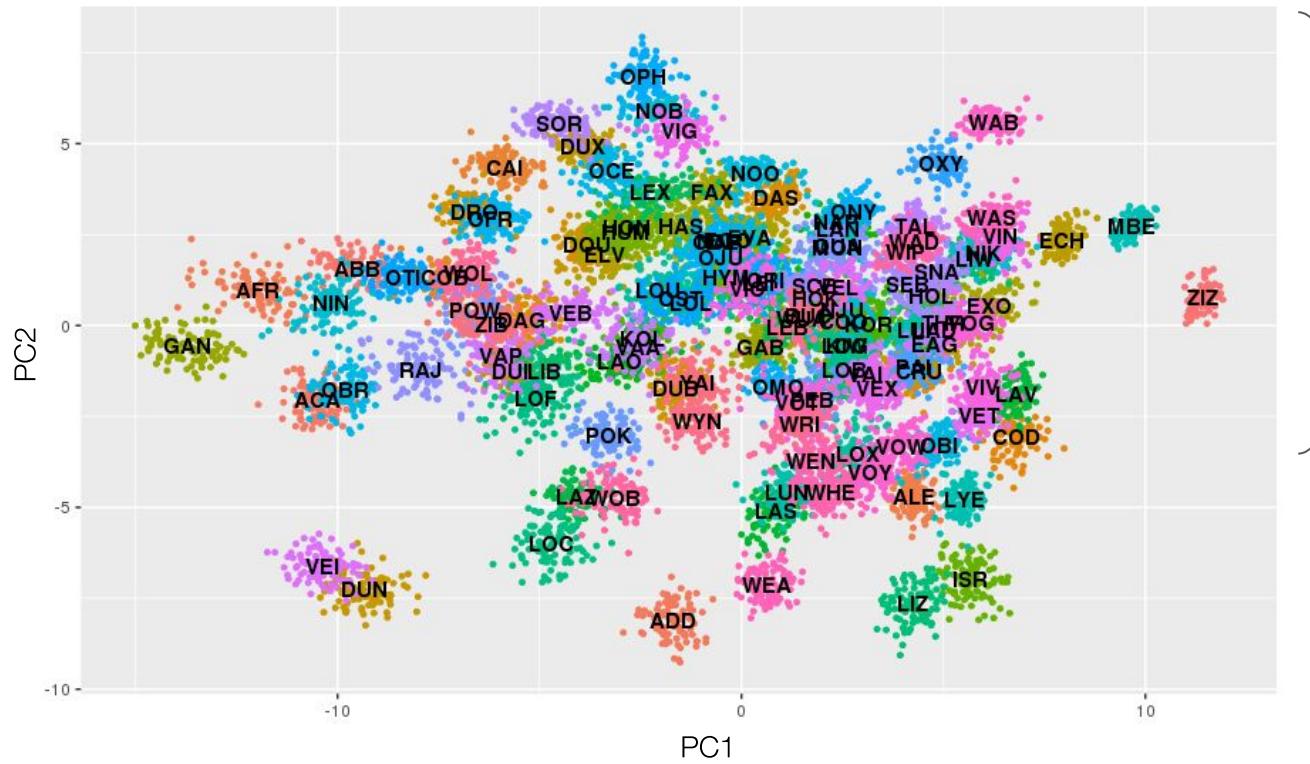
individual  
“WASP” from  
extreme lower  
end of PC 1



individual  
“DUIKER” from  
extreme upper  
end of PC 1



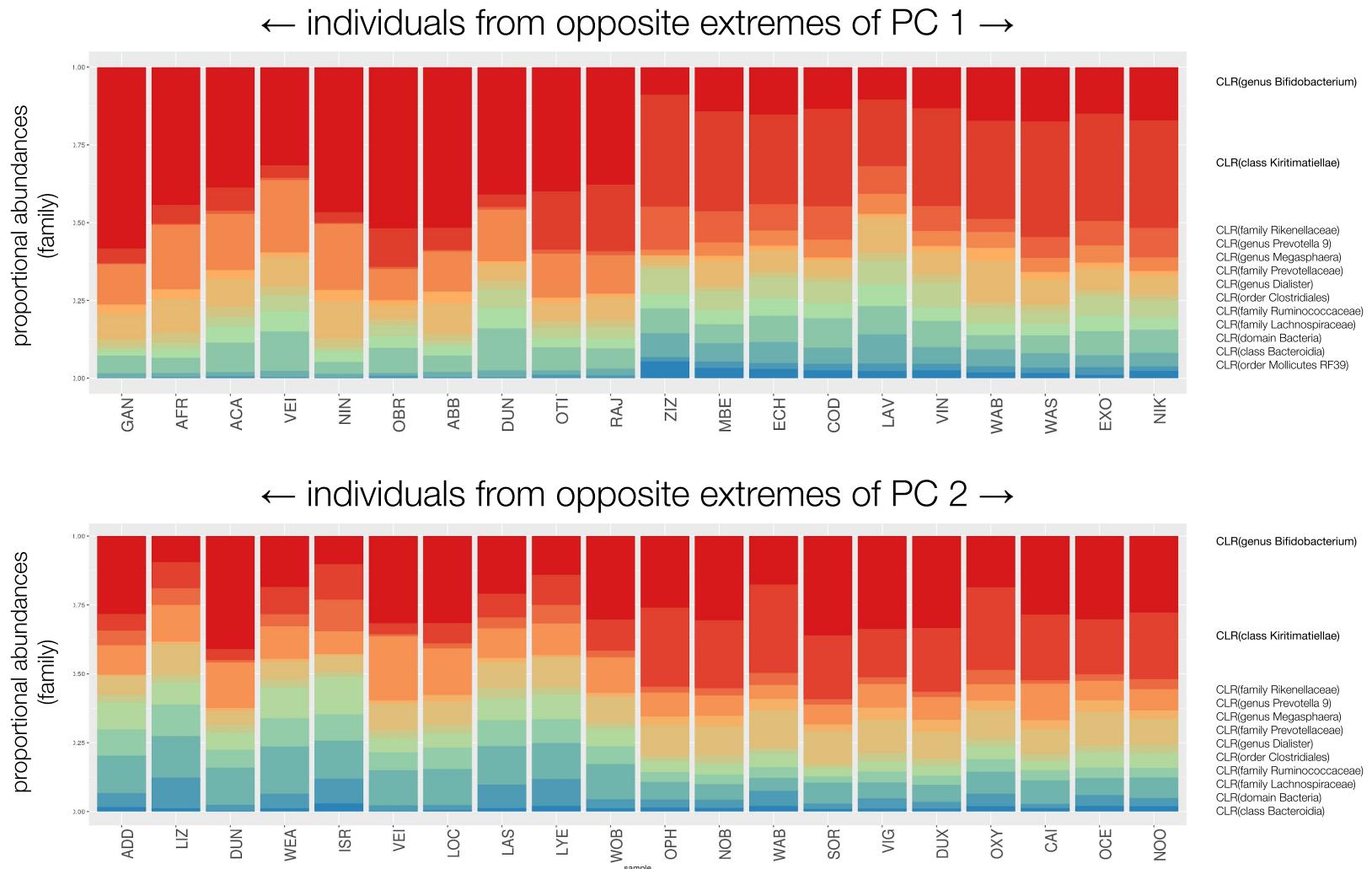
We can coarsely approximate a distance between individuals in terms of *compositional baseline* with this model



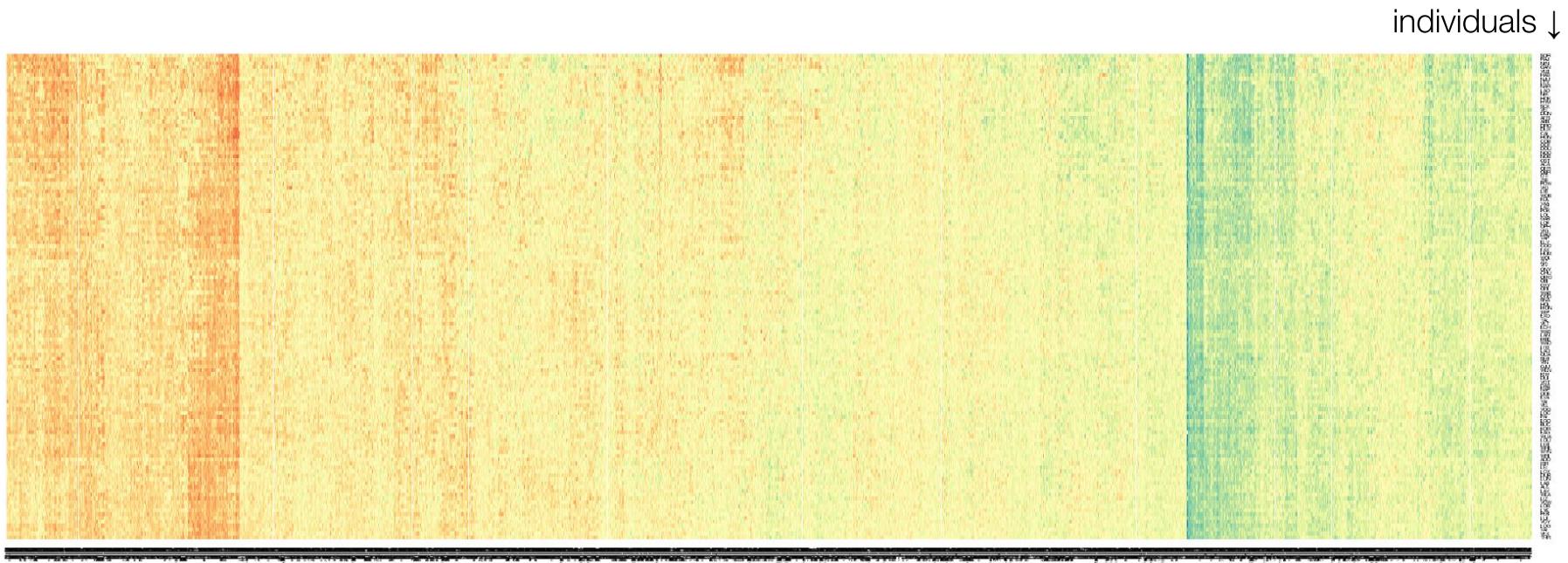
This whole quadrant is (seemingly) rich in Bifido + Kirimatiellae.; details follow

Each point is a posterior sample of  $\Lambda[X]$  from model in slide 40, averaged across time.

# Interpreting the extremes of difference in compositional baseline



Subsets of microbe-microbe interactions appear to be shared between (subsets of) individuals



microbe-microbe interaction (genus level) →

- strong negative covariation
- strong positive covariation

These interactions appear to be conserved between essentially all individuals.

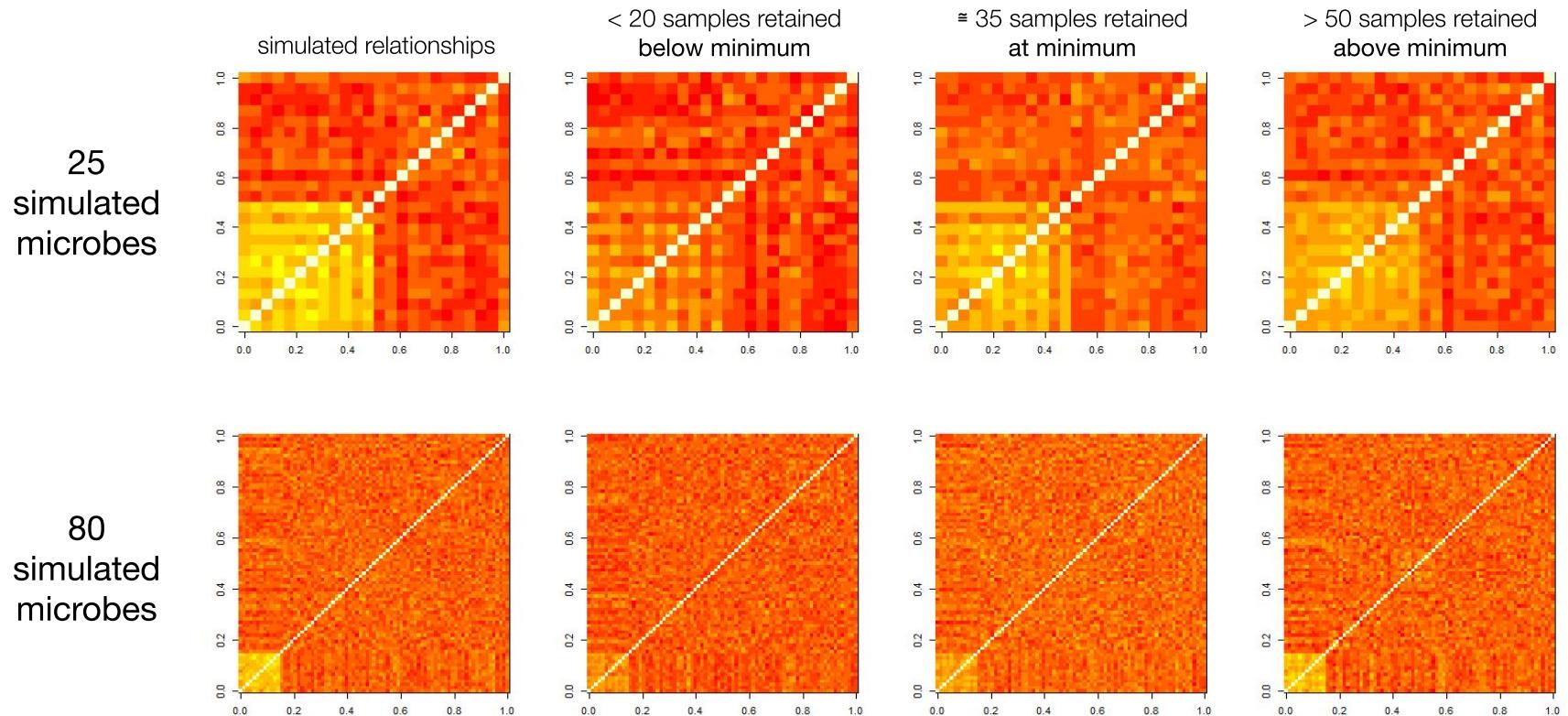
# Functional metagenomics sample selection by partitioning around medoids (PAM)

Functional metagenomics analysis will proceed on 500-2000 selected samples

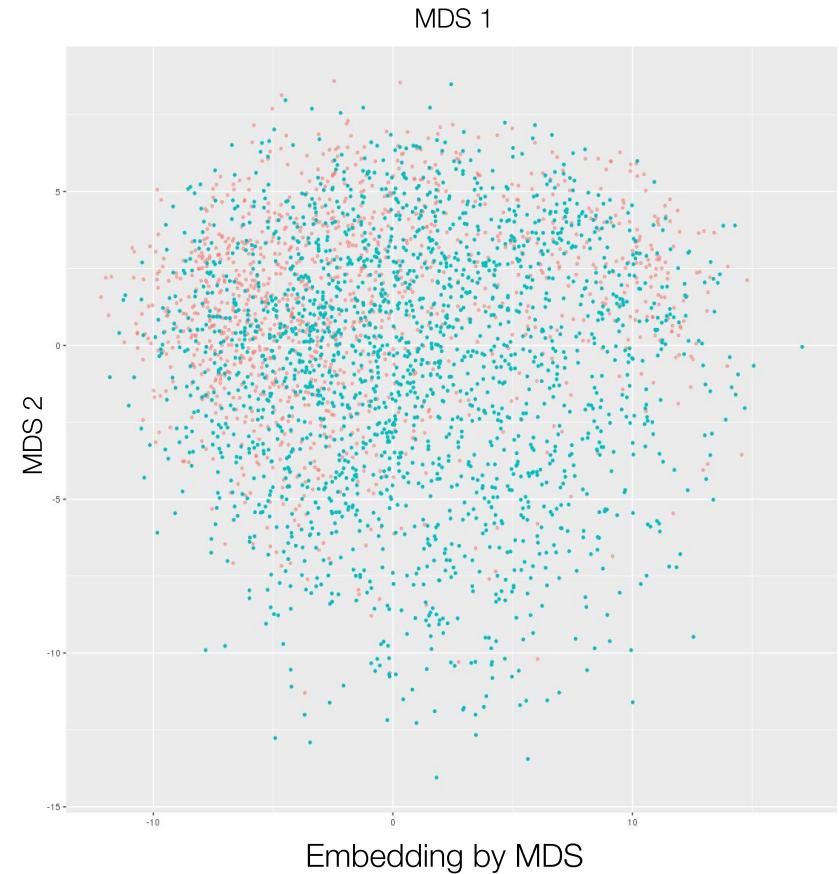
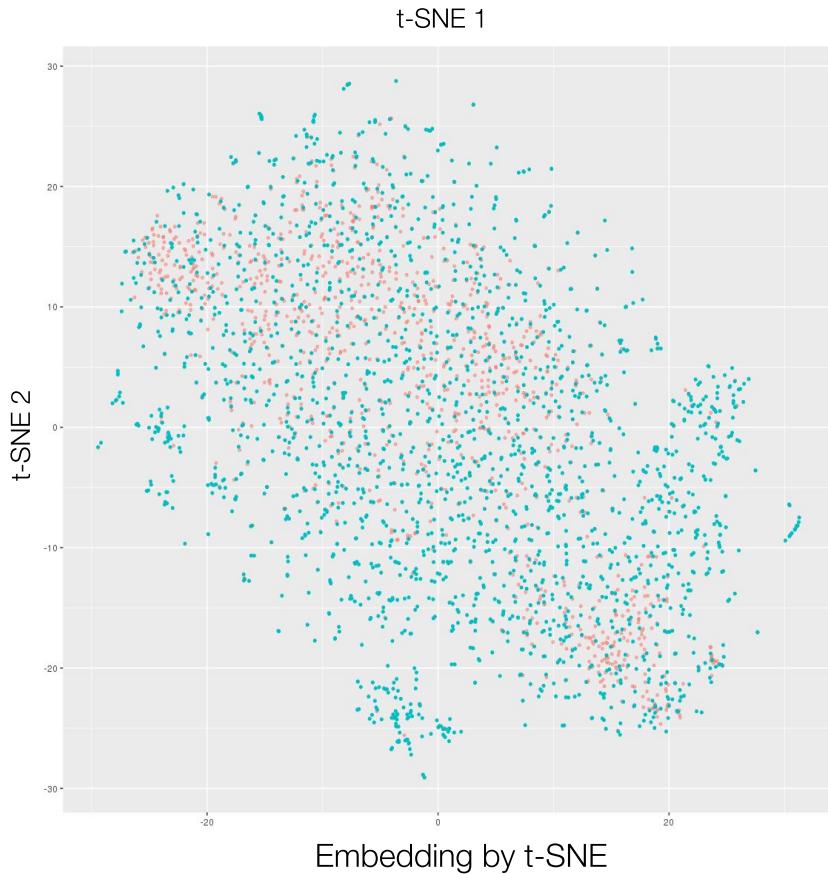
- 1 Collapse to *family* taxonomic level
- 2 Filter on reproductive annotation criteria (yields 86 individuals)
  - a. number (or rate) of live births recorded
  - b. number (or rate) of surviving births recorded
  - c. lifespan recorded
  - d. age at first live birth recorded
- 3 Filter on sampling inclusion criteria (yields 35 individuals); details follow

# Minimum informative sample number threshold can be guided by simulation

- 1 Simulate data with known microbial correlations, downsample data, and fit model.
- 2 A “minimum informative” sample number exists where microbial correlation can be recovered with some probability.



# Functional metagenomics sample selection results: selected samples span the space of variation



Points represent samples from selected individuals. Turquoise points are those selected for functional metagenomics by PAM. These span the space of individual variation in gut microbial composition.

# **A3** Supplement

Distances in terms of microbial dynamics do not appear to associate with differences in selected “fitness” characteristics

	PERMANOVA R <sup>2</sup>	PERMANOVA p-val
social group	0.047	0.000
maternal group	0.084	0.000
rank of mother	0.009	0.038
born in drought	0.008	0.580
born into large group	0.008	0.729
mother died	0.008	0.922
had competing sibling	0.009	0.019
early adversity	0.032	0.566
rate of live births (F)	0.174	0.271
rate of surviving births (F)	0.136	0.123